# Sample Complexity of Block-Sparse System Identification Problem

Salar Fattahi and Somayeh Sojoudi

Abstract-In this paper, we study the system identification problem for sparse linear time-invariant systems. We propose a sparsity promoting block-regularized estimator to identify the dynamics of the system with only a limited number of input-state data samples. We characterize the properties of this estimator under high-dimensional scaling, where the growth rate of the system dimension is comparable to or even faster than that of the number of available sample trajectories. In particular, using contemporary results on high-dimensional statistics, we show that the proposed estimator results in a small element-wise error, provided that the number of sample trajectories is above a threshold. This threshold depends polynomially on the size of each block and the number of nonzero elements at different rows of input and state matrices, but only logarithmically on the system dimension. A by-product of this result is that the number of sample trajectories required for sparse system identification is significantly smaller than the dimension of the system. Furthermore, we show that, unlike the recently celebrated least-squares estimators for system identification problems, the method developed in this work is capable of exact recovery of the underlying sparsity structure of the system with the aforementioned number of data samples. Extensive case studies on switching networks and power systems are offered to demonstrate the effectiveness of the proposed method.

# I. INTRODUCTION

With their ever-growing size and complexity, real-world dynamical systems are hard to model. Today's systems are complex and large, often with a massive number of unknown parameters, which render them doomed to the so-called *curse of dimensionality*. Therefore, system operators should rely on simple and tractable estimation methods to identify the dynamics of the system via a limited number of recorded input-output interactions, and then design control policies to ensure the desired behavior of the entire system. The area of *system identification* is created to address this problem.

In this work, our main goal is to characterize the sample complexity of learning block-sparse linear time-invariant (LTI) systems from noisy input-output trajectories. More specifically, we study the efficient learning of LTI systems in high-dimensional settings, where the system dimension is significantly larger than the number of collected samples. This type of dynamical system forms the basis of many classical control problems, such as Linear Quadratic Regulator and Linear Quadratic Gaussian problems. Our results are built

Email: fattahi@umich.edu and sojoudi@berkeley.edu.

Salar Fattahi is with the Department of Industrial and Operations Engineering, University of Michigan. Somayeh Sojoudi is with the Departments of Electrical Engineering and Computer Sciences and Mechanical Engineering as well as the Tsinghua-Berkeley Shenzhen Institute, University of California, Berkeley. This work was supported by grants from AFOSR, ONR and NSF.

upon the fact that, in many practical large-scale systems, the states and inputs exhibit sparse interactions with one another, which in turn translates into a block-sparse representation of the state-space equations of the system. Driven by the existing non-asymptotic results on the classical Lasso problem, the main focus of this paper is on the block-regularized estimators for the system identification problem, where the goal is to characterize the number of required sample trajectories to reliably estimate the block-sparse interactions of the system. To this goal, the  $\ell_\infty$ -norms of the blocks are penalized instead of their  $\ell_1$ -norms.

In many real-world systems, such as power networks and multi-agent systems, the local state and input behavior of the physical agents/subsystems can be captured and characterized via block matrices in their dynamical models. For instance, in the system identification problem for power systems, each block of the system matrices corresponds to the local states/inputs of an individual generator, and the goal is to learn the sparse interactions among generators given a limited number of measurements from phasor measurement units (PMUs) and supervisory control and data acquisition (SCADA) systems [1], [2]. In this context, it is reasonable to assume that the unknown dynamical interactions among the generators enjoy a block-sparse structure. As another example, consider the problem of planar vertical takeoff and landing (PVTOL) for a fleet of interconnected aircraft. In this context, the number of blocks in the state-space equation of the system corresponds to the number of aircraft that is known a priori, and the goal is to infer the time-varying and uncertain interactions among the aerial vehicles based on the local sensory data [3], [4]. Indeed, such local interactions can be captured via a block-sparse dynamical model.

# A. Related Works

Asymptotic Guarantees: System identification is a well-established area of research in control theory, with related preliminary results dating back to 1960s. Standard reference textbooks on the topic include [5], [6], all focusing on establishing asymptotic consistency of different types of estimators. Although these results shed light on the theoretical consistency of the existing methodologies, they are not applicable in the finite time/sample settings. In many applications, including neuroscience and transportation networks, the dimensionality of the system is overwhelmingly large, often surpassing the number of available input-output data [7], [8]. Under such circumstances, the classical approaches for checking the asymptotic consistency of an estimator face major breakdowns.

**Finite-Time Guarantees:** Contemporary results in statistical learning as applied to system identification seek to characterize *finite time and finite data* rates, relying heavily on tools from sample complexity analysis and concentration of measure. Such finite-time guarantees provide estimates of both system parameters and their uncertainty, which allows for a natural bridge to robust/optimal control. In [9], it was shown that under full state observation, if the system is driven by Gaussian noise, the ordinary least squares estimate of the system matrices constructed from independent data points achieves order optimal rates that are linear in the system dimension. This result was later generalized to the single trajectory setting for (i) marginally stable systems in [10], (ii) unstable systems in [11], and (iii) partially observed stable systems in [12].

Sparse System Identification: Recently, special attention has been devoted to the *sparse* system identification problem, where the states and inputs are assumed to possess localized or low-order interactions. These methods include, but are not restricted to, selective  $\ell_1$ -regularized estimator [13], identification based on compressive sensing [14]-[17], sparse estimation of polynomial system dynamics [18], kernel-based regularization [19], low rank estimation in frequency domain [20], and sparse system identification of time-varying systems [21]. On the other hand, with the unprecedented interest in datadriven control approaches, such as model-free reinforcement learning [22], a question arises as to what the minimum number of input-output data samples should be to guarantee a small error in the estimated model. Answering this question has been the subject of many recent studies on the sample complexity of the system identification problem [9], [23]. Most of these results are tailored to a specific type of dynamics, depend on the stability of the open-loop system, or do not exploit the a priori information about the structure of the system.

Autoregressive processes with sparse graphical models: Another closely-related line of research studies the inference of autoregressive processes whose structures can be captured via sparse graphical models. Earlier works on the inference of sparse autoregressive graphical models were based on hypothesis testing [24], [25]. More recently, the work [26] proposed an  $\ell_1$ -regularized maximum likelihood estimator for estimating the precision matrices of autoregressive Gaussian processes. A similar regularized estimator is also used in [27] to infer autoregressive processes with sparse latent-variable graphical models. Alternatively, the work [28] introduced a Bayesian approach for the inference of autoregressive graphical models. While being related to our proposed method, these works rely upon a different underlying generative model for the system, and hence, are not directly applicable to the system identification of linear time-invariant (LTI) systems.

### B. Contributions:

In this work, we introduce a regularized estimator for recovering the true block-sparsity of an LTI system. In particular, we use an  $\ell_1/\ell_\infty$ -regularized estimator, i.e., a least-squares estimator accompanied by a  $\ell_\infty$  regularizer on different blocks.

We show that the required number of sample trajectories to recover the nonzero blocks of the system matrices and to guarantee a small estimation error scales polynomially with the maximum block sizes and the number of row- and columnwise nonzero elements, but only logarithmically with the number of blocks in the system.

Our work makes a significant improvement over the recently-studied least-squares estimator whose sample complexity scales linearly with the system dimensions. Most interconnected systems consist of many smaller subsystems (blocks) with sparse or localized interactions. Under such circumstances, it may be costly, if not impossible, to collect as many samples as the system dimension. Another advantage of the proposed estimator over its least-squares analog is its exact recovery property. More specifically, we show that while the least-squares estimator is unable to identify the sparsity pattern of the input and state matrices for any finite number of samples, the proposed estimator recovers the true sparsity pattern of these matrices with a sublinear number of sample trajectories. It is worthwhile to mention that this work generalizes the results in [29], where the authors use a similar regularized estimator to learn the dynamics of a particular type of systems. However, [29] ignores the block structure of the system and assumes autonomy and inherent stability, all of which will be relaxed in this work.

This work is a significant extension of our previous conference papers on Lasso-type estimators for system identification [30] and non-asymptotic analysis of block-regularized linear regression problems [31]. In particular, by combining the properties of the block-regularized regression and the characteristics of LTI systems, we provide a unified sparsity-promoting framework for estimating the parameters of the system with arbitrary block structures. To this goal, we have generalized our theoretical results in [30] and [31] to account for partially-sparse structures. We explain the effect of different parameters of the problem—such as input energy and the length of the time horizon—on the sample complexity of the proposed estimator.

**Notations:** For a matrix M, the symbols  $||M||_F$ ,  $||M||_2$ ,  $||M||_0$ ,  $||M||_1$ , and  $||M||_{\infty}$  denote its Frobenius, operator, number of nonzero elements,  $\ell_1/\ell_1$ , and  $\ell_\infty/\ell_\infty$  norms, respectively. Furthermore,  $\kappa(M)$  refers to its 2-norm condition number, i.e., the ratio between its maximum and minimum singular values. Given integer sets I and J, the notation  $M_{IJ}$  refers to the submatrix of M whose rows and columns are indexed by I and J, respectively. The symbols  $M_{:,j}$  and  $M_{i,:}$  refer to the  $j^{th}$  column and  $i^{th}$  row of M, respectively. Given the sequences  $f_1(n)$  and  $f_2(n)$ , the notations  $f_1(n) = O(f_2(n))$ and  $f_1(n) = \Omega(f_2(n))$  imply that there exist  $c_1 < \infty$  and  $c_2 > 0$  such that  $f_1(n) \le c_1 f_2(n)$  and  $f_1(n) \ge c_2 f_2(n)$ , respectively. Furthermore,  $f_1(n) = \Theta(f_2(n))$  is used to imply that  $f_1(n) = O(f_2(n))$  and  $f_1(n) = \Omega(f_2(n))$ . Finally,  $f_1(n) = o(f_2(n))$  is used to show that  $f_1(n)/f_2(n) \to 0$  as  $n \to \infty$ . A zero-mean Gaussian distribution with covariance  $\Sigma$ is shown as  $N(0,\Sigma)$ . Given a function f(x), the expression arg min f(x) refers to its minimizer.

### II. PROBLEM FORMULATION

Consider the LTI system

$$x[t+1] = Ax[t] + Bu[t] + w[t],$$
 (1a)

where t is the time step,  $A \in \mathbb{R}^{n \times n}$  is the state matrix, and  $B \in \mathbb{R}^{n \times m}$  is the input matrix. Furthermore,  $x[t] \in \mathbb{R}^n$ ,  $u[t] \in \mathbb{R}^m$ , and  $w[t] \in \mathbb{R}^n$  are the state, input, and disturbance vectors at time t, respectively. The dimension of the system is defined as m+n. It is assumed that the input disturbance vectors are identically distributed and independent (i.i.d.) with distribution  $N(0, \Sigma_w)$  across different times. In this work, we assume that the matrices A and B are sparse and the goal is to estimate them based on a limited number of sample trajectories, i.e. a sequence  $\{(x^{(i)}[\tau], u^{(i)}[\tau])\}_{\tau=0}^T$  with i=1,2,...,d, where d is the number of available sample trajectories. The ith sample trajectory  $\{(x^{(i)}[\tau], u^{(i)}[\tau])\}_{\tau=0}^T$  is obtained by running the system from t=0 to t=T and collecting the input and state vectors. Note that in general, one may consider two general approaches to obtain the sample input-output trajectories for the system identification problem:

**Fixed** d, and variable T: In this approach, one sets the number of sample trajectories d to a fixed value (e.g., d=1) and instead, chooses a sufficiently long time horizon T to obtain enough information about the dynamics of the system. Notice that *this is only viable when the system is stable*. In other words, one needs to assume that either the system is inherently stable, or there exists an initial stabilizing controller in place to be able to use this approach. Note that this assumption of stability is necessary, as even a simple least-squares estimator may not be consistent if the system has unstable modes [11].

**Fixed** T, and variable d: In this approach, the length of the time horizon T is fixed and instead, the number of sample trajectories is chosen to be sufficiently large to collect enough information about the dynamics of the system. Notice that in this method, one needs to reset the initial state of the system at the beginning of each sample trajectory. However, unlike the previous method, its applicability is not contingent upon the stability of the true system.

Due to the aforementioned theoretical and practical limitations, one can only use the second approach for unstable systems. Such reset-and-run approach is possible and even crucial in many problems of practical relevance. For instance, having the ability to reset cyber-physical systems to a zero or safe state at any given time is deemed crucial to ensure the safety of the system and to protect it from malicious attacks [32], [33]. Moreover, the recent advances in Reinforcement Learning (RL) lends itself to the user's ability to run the system in different and independent sample trajectories (also known as *rollouts* or *episodes* in the RL literature), each with a controlled and independent initial state.

Given the sample trajectories  $\{(x^{(i)}[\tau], u^{(i)}[\tau])\}_{\tau=0}^{\mathsf{T}}$  for i = 1, 2, ..., d, one can obtain an estimate of (A, B) by solving the following least-squares optimization problem:

$$\min_{A,B} \sum_{i=1}^{d} \sum_{t=0}^{T-1} \|x^{(i)}[t+1] - (Ax^{(i)}[t] + Bu^{(i)}[t])\|_{2}^{2}.$$
 (2)

In order to describe the behavior of the least-squares estimator, define

$$Y^{(i)} = \begin{bmatrix} x^{(i)}[1]^{\mathsf{T}} \\ \vdots \\ x^{(i)}[T]^{\mathsf{T}} \end{bmatrix}, \quad X^{(i)} = \begin{bmatrix} x^{(i)}[0]^{\mathsf{T}} & u^{(i)}[0]^{\mathsf{T}} \\ \vdots & \vdots \\ x^{(i)}[T-1]^{\mathsf{T}} & u^{(i)}[T-1]^{\mathsf{T}} \end{bmatrix},$$

$$W^{(i)} = \begin{bmatrix} w^{(i)}[0]^{\mathsf{T}} \\ \vdots \\ w^{(i)}[T-1]^{\mathsf{T}} \end{bmatrix}, \tag{3}$$

for every sample trajectory i=1,2,...,d. Furthermore, let Y, X, and W be defined as vertical concatenations of  $Y^{(i)}$ ,  $X^{(i)}$ , and  $W^{(i)}$  for i=1,2,...,d, respectively. Finally, denote  $\Psi = \begin{bmatrix} A & B \end{bmatrix}^{\mathsf{T}}$  as the unknown system parameter and  $\Psi^*$  as its true value. Based on these definitions, it follows from (1) that

$$Y = X \cdot \Psi + W. \tag{4}$$

The system identification problem is then reduced to estimating  $\Psi$  based on the *observation matrix* Y and the *design matrix* X. Consider the following least-squares estimator:

$$\Psi_{ls} = \arg\min_{\Psi} \|Y - X\Psi\|_F^2. \tag{5}$$

One can easily verify the equivalence of (2) and (5). The optimal solution of (5) can be written as

$$\Psi_{ls} = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}Y = \Psi^* + (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}W. \tag{6}$$

Notice that  $\Psi_{ls}$  is well-defined and unique if and only if  $X^{T}X$  is invertible, which necessitates  $d \ge n+m$ . The estimation error is then defined as

$$E = \Psi_{ls} - \Psi^* = (X^{\mathsf{T}} X)^{-1} X^{\mathsf{T}} W. \tag{7}$$

Thus, one needs to study the behavior of  $(X^{T}X)^{-1}X^{T}W$  in order to control the estimation error of the least-squares estimator. However, since the state of the system at time t is affected by random input disturbances at times 0, 1, ... t-1, the matrices X and W are correlated, which renders (7) hard to analyze. In order to circumvent this issue, [9] simplifies the estimator and considers only the state of the system at time T in  $Y^{(i)}$ . By ignoring the first T-1 rows in  $Y^{(i)}$ ,  $X^{(i)}$ , and  $W^{(i)}$ , one can ensure that the random matrix  $(X^{T}X)^{-1}X^{T}$  is independent of W. Therefore, it is assumed in the sequel that

$$Y = \begin{bmatrix} x^{(1)}[T]^{\mathsf{T}} \\ \vdots \\ x^{(d)}[T]^{\mathsf{T}} \end{bmatrix}, \qquad X = \begin{bmatrix} x^{(1)}[T-1]^{\mathsf{T}} & u^{(1)}[T-1]^{\mathsf{T}} \\ \vdots & \vdots \\ x^{(d)}[T-1]^{\mathsf{T}} & u^{(d)}[T-1]^{\mathsf{T}} \end{bmatrix},$$

$$W = \begin{bmatrix} w^{(1)}[T-1]^{\mathsf{T}} \\ \vdots \\ w^{(d)}[T-1]^{\mathsf{T}} \end{bmatrix}. \tag{8}$$

With this simplification, [9] shows that, with input vectors  $u^{(i)}[t]$  chosen randomly from  $N(0, \Sigma_u)$  for every t=1,2,...,T-1 and i=1,2,...,d, the least-squares estimator requires at least  $d=\Omega(m+n+\log(1/\delta))$  sample trajectories to guarantee  $\|E\|_2 = \mathcal{O}\left(\sqrt{(m+n)\log(1/\delta)/d}\right)$  with probability of at least  $1-\delta$ . In what follows, a block-regularized estimator will be introduced that exploits the underlying sparsity structure of the system dynamics to significantly reduce the number of sample trajectories for an accurate estimation of the parameters. To streamline the presentation, the main technical proofs are deferred to Section A.

### III. MAIN RESULTS

Suppose that A and B can be partitioned as  $A = [A^{(i,j)}]$ and  $B = [B^{(k,l)}]$  where  $(i,j) \in \{1,...,\bar{n}\} \times \{1,...,\bar{n}\}$  and  $(k,l) \in \{1,...,\bar{n}\} \times \{1,...,\bar{m}\}$ .  $A^{(i,j)}$  is the  $(i,j)^{\text{th}}$  block of Awith size  $n_i \times n_j$ . Similarly,  $B^{(k,l)}$  is the (k,l)<sup>th</sup> block of Bwith size  $n_k \times m_l$ . Note that  $\sum_{i=1}^{\bar{n}} n_i = n$  and  $\sum_{i=1}^{\bar{m}} m_i = m$ . Suppose that it is known a priori that all elements in each block  $A^{(i,j)}$  or  $B^{(k,l)}$  are simultaneously zero or nonzero. This implies that, as long as one element in  $A^{(i,j)}$  or  $B^{(k,l)}$ is nonzero, there is no reason to promote sparsity in the remaining elements of the corresponding block. Clearly, this kind of block-sparsity constraint is not correctly reflected in (2). To simplify the presentation, we use the notation  $\Psi = \begin{bmatrix} A & B \end{bmatrix}^{\mathsf{T}}$ . Note that  $\Psi^{(i,j)} = (A^{(j,i)})^{\mathsf{T}}$  for  $i \in \{1,...,\bar{n}\}$ and  $\Psi^{(i,j)} = (B^{(j,i-\bar{n})})^{\top}$  for  $i \in \{\bar{n}+1,...,\bar{n}+\bar{m}\}$ . In order to recover the true block-sparsity of A and B, one can resort to an  $\ell_1/\ell_{\infty}$  variant of the Lasso problem—known as the block-regularized least-squares (or simply block-regularized) problem:

$$\hat{\Psi} = \arg\min_{\Psi} \frac{1}{2d} \|Y - X\Psi\|_F^2 + \lambda_d \|\Psi\|_{\text{block}}, \tag{9}$$

where  $\|\Psi\|_{\text{block}}$  is defined as the summation of  $\|\Psi^{(i,j)}\|_{\infty}$ over  $(i,j) \in \{1,...,\bar{n} + \bar{m}\} \times \{1,...,\bar{n}\}$ . D is used to denote the maximum size of the blocks of  $\Psi$ . Under the sparsity assumption on (A, B), we will show that the non-asymptotic statistical properties of  $\Psi$  significantly outperform those of  $\Psi_{ls}$ . In particular, the primary objective is to prove that  $\|\hat{\Psi} - \Psi^*\|_{\infty}$ decreases at the rate  $\mathcal{O}(\sqrt{D\log(n+m)} + D^2\log(1/\delta)/d)$ with probability of at least  $1 - \delta$  and with an appropriate scaling of the regularization coefficient, provided that d = $\Omega\left(k_{\max}^2\left(D\log(\bar{n}+\bar{m})+D^2\log(1/\delta)\right)\right)$ . Here,  $k_{\max}$  is the maximum number of nonzero elements in the columns of  $[A \ B]^{\mathsf{T}}$ . Comparing this number with the required lower bound  $\Omega(n+m+\log(1/\delta))$  on the number of sample trajectories for the least-squares estimator, we conclude that the proposed method needs significantly fewer samples when A and B are sparse. The third objective is to prove that this method is able to find the correct block-sparsity structure of A and B with high probability. In contrast, it will be shown that the solution of the least-squares estimator is fully dense for any finite number of sample trajectories, and hence, it cannot correctly extract the sparsity structures of A and B. We will showcase the superior performance of the block-regularized estimator both in sparsity identification and estimation accuracy in simulations.

To present the main results of this work, first note that

$$x^{(i)}[T-1] = A^{T-2}Bu^{(i)}[0] + A^{T-3}Bu^{(i)}[1] + \dots + Bu^{(i)}[T-2] + A^{T-2}w^{(i)}[0] + A^{T-3}w^{(i)}[1] + \dots + w^{(i)}[T-2] + A^{T-1}x[0].$$
(10)

Suppose that  $u^{(i)}[t]$  and  $w^{(i)}[t]$  are i.i.d samples of  $N(0, \Sigma_u)$  and  $N(0, \Sigma_w)$ , respectively. Moreover, we assume that the initial state is random with a Gaussian distribution  $N(0, \Sigma_x)$ . Therefore, (10) and (8) imply that

$$X_{i,:}^{\mathsf{T}} \sim N\left(0, \tilde{\Sigma}\right),$$
 (11)

where  $X_{i,:}$  is the  $i^{th}$  row of X and

$$\tilde{\Sigma} = \begin{bmatrix} C^{\mathsf{T}}C & 0\\ 0 & \Sigma_u \end{bmatrix}, C = \begin{bmatrix} F_T^{\mathsf{T}}\\ G_T^{\mathsf{T}} \end{bmatrix}$$
 (12a)

$$F_T = \begin{bmatrix} A^{T-2}B\Sigma_u^{1/2} & A^{T-3}B\Sigma_u^{1/2} & \dots & B\Sigma_u^{1/2} \end{bmatrix}$$
 (12b)

$$G_T = \begin{bmatrix} A^{T-1} \Sigma_x^{1/2} & A^{T-2} \Sigma_w^{1/2} & A^{T-3} \Sigma_w^{1/2} & \dots & \Sigma_w^{1/2} \end{bmatrix}.$$
(12c)

The matrix C is referred to as the *combined controllability* matrix in the sequel. Define  $\mathcal{A}_j(\Psi) = \{i : \Psi^{(i,j)} \neq 0\}$ . Unless stated otherwise,  $\mathcal{A}_j$  is used to refer to  $\mathcal{A}_j(\Psi^*)$ . Define  $\mathcal{A}_j^c$  as the complement of  $\mathcal{A}_j$ . For  $\mathcal{T} \subseteq \{1,...,\bar{n}+\bar{m}\}$ , denote  $I(\mathcal{T})$  as the index set of rows in  $\Psi^*$  corresponding to the blocks  $\{\Psi^{*(i,i)} : i \in \mathcal{T}\}$ . For an index set  $\mathcal{U}$ , define  $X_{\mathcal{U}}$  as a  $d \times |\mathcal{U}|$  submatrix of X after removing the columns with indices not belonging to  $\mathcal{U}$ . With a slight abuse of notation,  $X_{(i)}$ ,  $X_{\mathcal{A}_j}$ , and  $X_{\mathcal{A}_j^c}$  are used to denote  $X_{I(\{i\})}$ ,  $X_{I(\mathcal{A}_j)}$ , and  $X_{I(\mathcal{A}_j^c)}$  when there is no ambiguity. Similarly,  $\tilde{\Sigma}_{(i),\mathcal{A}_j}$  and  $\tilde{\Sigma}_{\mathcal{A}_j,\mathcal{A}_j}$  are used in lieu of  $\tilde{\Sigma}_{I(\{i\}),I(\mathcal{A}_j)}$  and  $\tilde{\Sigma}_{I(\mathcal{A}_j),I(\mathcal{A}_j)}$ , respectively. Denote  $k_j$  as the maximum number of nonzero elements in any column of  $\Psi^{*(:,j)}$  which is the  $j^{\text{th}}$  block column of  $\Psi^*$ . Finally, define

$$n_{\max} = \max_{1 \le i \le \bar{n}} n_i, \qquad m_{\max} = \max_{1 \le i \le \bar{m}} m_i,$$

$$p_{\max} = \max \left\{ n_{\max}, m_{\max} \right\}, \qquad k_{\max} = \max_{1 \le j \le \bar{n}} k_j,$$

$$\sigma_{\max}^2 = \max_{1 \le i \le n+m} \tilde{\Sigma}_{ii}. \tag{13}$$

The following set of assumptions plays a key role in deriving the main result of this paper:

**Assumption 1.** By fixing the time horizon T, we assume that the following conditions hold for all finite system dimensions:

A1. (Mutual Incoherency Property): There exists a number  $\gamma \in (0,1]$  such that

$$\max_{j=1,\dots,\bar{n}} \left\{ \max_{i \in \mathcal{A}_j^c} \left\| \tilde{\Sigma}_{(i),\mathcal{A}_j} (\tilde{\Sigma}_{\mathcal{A}_j,\mathcal{A}_j})^{-1} \right\|_1 \right\} \le 1 - \gamma.$$
 (14)

A2. (Bounded eigenvalue): There exist numbers  $0<\Lambda_{\min}<\infty$  and  $0<\Lambda_{\max}<\infty$  such that

$$\Lambda_{\min} \le \lambda_{\min}(\tilde{\Sigma}) \le \lambda_{\max}(\tilde{\Sigma}) \le \Lambda_{\max}.$$
(15)

A3. (Bounded minimum value): There exists a number  $t_{\rm min}$  > 0 such that

$$t_{\min} \le \min_{1 \le j \le \bar{n}} \min_{i \in \mathcal{A}_j} \left\| \Psi^{*(i,j)} \right\|_{\infty}. \tag{16}$$

A4. (Block sizes): There exist numbers  $\alpha_n, \alpha_m < \infty$  such that

$$n_{\max} = O\left(\left(\bar{n} + \bar{m}\right)^{\alpha_n}\right), \ m_{\max} = O\left(\left(\bar{n} + \bar{m}\right)^{\alpha_m}\right). \tag{17a}$$

The mutual incoherency property in Assumption A1 is a commonly known condition for the exact recovery of unknown parameters in compressive sensing and classical Lasso problems [34], [35]. This assumption entails that the effect of those submatrices of  $\tilde{\Sigma}$  corresponding to zero (unimportant) elements of  $\Psi$  on the remaining entries of  $\tilde{\Sigma}$  should not be large. Roughly speaking, this condition guarantees that the

unknown parameters are *recoverable* in the noiseless scenario, i.e. when W=0. It is also worth noting that this condition can be further relaxed under additional conditions [36]. If the recovery cannot be guaranteed in the noise-free setting, then there is little hope for the block-regularized estimator to recover the true structure of A and B when the system is subject to noise.

The bounded eigenvalue condition in Assumption A2 entails that the condition number of  $\tilde{\Sigma}$  is bounded away from 0 and  $\infty$  for all finite system dimensions. Assuming that the eigenvalues of  $\Sigma_u$  and  $\Sigma_w$  do not scale with the system dimension, it is easy to verify that  $\min\{\lambda_{\min}(\Sigma_u),\lambda_{\min}(\Sigma_w)\} \leq \Lambda_{\min} \leq \lambda_{\min}(\Sigma_w)$ . However, as will be shown later, the value of  $\Lambda_{\max}$  can change with respect to the time horizon T. In particular, it will be later shown that for highly unstable systems,  $\tilde{\Sigma}$  becomes severely ill-conditioned as the time horizon increases, which in turn makes the system identification problem difficult to solve. Furthermore, this assumption implies that there exists a constant  $\bar{\sigma}_{\max}^2 < \infty$  such that  $\max_{1 \leq i \leq n+m} \tilde{\Sigma}_{ii} \leq \bar{\sigma}_{\max}^2$ .

Assumption A3 implies that, independent of the system dimensions, there always exists a strictly positive gap between the zero and nonzero elements of A and B. This assumption holds in almost all practical settings and will facilitate the exact sparsity recovery of the parameters of the system.

Finally, Assumption A4 requires that the maximum size of the blocks in  $\Psi^*$  be polynomially bounded by the number of its block columns. For instance,  $\bar{n} = O(1)$  and  $\bar{m} = O(1)$  violate this assumption since it implies that  $n_{\max} = \Omega((\bar{n} + \bar{m})^{\log n})$  and  $m_{\max} = \Omega((\bar{n} + \bar{m})^{\log m})$ . It is worthwhile to mention that Assumption A4 results in  $k_{\max} = O((\bar{n} + \bar{m})^{\alpha_k})$  for some number  $\alpha_k < \infty$ ; this will be used later in the derivations.

Define  $D = p_{\text{max}} n_{\text{max}}$ , which is the maximum size of the blocks in  $\Psi$ .

Theorem 1 (block-wise regularization). Upon choosing

$$\lambda_{d} = \Theta\left(\sigma_{\max}\sqrt{\frac{D\log(\bar{n} + \bar{m}) + D^{2}\log(1/\delta)}{d}}\right), \tag{18a}$$

$$d = \Omega\left(\kappa(\tilde{\Sigma})^{2}k_{\max}\left(D\log(\bar{n} + \bar{m}) + D^{2}\log(1/\delta)\right)\right), \tag{18b}$$

the following statements hold with probability of at least  $1-\delta$ :

- 1.  $\hat{\Psi}$  is unique and has the same nonzero blocks as  $\Psi^*$ .
- 2. We have

$$g = \|\hat{\Psi} - \Psi^*\|_{\infty}$$

$$= O\left(\kappa(\tilde{\Sigma}) \left(1 + \sqrt{\frac{k_{\max}(k_{\max}n_{\max} + \log(\bar{n} + \bar{m}) + \log(1/\delta))}{d}}\right) \times \sqrt{\frac{D\log(\bar{n} + \bar{m}) + D^2\log(1/\delta)}{d}}\right). \tag{19}$$

Theorem 1 shows that the minimum number of required sample trajectories is a quadratic function of the maximum block size. Therefore, only a small number of samples is enough to guarantee the uniqueness, exact block-sparsity recovery, and small estimation error for sparse systems, assuming that the sizes of the blocks are significantly smaller than the system dimensions.

**Corollary 1.** Assume that  $n_{\text{max}} = O(n^{\beta_n})$  and  $m_{\text{max}} = O(m^{\beta_m})$  for some  $\beta_n > 0$  and  $\beta_m > 0$ . Then,

$$\lambda_d = \Theta\left(\sigma_{\max}(n+m)^{(\beta_n+\beta_m)}\sqrt{\frac{\log(1/\delta)}{d}}\right),\tag{20a}$$

$$d = \Omega(\kappa(\tilde{\Sigma})^2 k_{\max}^2 (n+m)^{2(\beta_n+\beta_m)} \log(1/\delta)), \tag{20b}$$

is enough to guarantee the exact sparsity recovery of  $\Psi^*$  and

$$\|\hat{\Psi} - \Psi^*\|_{\infty} = O\left(\kappa(\tilde{\Sigma})(n+m)^{(\beta_n + \beta_m)} \sqrt{\frac{\log(1/\delta)}{d}}\right), \quad (21)$$

with probability of at least  $1 - \delta$ .

*Proof.* The proof follows from Theorem 1. The details are omitted for brevity.  $\Box$ 

Corollary 1 analyzes the behavior of the proposed estimator for the *polynomial scaling* of the block size. It can be seen that the size of the required sample trajectories heavily depends on the growth rate of the maximum block size of  $\Psi$ . Although the sampling rate is still sublinear when  $\beta_n + \beta_m < 1/2$ , it may surpass the system dimension if  $\beta_n + \beta_m > 1/2$ . A question arises as to whether one can resort to the ordinary least-squares estimator in lieu of the proposed block-regularized estimator for the cases where  $\beta_n + \beta_m > 1/2$  since the proposed estimator requires  $d = \Omega((n+m)^{1+\epsilon}\log(1/\delta))$  for some  $\epsilon > 0$  whereas  $d = \Theta(n+m+\log(1/\delta))$  is enough to guarantee the uniqueness of the least-squares estimator. This will be addressed in the next subsection.

**Remark 1.** In this paper, we assume that A and B are partitioned into blocks with known sizes, each with a maximum size of D. If the blocks sizes are unknown, an alternative approach is to treat A and B as sparse matrices where each block is of size D = 1. This lack of prior knowledge on the block sizes of the system matrices can be compensated with a higher number of collected sample trajectories from the system. In particular, as it is shown in [30], an element-wise regularized estimator (i.e. vanilla Lasso) can still recover the correct sparsity pattern of the true system matrices with no prior knowledge on the block sizes, albeit with potentially a higher number of sample trajectories and worse estimation error. In Section IV, we showcase the performance of these regularized estimators with and without prior knowledge on the block sizes.

It is worth noting that, based on Theorem 1, one may speculate that setting D=1 (i.e., not using the prior information on the block sizes) may lead to a better statistical guarantee. However, note that the derived bound is based on a customized  $\lambda_d$  that is designed to obtain a logarithmic dependency on  $\bar{n}+\bar{m}$ . This  $\lambda_d$  is specifically designed to offer a small value in terms of  $\bar{n}+\bar{m}$  without optimizing its dependency on D. To obtain a tighter bound with respect to D (instead of  $\bar{n}+\bar{m}$ ), one may need to select another  $\lambda_d$  that (1) would depend on D in a more sophisticated way; and (2) similar to [37], would potentially depend on the level of "overlap" in the block-wise support of the unknown parameters. We consider obtaining a better dependency on D as an enticing challenge for future research.

Remark 2. Similar to the classical results on the regularized linear regression [38], [39], the particular choice of the regularization coefficient  $\lambda_d$  in our analysis depends on the unknown parameters of the true system, such as  $\sigma_w$ ,  $\sigma_{\max}$ , and  $\gamma$ . As will be shown in the next section, in practice we do not rely on these unknown parameters. In particular, the chosen value for  $\lambda_d$  in our simulations will merely depend on the known parameters of the system, such as d,  $\bar{n} + \bar{m}$ , and D when we know the block sizes, or d and n+m when the block sizes are unknown.

**Remark 3.** Another alternative approach to promote the block sparsity in the identification of dynamical systems is  $\ell_1/\ell_2$ regularized estimator (also known as group Lasso), where the  $\ell_{\infty}$  regularization on different blocks is replaced by a  $\ell_2$ regularization [40], [41]. In Section IV, it is empirically shown that these estimators offer a similar performance in terms of the estimation error. However, an important advantage of the  $\ell_1/\ell_{\infty}$ -regularized estimator over the group Lasso is in terms of its computational complexity. As pointed out in [42], one of the main benefits of  $\ell_1/\ell_{\infty}$ -regularized estimator lies in the efficient computation of its entire solution path over a compact range of regularization coefficients (as opposed to a single regularization coefficient). In particular, contrary to the group Lasso, the solution path for  $\ell_1/\ell_{\infty}$ -regularized estimator is piecewise linear with easily computable breakpoints. This in turn can be used in sensitivity analysis and boosting methods [42], [43].

### A. Comparison to Least-Squares

In this subsection, we prove that the least-squares estimator does not extract the correct sparsity structure of  $\Psi$  for any finite number of sample trajectories.

**Theorem 2.** If A and B are not fully dense matrices,  $\Psi_{ls}$  does not recover the support of  $\Psi^*$  for any finite number of sample trajectories with probability 1.

*Proof.* The proof is omitted for brevity and can be found in [44].

Define  $h(n,m) = \sqrt{(n+m)\log(1/\delta)/d}$  and recall that  $\|\Psi_{ls} - \Psi^*\|_2 = O(h(n,m))$ . In the next corollary, we show that, under additional sparsity conditions, the operator norm of the estimation error for  $\hat{\Psi}$  becomes arbitrarily smaller than h(n,m) as the system dimension grows.

**Corollary 2.** Assume that the number of nonzero elements at different rows and columns of  $\Psi^*$  is upper bounded by  $k_{\max}$ . Furthermore, suppose that  $\lambda_d$  satisfies (18a) and

$$d = \Omega\left(\kappa(\tilde{\Sigma})^2 k_{\text{max}}^2 \left(D\log(\bar{n} + \bar{m}) + D^2\log(1/\delta)\right)\right). \tag{22}$$

Then, we have

$$\|\hat{\Psi} - \Psi^*\|_2 = O\left(\underbrace{\kappa(\tilde{\Sigma})k_{\max}\sqrt{\frac{D\log(\bar{n} + \bar{m}) + D^2\log(1/\delta)}{d}}}_{v(n,m)}\right),$$

with probability of at least  $1 - \delta$ . Furthermore, we have

$$\frac{v(n,m)}{h(n,m)} \to 0 \quad as \quad (n,m) \to \infty, \tag{24}$$

provided that

$$k_{\max}D = o\left(\sqrt{\frac{n+m}{\log(n+m)}}\right). \tag{25}$$

*Proof.* The proof is omitted for brevity and can be found in [44].

Corollary 2 describes the settings under which our proposed method significantly outperforms the least-squares estimator in terms of the operator norm of the errors. This improvement is more evident for those systems where the states and inputs have sparse interactions and the block sizes in A and B are smaller than the system dimensions. A class of such systems is multi-agent networks where the agents interact only locally and their total number dominates the dimension of each individual agent.

# B. Controllability and the Effect of T

Notice that the minimum number of required sample trajectories and the element-wise error of the estimated parameters depend on  $\kappa(\tilde{\Sigma})$ . Recall that  $\min\{\lambda_{\min}(\Sigma_u), \lambda_{\min}(\Sigma_w)\} \leq$  $\Lambda_{\min} \leq \lambda_{\min}(\Sigma_w)$ , independent of T. Therefore, the value of  $\kappa(\hat{\Sigma})$  is governed by the maximum eigenvalue of  $C^{\mathsf{T}}C$ . Roughly speaking,  $\lambda_{\max}(C^{\mathsf{T}}C)$  quantifies the easiest-toidentify mode of the dynamical system. Therefore, Theorem 1 imply that the sample complexity of the proposed blockregularized estimator depends on the modes of the system, as well as the expected energy of the input and disturbance *noise*. In particular, by fixing  $\Sigma_u$  and  $\Sigma_w$ , only a small number of samples is required to accurately identify the dynamics of the system if all of its modes are easily excitable. The dependency of the estimation error on the modes of the system is also reflected in the non-asymptotic error bound of the least-squares estimator in [9]. This is completely in line with the conventional results on the identifiability of dynamical systems: independent of the method in use, it is significantly harder to identify the parameters of the system accurately if it possesses nearly-hidden modes.

On the other hand, for fixed  $\sigma_w$ , the performance of the estimator deteriorates as the expected energy of the input decreases. In the extreme case of zero input, we inevitably have  $\Lambda_{\min} = 0$ , which in turn implies that the proposed estimator provides no guarantee on the accuracy of the estimated parameters.

Furthermore, notice that  $F_T$ ,  $G_T$ , and, hence,  $\lambda_{\max}(C^{\mathsf{T}}C)$  depend directly on the length of the time horizon T for each sample trajectory. In what follows, we will show that for highly unstable systems,  $\lambda_{\max}(C^{\mathsf{T}}C)$  can grow exponentially fast in terms of T and, hence, short sample trajectories are more desirable in estimating the parameters of such unstable systems. To better understand this, assume that the spectral radius of A—shown as  $\rho(A)$ —is greater than one, it is diagonalizable, and n is fixed. One can easily verify that the

following chain of inequalities holds:

$$\lambda_{\max}(\tilde{\Sigma}) \geq \lambda_{\max}(\sigma_{u}^{2} F_{T} F_{T}^{\mathsf{T}} + \sigma_{w}^{2} G_{T} G_{T}^{\mathsf{T}})$$

$$\geq \lambda_{\min}(\Sigma_{w}) \lambda_{\max} \left(A^{T-2} (A^{T-2})^{\mathsf{T}}\right)$$

$$\geq \lambda_{\min}(\Sigma_{w}) \max_{i} \left\{ \left(\left(A^{T-2} \left(A^{T-2}\right)^{\mathsf{T}}\right)_{ii}\right)^{2} \right\}$$

$$\geq \frac{\lambda_{\min}(\Sigma_{w})}{n} \|A^{T-2}\|_{\infty} \geq \frac{\lambda_{\min}(\Sigma_{w})}{n} \rho(A)^{T-2}. \tag{26}$$

This exponential dependency is also empirically observed in our numerical experiments.

# C. Mutual Incoherency

In this subsection, we will analyze the mutual incoherency condition (14). In particular, we will show that the proposed mutual incoherency condition is tightly related to the so-called *identifiability* condition, and hence, cannot be relaxed for specific classes of problems. For simplicity of the subsequent arguments, assume that the size of each block is equal to 1, and that the oracle estimator can measure the disturbance matrix W. Furthermore, suppose that the estimator can collect and work with an infinite number of sample trajectories. Under these assumptions, the oracle estimator should solve the following optimization problem to estimate the parameters of the system:

$$\min_{\Psi} \|\Psi\|_0 \tag{27a}$$

s.t. 
$$X\Psi = Y - W$$
. (27b)

Notice that the oracle estimator cannot be obtained in practice since: 1) the exact value of the disturbance noise is not available, 2) only a finite number of sample trajectories can be collected, and 3) the corresponding optimization is non-convex and NP-hard in its worst case.

As mentioned before, there are fundamental limits on the performance of the introduced oracle estimator. To explain this, we introduce the mutual-coherence metric for a matrix. For a given matrix  $A \in \mathbb{R}^{t_1 \times t_2}$ , its mutual-coherence  $\mu(A)$  is defined as

$$\mu(A) = \max_{1 \le i < j \le t_2} \frac{|A_{:,i}^{\top} A_{:,j}|}{\|A_{:,i}\|_2 \|A_{:,j}\|_2}.$$
 (28)

In other words,  $\mu(A)$  measures the maximum correlation between distinct columns of A (with a slight abuse of notation, we assume that  $\frac{1}{\mu(A)} = +\infty$  if  $\mu(A) = 0$ ). Reminiscent of the classical results in the compressive sensing literature, it is well-known that the optimal solution  $\Psi^*$  of (27) is unique if the *identifiability* condition

$$\|\Psi_{:,j}^*\|_0 < \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right) \tag{29}$$

holds for every j=1,2,...,n (see, e.g., Theorem 2.5 in [45]). Furthermore, this bound cannot be tightened, since there exist instances of the problem for which the violation of  $\|\Psi_{:,j}^*\|_0 < \frac{1}{2}\left(1+\frac{1}{\mu(X)}\right)$  for some j results in the non-uniqueness of the optimal solution. On the other hand, one can invoke the Central Limit Theorem to show that  $\frac{1}{d}X^{\mathsf{T}}X=\tilde{\Sigma}$  almost surely as  $d\to\infty$ . Furthermore, recall the definition of the

combined controllability matrix C in (12a). This, together with the definition of  $\tilde{\Sigma}$ , implies that

$$\mu(X) = \max_{1 \le i < j \le m+n} \frac{|X_{:,i}^{\top} X_{:,j}|}{\|X_{:,i}\|_2 \|X_{:,j}\|_2}$$

$$= \max_{1 \le i < j \le n} \frac{|C_{:,i}^{\top} C_{:,j}|}{\|C_{:,i}\|_2 \|C_{:,j}\|_2} = \mu(C). \tag{30}$$

According to the above equality, the correlation between different columns of C plays a crucial role in the identifiability of the true parameters: as  $\mu(C)$  becomes smaller, the oracle estimator can correctly identify the structure of  $\Psi$  for a wider range of sparsity levels.

Revisiting Assumption A1, one can verify that the mutual incoherency condition is reduced to the following inequality when the size of each block is equal to one:

$$\left\| (C_{:,\mathcal{A}_{j}}^{\mathsf{T}} C_{:,\mathcal{A}_{j}})^{-1} C_{:,\mathcal{A}_{j}}^{\mathsf{T}} C_{:,k} \right\|_{1} \le 1 - \alpha,$$

$$\forall k \in \mathcal{A}_{j}^{c}, \ j = 1, 2, \dots, n \quad (31)$$

where, with a slight abuse of notation, we use  $A_j$  to denote the set  $\{i: A_{ij} \neq 0\}$ . Notice that, similar to (29), the above condition is expected to be satisfied when different columns of C are nearly orthogonal, i.e., when the elements in  $C_{:,A_j}^{\mathsf{T}}C_{:,k}$  have small magnitudes. In particular, we introduce a class of k-sparse dynamical systems for which the above condition is equivalent to (29) (modulo a constant factor).

k-sparse systems: Consider a class of problems where each row or column of A has at most k nonzero entries and B is diagonal. Without loss of generality and to simplify the subsequent derivations, suppose that the following assumptions hold:

- B is equal to identity matrix and diagonal entries of A are equal to 1. Moreover, the magnitude of each off-diagonal entry of A is upper bounded by  $\varphi > 0$ .
- T is set to 3.
- $\Sigma_u = \sigma_u^2 I$  and  $\Sigma_w = \sigma_w^2 I$ , where  $\sigma_u$  and  $\sigma_w$  are less than or equal to 1. Moreover  $\Sigma_x = 0$ .

**Proposition 1.** For k-sparse systems with  $k \ge 3$ , the following statements hold:

- There exists an instance for which the identifiability condition fails to hold for the oracle estimator if  $\varphi \ge \frac{3}{k}$ .
- The mutual incoherency condition holds if  $\varphi < \frac{\sigma_u + \sigma_w}{9k}$ .

*Proof.* The proof is omitted for brevity and can be found in [44].

The tightness of the identifiability condition 29 together with the Proposition 1 implies that for some specific classes of problems, it is not possible to a have a consistent sparsity promoting technique with significantly more relaxed conditions than the ones introduced in this paper. However, we point out that, in general, the mutual incoherency condition may be improved by resorting to more sophisticated (and potentially nonconvex) estimators [36]. Moreover, it will be shown in Section IV that the incoherency condition is expected to hold in many cases of practical relevance.

### IV. NUMERICAL RESULTS

In this section, we illustrate the performance of the block-regularized estimator and compare it with its least-squares counterpart. We consider two case studies on switching networks and power systems.

Define the (block) mismatch error as the total number of false positives and false negatives in the (block) sparsity pattern of the estimator. Moreover, define *relative number of sample trajectories* (RST) as the number of sample trajectories normalized by the dimension of the system, and *relative* (block) *mismatch error* (RME) as the mismatch error normalized by total number of elements (blocks) in  $\Psi$ .

# A. Case Study 1: Switching Networks

In this case study, we study a network of multi-agent systems that are interconnected through a switching information exchange topology. Recently, a special attention has been devoted to multi-agent systems with a time-varying network topology; in many communication networks, each sensor has access only to the information of its neighbors. Therefore, when the location of these sensors changes over time, so does the topology of the interconnecting links [46]. The *dwell time* is defined as the time interval in which the network topology is unchanged. The goal is to identify the structure of the network within the dwell time. The state-space equation of agent *i* admits the following general form:

$$\dot{x}_{i}(t) = \sum_{(i,j)\in\mathcal{N}_{x}(i)} A^{(i,j)} x_{j}(t) + \sum_{(i,j)\in\mathcal{N}_{u}(i)} B^{(i,j)} u_{j}(t) + w_{i}(t),$$
(32)

where, as before,  $A^{(i,j)} \in \mathbb{R}^{n_i \times n_i}$  and  $B^{(i,j)} \in \mathbb{R}^{n_i \times m_i}$  are the  $(i,j)^{\text{th}}$  blocks of A and B. Furthermore,  $\mathcal{N}_x(i)$  and  $\mathcal{N}_u(i)$  are the sets of neighbors of agent i whose respective state and input actions affect the state of agent i.

We consider 200 agents connected through a randomly generated sparse network. In particular, we assume that each agent is connected to 5 other agents. If  $j \in \mathcal{N}_x(i)$  or  $j \in \mathcal{N}_u(i)$ , then each element of  $A^{(i,j)}$  or  $B^{(i,j)}$  is randomly selected from  $[-0.4 - 0.3] \cup [0.3 \ 0.4]$ . Moreover, the regularization coefficient  $\lambda_d$  is set to

$$\sqrt{\frac{2(D^2 + D\log(\bar{n} + \bar{m}))}{d}}.$$
 (33)

Note that this choice of  $\lambda_d$  does not rely on the unknown parameters of the system, and it does not require any additional fine-tuning. The behavior of the proposed block-regularized estimator will be examined for different dimensions of the agents. In particular, we investigate the performance of this estimator in comparison with the Lasso for which the sparsity of the system matrices is promoted on different elements independent of the block structures. In these experiments,  $(n_i, m_i)$  is chosen from  $\{(5,5),(8,8),(11,11)\}$ . This entails that  $D \in \{25,64,121\}$  and  $(n,m) \in \{(1000,1000),(1600,1600),(2200,2200)\}$ . Furthermore, T is set to 3 and the system is discretized using the forward Euler method with the sampling time of 0.2 seconds. This implies that each sample trajectory is collected

within 0.6 seconds. The number of block mismatch and 2-norm estimation errors are depicted in Figures 1a and 1b with respect to the dwell time. As can be seen in these figures, the incorporation of the block sizes in the estimation procedure can significantly improve the accuracy.

Figure 1a shows the number of block mismatch error for the block-regularized and Lasso estimators. Evidently, the former substantially outperforms the latter in terms of the correct sparsity recovery. In particular, 252, 260, and 302 sample trajectories are enough to achieve RME  $\leq 0.1\%$  when D is equal to 25, 64, and 121, respectively (notice that the largest instance has more than 9 million parameters to be estimated). However, the Lasso estimator cannot achieve this accuracy with even 2000 sample trajectories.

Figure 1b demonstrates the 2-norm of the estimation error for these estimators. Although the Lasso has a smaller estimation error for d < 200, it is strictly dominated by that of the block-regularized estimator when  $d \ge 200$ .

Finally, we compare the proposed estimator with group Lasso, where the  $\ell_{\infty}$  regularization is replaced by a  $\ell_2$  regularization. Suppose that D=25, and the regularization coefficient for the group Lasso (i.e.  $\ell_1/\ell_2$ -regularized estimator) is chosen as  $\lambda = \sqrt{\frac{0.5(D^2 + D \log(\bar{n} + \bar{m}))}{d}}$  (the constant factor is fine-tuned for this case study). According to Figure 1c, the proposed  $\ell_1/\ell_{\infty}$  slightly outperforms group Lasso in terms of the mismatch error. On the other hand, Figure 1d illustrates that neither of the estimators is superior in terms of the estimation error. As a future research direction, we plan to conduct a more comprehensive study on the group Lasso, and its statistical performance in the context of system identification.

# B. Case Study 2: Power Systems

For the second case study, we consider the frequency control problem for power systems, where the goal is to control the governing frequency of the entire network based on the so-called *swing* equations [47]. Assume that there exist  $N_g$  generators in the system. It is easy to describe the swing equations using the well-known direct current (DC) approximation:

$$M_i\ddot{\theta}_i + D_i\dot{\theta}_i = P_{M_i} - P_{E_i}$$

where  $\theta_i$  is the voltage angle at generator i,  $P_{M_i}$  is the mechanical power input at generator i, and  $P_{E_i}$  denotes the active power injection at the bus connected to generator i. Furthermore,  $M_i$  and  $D_i$  are the inertia and damping coefficients at generator i, respectively. Under the DC approximation, the relationship between active power injection and voltage can be written as:

$$P_{E_i} = \sum_{j \in \mathcal{N}_i} B_{ij} (\theta_i - \theta_j),$$

where  $\mathcal{N}_i$  collects the neighbors of generator i, and  $B_{ij}$  is the susceptance of the line (i, j). After discretization with the sampling time dt, the system of swing equations is reduced to the following dynamical system:

$$x_{i}[t+1] = \left(A_{ii}x_{i}[t] + \sum_{j \in \mathcal{N}_{i}} A_{ij}x_{j}[t]\right) + B_{ii}u_{i}[t] + w_{i}[t],$$

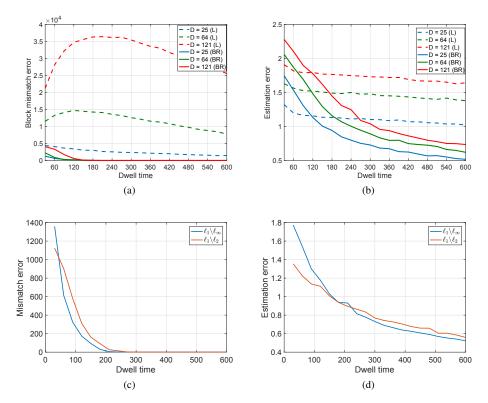


Fig. 1: (a) The block mismatch error for the block-regularized (abbreviated as BR) and Lasso (abbreviated as L) estimators, (b) the estimation error for the block-regularized and Lasso estimators, (d) the normalized estimation error for the block-regularized and Lasso estimators.

where 
$$x_i = \begin{bmatrix} \theta_i & \dot{\theta}_i \end{bmatrix}^\mathsf{T}$$
,  $u_i(t) = P_{M_i}$ , and 
$$A_{ii} = \begin{bmatrix} 1 & dt \\ -\frac{\sum_{j \in \mathcal{N}_i} B_{ij}}{M_i} dt & 1 - \frac{D_i}{M_i} dt \end{bmatrix}$$
,  $A_{ij} = \begin{bmatrix} 0 & 0 \\ \frac{B_{ij}}{M_i} dt & 0 \end{bmatrix}$ ,  $B_{ii} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ .

Realistic power systems are often equipped with an initial distributed controller whose sensing and actuation communication topology is limited by the underlying physical structure of the system [48]. In particular, consider a static distributed controller as follows:

$$u_i[t] = K_{ii}x_i[t] + \sum_{j \in \mathcal{N}_i} K_{ij}x_j[t] + v_i[t],$$
 (34)

where K is a matrix with  $(i,j)^{\text{th}}$  block equal to zero if the generators i and j are not connected. Moreover,  $v_i[t]$  is an exogenous input. Therefore, the closed-loop dynamics of the power system can be written as

$$x_{i}[t+1] = \left(A_{ii}^{c}x_{i}[t] + \sum_{j \in \mathcal{N}_{i}} A_{ij}^{c}x_{j}[t]\right) + B_{ii}v_{i}[t] + w_{i}[t],$$

where

$$A_{ii}^{c} = \begin{bmatrix} 1 & dt \\ -\frac{\sum_{j \in \mathcal{N}_{i}} B_{ij}}{M_{i}} dt + K_{ii}^{1} & 1 - \frac{D_{i}}{M_{i}} dt + K_{ii}^{2} \end{bmatrix},$$

$$A_{ij}^{c} = \begin{bmatrix} 0 & 0 \\ \frac{B_{ij}}{M_{i}} dt + K_{ij}^{1} & K_{ij}^{2} \end{bmatrix},$$
(35)

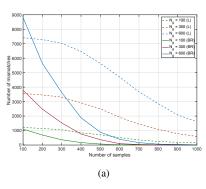
and  $K_{ij} = \begin{bmatrix} K_{ij}^1 & K_{ij}^2 \end{bmatrix}$  for every block  $(i, j)^{\text{th}}$ . Our goal is to identify the closed-loop dynamics of the power system

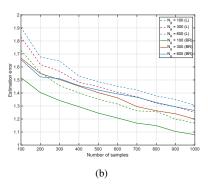
and the underlying topology of the network, based on the sample trajectories collected from the system. Note that the underlying topology structure of the network can be naturally obtained from the block-sparsity structure of  $A^c$ : the block  $A^c_{ij}$  is equal to zero if and only if the generators i and j are not connected. Therefore, the topology inference problem reduces to obtaining the correct block-sparsity pattern of the system matrices  $A^c$  and  $B^c$ . To assess the performance of the proposed method, we generate different instances of the problem according to the following rules:

- the generators are connected via a randomly generated graph with the average degree of 6.
- the parameters  $B_{ij}$ ,  $M_i$ ,  $D_i$  are uniformly chosen from the intervals [0.5, 1], [1, 2], [0.5, 1.5], respectively.
- The nonzero elements of K are uniformly chosen from the interval [0.1, 0.2].

The sampling time dt is set to 0.1. We assume that the disturbance noise has a zero-mean Gaussian distribution with variance 0.01. The mechanical input  $v_i(t)$  is randomly generated according to a zero-mean Gaussian distribution with variance 0.05. In this case study, we compare the performance of the block-regularized and Lasso estimators. The regularization coefficients for these estimators are chosen as  $\sqrt{0.1(D^2 + D\log(\bar{n} + \bar{m}))/d}$  with D = 4 (i.e., the maximum block size), and  $\sqrt{0.01(1 + \log(n + m))/d}$ , respectively.

Figure 2a illustrates the mismatch error of these estimators for different numbers of generators  $N_g$  chosen from  $\{100, 300, 600\}$ . Not surprisingly, the learning time needed to





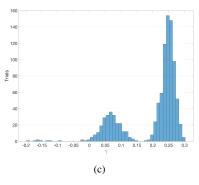


Fig. 2: (a) The block mismatch error for the block-regularized and Lasso estimators, (b) the estimation error for the block-regularized and Lasso estimators, (c) The distribution of mutual incoherency parameter  $\gamma$  over 1000 instances of the problem.

achieve a small mismatch error for both estimators increases as the dimension of the system grows. Conversely, a smaller value for RLT is needed to achieve infinitesimal RME for larger systems. In particular, when  $N_g$  is equal to 100, 300, and 600, the minimum RLT for the proposed block-regularized estimator to guarantee RME  $\leq 0.1\%$  is equal to 0.67, 0.45, and 0.29, respectively. On the other hand, the minimum RLT for the Lasso to achieve the same RME is on average 2.45 times larger than that of the block-regularized estimator.

Figure 2b depicts the 2-norm of the estimation error of the block-regularized and Lasso estimators. In can be seen that the estimation error of the block-regularized estimator is strictly smaller than that of the Lasso, highlighting its superior performance in the block-sparse systems.

Finally, 2c illustrates the distribution of the mutual incoherency parameter  $\gamma$  for 1000 randomly generated instances of power systems with 300 generators. It can be seen that only 1.2% of the instances violate the mutual incoherency condition 14 due to the negative values of  $\gamma$ . This highlights the nonconservativeness of this condition in practice.

# V. CONCLUSION

We consider the problem of identifying the parameters of linear time-invariant (LTI) systems. In many real-world problems, the state-space equation of the system admits a blocksparse representation due to localized or internally limited interactions of its states and inputs. In this work, we leverage this property and introduce a block-regularized estimator to identify the sparse representation of the system. We derive sharp non-asymptotic bounds on the minimum number of input-state data samples to guarantee a small element-wise estimation error. In particular, we show that the number of available sample trajectories can be significantly smaller than the system dimension and yet, the proposed block-regularized estimator can correctly recover the block-sparsity of the state and input matrices and result in a small element-wise error. Through different case studies on switching networks and power systems, we demonstrate the performance of the proposed estimator.

## REFERENCES

- M. Pai, Energy function analysis for power system stability. Springer Science & Business Media, 2012.
- [2] Y. Wang, D. J. Hill, and G. Guo, "Robust decentralized control for multimachine power systems," *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, vol. 45, no. 3, pp. 271–279, 1998.
- [3] I. Sa, M. Kamel, R. Khanna, M. Popovic, J. Nieto, and R. Siegwart, "Dynamic system identification, and control for a cost effective open-source vtol may," arXiv preprint arXiv:1701.08623, 2017.
- [4] M. Arcak, "Synchronization and pattern formation in diffusively coupled systems," in 2012 IEEE 51st IEEE Conference on Decision and Control (CDC). IEEE, 2012, pp. 7184–7192.
- [5] L. Ljung, "System identification," Wiley Encyclopedia of Electrical and Electronics Engineering, pp. 1–19, 1999.
- [6] T. Söderström and P. Stoica, System identification. Prentice-Hall International, 1989.
- [7] P. E. Vértes, A. F. Alexander-Bloch, N. Gogtay, J. N. Giedd, J. L. Rapoport, and E. T. Bullmore, "Simple models of human brain functional networks," *Proceedings of the National Academy of Sciences*, vol. 109, no. 15, pp. 5868–5873, 2012.
- [8] S. Sun, R. Huang, and Y. Gao, "Network-scale traffic modeling and forecasting with graphical lasso and neural networks," *Journal of Trans*portation Engineering, vol. 138, no. 11, pp. 1358–1367, 2012.
- [9] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu, "On the sample complexity of the linear quadratic regulator," Foundations of Computational Mathematics, pp. 1–47, 2019.
- [10] M. Simchowitz, H. Mania, S. Tu, M. I. Jordan, and B. Recht, "Learning without mixing: Towards a sharp analysis of linear system identification," in *Conference On Learning Theory*, 2018, pp. 439–473.
- [11] T. Sarkar and A. Rakhlin, "How fast can linear dynamical systems be learned?" *arXiv preprint arXiv:1812.01251*, 2018.
- [12] S. Oymak and N. Ozay, "Non-asymptotic identification of lti systems from a single trajectory," arXiv preprint arXiv:1806.05722, 2018.
- [13] V. L. Le, F. Lauer, and G. Bloch, "Selective ℓ<sub>1</sub> minimization for sparse recovery," *IEEE Transactions on Automatic Control*, vol. 59, no. 11, pp. 3008–3013, 2014.
- [14] B. M. Sanandaji, T. L. Vincent, M. B. Wakin, R. Tóth, and K. Poolla, "Compressive system identification of lti and ltv arx models," in 2011 50th IEEE Conference on Decision and Control and European Control Conference. IEEE, 2011, pp. 791–798.
- [15] R. Tóth, B. M. Sanandaji, K. Poolla, and T. L. Vincent, "Compressive system identification in the linear time-invariant framework," in 2011 50th IEEE Conference on Decision and Control and European Control Conference. IEEE, 2011, pp. 783–790.
- [16] X. Jiang, Y. Yao, H. Liu, and L. Guibas, "Compressive network analysis," *IEEE transactions on automatic control*, vol. 59, no. 11, pp. 2946–2961, 2014.
- [17] V. Cerone, S. M. Fosson, and D. Regruto, "Sparse linear regression with compressed and low-precision data via concave quadratic programming," in 2019 IEEE 58th Conference on Decision and Control (CDC). IEEE, 2019, pp. 6971–6976.
- [18] C. R. Rojas, R. Tóth, and H. Hjalmarsson, "Sparse estimation of polynomial and rational dynamical models." *IEEE Trans. Automat. Contr.*, vol. 59, no. 11, pp. 2962–2977, 2014.

- [19] T. Chen, M. S. Andersen, L. Ljung, A. Chiuso, and G. Pillonetto, "System identification via sparse multiple kernel-based regularization using sequential convex optimization techniques," *IEEE Transactions* on Automatic Control, vol. 59, no. 11, pp. 2933–2945, 2014.
- [20] R. S. Smith, "Frequency domain subspace identification using nuclear norm minimization and hankel matrix realizations," *IEEE Transactions* on Automatic Control, vol. 59, no. 11, pp. 2886–2896, 2014.
- [21] S. M. Fosson, "Online optimization in dynamic environments: a regret analysis for sparse problems," in 2018 IEEE Conference on Decision and Control (CDC). IEEE, 2018, pp. 7225–7230.
- [22] S. Ross and J. A. Bagnell, "Agnostic system identification for model-based reinforcement learning," arXiv preprint arXiv:1203.1007, 2012.
- [23] E. Weyer, "Finite sample properties of system identification of arx models under mixing conditions," *Automatica*, vol. 36, no. 9, pp. 1291– 1299, 2000.
- [24] R. Dahlhaus, "Graphical interaction models for multivariate time series," *Metrika*, vol. 51, no. 2, pp. 157–172, 2000.
- [25] M. Eichler, "Testing nonparametric and semiparametric hypotheses in vector stationary processes," *Journal of Multivariate Analysis*, vol. 99, no. 5, pp. 968–1009, 2008.
- [26] J. Songsiri and L. Vandenberghe, "Topology selection in graphical models of autoregressive processes," *The Journal of Machine Learning Research*, vol. 11, pp. 2671–2705, 2010.
- [27] M. Zorzi and R. Sepulchre, "Ar identification of latent-variable graphical models," *IEEE Transactions on Automatic Control*, vol. 61, no. 9, pp. 2327–2340, 2015.
- [28] D. F. Ahelegbey, M. Billio, and R. Casarin, "Sparse graphical vector autoregression: a bayesian approach," *Annals of Economics and Statis*tics/Annales d'Économie et de Statistique, no. 123/124, pp. 333–361, 2016.
- [29] J. Pereira, M. Ibrahimi, and A. Montanari, "Learning networks of stochastic differential equations," in *Advances in Neural Information Processing Systems*, 2010, pp. 172–180.
- [30] S. Fattahi and S. Sojoudi, "Data-driven sparse system identification," in 2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2018, pp. 462–469.
- [31] S. Fattahi and S. Sojoudi, "Non-asymptotic analysis of block-regularized regression problem," 2018 IEEE Conference on Decision and Control (CDC), 2018.
- [32] F. Abdi, M. Hasan, S. Mohan, D. Agarwal, and M. Caccamo, "Resecure: A restart-based security protocol for tightly actuated hard real-time systems," *IEEE CERTS*, pp. 47–54, 2016.
- [33] K. Andersson, B. Lennartson, and M. Fabian, "Restarting manufacturing systems; restart states and restartability," *IEEE Transactions on Automa*tion Science and Engineering, vol. 7, no. 3, pp. 486–499, 2010.
- [34] P. Zhao and B. Yu, "On model selection consistency of lasso," *Journal of Machine learning research*, vol. 7, no. Nov, pp. 2541–2563, 2006.
- [35] E. Candes and J. Romberg, "Sparsity and incoherence in compressive sampling," *Inverse problems*, vol. 23, no. 3, p. 969, 2007.
- [36] V. Cerone, S. Fosson, D. Regruto, and A. Salam, "Sparse learning with concave regularization: relaxation of the irrepresentable condition," in 2020 59th IEEE Conference on Decision and Control (CDC). IEEE, 2020, pp. 396–401.
- [37] S. N. Negahban and M. J. Wainwright, "Simultaneous support recovery in high dimensions: Benefits and perils of block ℓ<sub>1</sub>/ℓ<sub>∞</sub>-regularization," *IEEE Transactions on Information Theory*, vol. 57, no. 6, pp. 3841– 3863, 2011.
- [38] M. J. Wainwright, "Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ₁-constrained quadratic programming (lasso)," IEEE transactions on information theory, vol. 55, no. 5, pp. 2183–2202, 2009.
- [39] S. N. Negahban and M. J. Wainwright, "Simultaneous support recovery in high dimensions: Benefits and perils of block ℓ<sub>1</sub>/ℓ<sub>∞</sub>-regularization," *IEEE Transactions on Information Theory*, vol. 57, no. 6, pp. 3841– 3863, 2011
- [40] G. Obozinski, M. J. Wainwright, and M. I. Jordan, "Union support recovery in high-dimensional multivariate regression," in 2008 46th Annual Allerton Conference on Communication, Control, and Computing. IEEE, 2008, pp. 21–26.
- [41] J. Huang, T. Zhang et al., "The benefit of group sparsity," The Annals of Statistics, vol. 38, no. 4, pp. 1978–2004, 2010.
- [42] P. Zhao, G. Rocha, and B. Yu, "Grouped and hierarchical model selection through composite absolute penalties," *Department of Statistics*, UC Berkeley, Tech. Rep, vol. 703, 2006.
- [43] S. Rosset, "Topics in regularization and boosting," Ph.D. dissertation, stanford university, 2003.

- [44] S. Fattahi and S. Sojoudi, "Sample complexity of sparse system identification problem," arXiv preprint arXiv:1803.07753, 2018.
- [45] M. Elad, Sparse and redundant representations: from theory to applications in signal and image processing. Springer Science & Business Media. 2010.
- [46] M. Mesbahi and M. Egerstedt, Graph theoretic methods in multiagent networks. Princeton University Press, 2010.
- [47] X.-F. Wang, Y. Song, and M. Irving, Modern power systems analysis. Springer Science & Business Media, 2010.
- [48] M. Andreasson, D. V. Dimarogonas, K. H. Johansson, and H. Sandberg, "Distributed vs. centralized power systems frequency control," in 2013 European Control Conference (ECC). IEEE, 2013, pp. 3524–3529.
- [49] G. Obozinski, M. J. Wainwright, M. I. Jordan et al., "Support union recovery in high-dimensional multivariate regression," The Annals of Statistics, vol. 39, no. 1, pp. 1–47, 2011.
- [50] N. Simon, J. Friedman, T. Hastie, and R. Tibshirani, "A sparse-group lasso," *Journal of computational and graphical statistics*, vol. 22, no. 2, pp. 231–245, 2013.
- [51] S. N. Negahban, P. Ravikumar, M. J. Wainwright, B. Yu et al., "A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers," *Statistical Science*, vol. 27, no. 4, pp. 538– 557, 2012.

### APPENDIX

# A. Proof of Main Theorem

Let  $\hat{S}_{\mathcal{A}}$  and  $\hat{S}_{\mathcal{A}^c}$  be obtained by removing those blocks of  $\hat{S}$  with indices not belonging to  $\mathcal{A}$  and  $\mathcal{A}^c$ , respectively. The equation (4) can be reformulated as the set of linear equations

$$Y^{(:,j)} = X\Psi^{(:,j)} + W^{(:,j)} \quad \forall j \in \{1, ..., \bar{n}\},$$
 (36)

where  $Y^{(:,j)}$ ,  $\Psi^{(:,j)}$ , and  $W^{(:,j)}$  are the  $j^{th}$  block column of Y,  $\Psi$ , and W, respectively. Based on this definition, consider the following set of block-regularized subproblems:

$$\hat{\Psi}^{(:,j)} = \arg\min \frac{1}{2d} \|Y^{(:,j)} - X\Psi^{(:,j)}\|_{2}^{2} + \lambda_{d} \|\Psi^{(:,j)}\|_{\text{block}}.$$
(37)

Define  $D_j = p_{\text{max}} n_j$ . The next two lemmas are at the core of our proof for Theorem 1. Due to space restrictions, we have deferred their proofs to the extended version of the paper [44].

**Lemma 1** (No false positives). Given arbitrary constants  $c_1, c_2 > 1$ , suppose that  $\lambda_d$  and d are chosen such that

$$\lambda_d \ge \sqrt{\frac{32c_1\lambda_{\max}(\Sigma_w)^2\sigma_{\max}^2}{\gamma^2} \cdot \frac{(D_j)^2 + D_j\log(\bar{n} + \bar{m})}{d}},$$
(38a)

$$d \ge \frac{72c_2\sigma_{\max}^2}{\gamma^2\Lambda_{\min}} \cdot k_j \left(D_j^2 + D_j \log(\bar{n} + \bar{m})\right). \tag{38b}$$

Then, with probability of at least

$$1 - 3\exp\left(-(c_1 - 1)(D_j + \log(\bar{n} + \bar{m}))\right) - 4\exp\left(-(c_2 - 1)(D_j + \log(\bar{n} + \bar{m}))\right)$$
(39)

 $\hat{\Psi}^{(:,j)}$  is unique and its nonzero blocks exclude the zero blocks of  $\Psi^{*(:,j)}$ .

Due to Assumption A4, we have  $n_{\max} = O((\bar{n} + \bar{m})^{\alpha_n})$  and  $k_{\max} = O((\bar{n} + \bar{m})^{\alpha_k})$  for some  $\alpha_n \ge 0$  and  $\alpha_k \ge 0$ .

**Lemma 2** (Element-wise error). Given arbitrary constants  $c_3 > 0$  and  $c_4 > 1$ , suppose that  $\hat{\Psi}$  is unique and the set of its nonzero blocks excludes the zero blocks of  $\Psi^*$ . Then, with probability of at least

$$1 - 2\exp(-(k_j n_j + c_3 \log(\bar{n} + \bar{m}))/2) - 2\exp(-d/2) - 2\exp(-2(c_4 - 1)(\alpha_n + \alpha_k)\log(\bar{n} + \bar{m})))$$
(40)

we have

$$\|\hat{\Psi}^{(:,j)} - \Psi^{*(:,j)}\|_{\infty} \leq \sqrt{\frac{36c_4(\alpha_n + \alpha_k)\lambda_{\max}(\Sigma_w)^2 \log(\bar{n} + \bar{m})}{\Lambda_{\min}d}} + \frac{\lambda_d}{\Lambda_{\min}} \left(8\sqrt{k_j}\sqrt{\frac{k_j n_j + c_3 \log(\bar{n} + \bar{m})}{d}} + 1\right) = g_j. \quad (41)$$

Furthermore, the zero blocks of  $\hat{\Psi}^{(:,j)}$  exclude the nonzero blocks of  $\Psi^{*(:,j)}$  if  $\min_{i \in \mathcal{A}_j} \|\Psi^{(i,j)}\|_{\infty} > g_j$ .

Most existing block-sparsity methods in linear regression focus on the problems where the blocks have row or column dimension of one [39], [41], [49]–[51], and hence, are not applicable to problems with arbitrary block sizes. On the other hand, recall that many large-scale dynamical systems are composed of interacting subsystems, each with its own local states/inputs with potentially different sizes. This imposes a general block structure on different rows and columns of the matrices A and B, and hence, the existing results on blockregularized estimators cannot be readily used in these settings. Lemmas 1 and 2 are precisely aimed to address this issue, and will play key roles in proving the main theorem of this paper. The proofs of Lemmas 1 and 2 are based on the extended version of the so-called primal-dual witness (PDW) approach, which was initially proposed in [38] for element- or row-wise sparse structures. The details of this generalization can be found in the extended version of the paper [44].

## B. Proof of Theorem 1:

First, we present the proof in a few steps: *Step 1:* (9) can be rewritten as follows:

$$\hat{\Psi} = \arg\min_{\Psi} \sum_{j=1}^{n} \left( \frac{1}{2d} \| Y^{(:,j)} - X \Psi^{(:,j)} \|_{2}^{2} + \lambda \| \Psi^{(:,j)} \|_{\text{block}} \right). \tag{42}$$

The above optimization problem can be decomposed into  $\bar{n}$  disjoint block-regularized subproblems in the form of (37).

Step 2: Assume that (38b) and (38a) hold for every  $1 \le j \le \bar{n}$ . Upon defining  $\mathcal{T}_j$  as the event that Lemmas 1 and 2 hold, one can write:

$$\mathbb{P}(\mathcal{T}_{j}) \geq 1 - 5 \exp\left(-(c_{1} - 1)(D_{j} + \log(\bar{n} + \bar{m}))\right)$$

$$-4 \exp\left(-(c_{2} - 1)(D_{j} + \log(\bar{n} + \bar{m}))\right)$$

$$-2 \exp\left(-(k_{j}n_{j} + c_{3}\log(\bar{n} + \bar{m}))/2\right)$$

$$-2 \exp\left(-2(c_{4} - 1)(\alpha_{n} + \alpha_{k})\log(\bar{n} + \bar{m})\right),$$
(43)

for every  $1 \le j \le \bar{n}$ .

Step 3: Assume that  $c_1, c_2, c_4 > 2$  and  $c_3 > 1$ . Consider the event  $\mathcal{T} = \mathcal{T}_1 \cap \mathcal{T}_2 \cap \cdots \cap \mathcal{T}_n$ . Based on (43) and a simple union bound, one can write:

$$\mathbb{P}(\mathcal{T}) \ge 1 - \underbrace{K_{1}(\bar{n} + \bar{m})^{-(c_{1}-2)}}_{(a)} - \underbrace{K_{2}(\bar{n} + \bar{m})^{-(c_{2}-2)}}_{(b)} - \underbrace{K_{3}(\bar{n} + \bar{m})^{-(\frac{c_{3}}{2}-1)}}_{(c)} - \underbrace{K_{4}(\bar{n} + \bar{m})^{-(2(\alpha_{n} + \alpha_{k})(c_{4}-1)-1)}}_{(d)}$$

for some constants  $K_1, K_2, K_3, K_4$ . One can easily verify that the following equalities are enough to guarantee that the right hand side of (44) is equal to  $1 - \delta$ :

$$c_{1} = \frac{\log(4K_{1}/\delta)}{\log(\bar{n} + \bar{m})} + 2, \quad c_{2} = \frac{\log(4K_{2}/\delta)}{\log(\bar{n} + \bar{m})} + 2,$$

$$c_{3} = \frac{2\log(4K_{3}/\delta)}{\log(\bar{n} + \bar{m})} + 2,$$

$$c_{4} = \frac{\log(4K_{4}/\delta)}{2(\alpha_{n} + \alpha_{k})\log(\bar{n} + \bar{m})} + \frac{1}{2(\alpha_{n} + \alpha_{k})} + 1.$$
(45)

Substituting (45) in Lemmas 1 and 2 leads to two observations:

- If  $\lambda_d$  and d satisfy (18a) and (18b), then they also satisfy (38a) and (38b).
- The parameter g defined in (19) is greater than or equal to  $g_j$  for every  $j = 1, ..., \bar{n}$ .

Therefore, (18a) and (18b) guarantee that: 1)  $\hat{\Psi}$  is unique and does not have any false positive in its blocks, and 2) its element-wise error is upper bounded by (19). Now, it only remains to show that  $\hat{\Psi}$  excludes false negatives (the blocks that are mistakenly estimated to have nonzero values). To this goal, it suffices to show that (18b) guarantees  $g < t_{\min}$ . Suppose that

$$d = \Omega\left(C_{\Psi}\kappa(\tilde{\Sigma})^2 k_{\max}\left(D\log(\bar{n} + \bar{m}) + D^2\log(1/\delta)\right)\right).$$
(46)

In what follows, we will show that  $C_{\Psi}$  = O(1) is enough to have  $g < t_{\min}$ . The lower bound on d in (18b) yields that

$$g \le K \left( \frac{1}{\sqrt{C_{\Psi} k_{\max}}} + \frac{1}{C_{\Psi} \kappa(\tilde{\Sigma})} \right),$$
 (47)

for some constant K. Therefore,

$$C_{\Psi} = \frac{2/K}{t_{\min}\kappa(\tilde{\Sigma})} + \frac{4/K}{t_{\min}^2 k_{\max}} = O(1)$$
 (48)

is enough to ensure  $g < t_{\min}$ . This completes the proof.  $\square$ 



Salar Fattahi is an assistant professor in the University of Michigan. Salar's research lies at the intersection of optimization, statistics, and control. He was the recipient of several awards, including 2020 INFORMS ENRE Best Student Paper Award (as a co-author), 2020 Power & Energy Society General Meeting Best-of-the-Best Paper Award, and 2018 INFORMS Data Mining Best Paper Award. He was also a finalist for the 2018 American Control Conference Best Paper Award.



Somayeh Sojoudi Somayeh Sojoudi is an Assistant Professor in the Departments of Electrical Engineering & Computer Sciences and Mechanical Engineering at the University of California, Berkeley. She is an Associate Editor for the journals of the IEEE Transactions on Smart Grid, IEEE Access, and Systems & Control Letters. She is also a member of the conference editorial board of the IEEE Control Systems Society. She received many awards and honors, including INFORMS Optimization Society Prize for Young Researchers, INFORMS Energy

Best Publication Award, INFORMS Data Mining Best Paper Award, NSF CAREER Award, and ONR Young Investigator Award. She has also received several best student conference paper awards (as advisor or co-author) from the Control Systems Society.