# Exact Recovery for System Identification with More Corrupt Data than Clean Data

Baturalp Yalcin, Haixiang Zhang, Javad Lavaei, and Murat Arcak

Abstract—This paper investigates the system identification problem for linear discrete-time systems under adversaries and analyzes two lasso-type estimators. We examine both asymptotic and non-asymptotic properties of these estimators in two separate scenarios, corresponding to deterministic and stochastic models for the attack times. Since the samples collected from the system are correlated, the existing results on lasso are not applicable. We prove that when the system is stable and attacks are injected periodically, the sample complexity for exact recovery of the system dynamics is linear in terms of the dimension of the states. When adversarial attacks occur at each time instance with probability p, the required sample complexity for exact recovery scales polynomially in the dimension of the states and the probability p. This result implies almost sure convergence to the true system dynamics under the asymptotic regime. As a by-product, our estimators still learn the system correctly even when more than half of the data is compromised. We highlight that the attack vectors are allowed to be correlated with each other in this work, whereas we make some assumptions about the times at which the attacks happen. This paper provides the first mathematical guarantee in the literature on learning from correlated data for dynamical systems in the case when there is less clean data than corrupt data.

Index Terms—System Identification, Robust Control, Statistical Learning, Linear Systems, Uncertain Systems

#### I. INTRODUCTION

Dynamical systems serve as the fundamental components in reinforcement learning and control systems. The system dynamics may not be known exactly when the system is complex. Therefore, learning the underlying system dynamics, named the system identification problem, and using the data collected from the system are essential in robotics, control theory, time-series, and reinforcement learning applications. The system identification problem with small disturbances using the least square estimator has been ubiquitously studied, and the literature for this problem is overly rich [1]. Despite several advances in this field, most results in system identification focus on the asymptotic properties, i.e., properties

This work was supported by grants from AFOSR, ARO, ONR, and NSF.

B. Yalcin and J. Lavaei are with the Department of Industrial Engineering and Operations Research, University of California, Berkeley, CA, USA, 94720 (e-mail: baturalp\_yalcin@berkeley.edu; lavaei@berkeley.edu).

H. Zhang is with the Department of Mathematics, University of California, Berkeley, CA, USA, 94720 (e-mail: haixiang\_zhang@berkeley.edu).

M. Arcak is with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA, USA, 94720 (email: arcak@berkeley.edu). of the estimators at infinity, of the proposed estimators only. Nonetheless, the non-asymptotic analysis of the system identification problem has gained interest in recent years [2]–[5]. Although non-asymptotic analysis is harder, it is crucial to understand the required sample complexity for online control problems.

The robust learning of dynamical systems is crucial for safety-critical applications, such as autonomous driving [6], unmanned aerial vehicles [7], and robotic arms [8]. While recent papers have addressed online non-asymptotic control of linear time-invariant (LTI) systems, their applicability often hinges on the assumption of small noise in measurements, neglecting scenarios involving large magnitudes of noise indicative of adversarial attacks or data corruption [9]-[11]. These papers utilize recent advances in high-dimensional statistics and learning theory to analyze the properties of the solution even when the data samples are correlated. The work [12] provides a tutorial on proof techniques. Leastsquare estimators are the main tool in those works, which are susceptible to outliers and large noise in the system. Consequently, we propose two new non-smooth estimators inspired by the lasso problem and robust regression literature [13]. We study the required sample complexity for the exact recovery of LTI systems using these estimators when there are sporadic large disturbance injections to the system.

The robust regression and learning problems under adversaries are ubiquitously studied in the literature [14]-[17]. However, existing methods for analyzing the estimators cannot be directly generalized to control problems due to the correlation between the samples. Therefore, different strategies have been developed recently to tackle this challenge. Firstly, the system is initiated multiple times, and the data point at the end of each run is used to obtain uncorrelated data points, as in [18]. However, obtaining multiple trajectories is not viable and costefficient for most safety-critical applications. One method with a single trajectory relies on the persistent excitation of the states so that the dynamics can be explored thoroughly. This is achieved by injecting a Gaussian noise input into the system. Small ball techniques are used to analyze the properties of the estimator [9], [19], [20]. This technique employs normalized martingale bounds for the estimation error when the excitation is large enough [9].

Unlike the non-asymptotic analysis of correlated data, the least-squares estimator offers a closed-form solution when the system is subjected to small white noise [21]–[23]. As long as the noise magnitudes are not large, the least-squares estimator performs relatively well. The estimation error asymptotically

converges to zero with the optimal rate of  $T^{-1/2}$ , where T is the number of samples collected from the system [9]. However, it is not robust to adversarial attacks, and the literature on robust learning of dynamical systems is limited. The work by [24] defines the null space property (NSP) to analyze a lasso-type estimator for the system. It provides necessary and sufficient conditions for exact recovery when NSP is satisfied, which is NP-hard to check. To circumvent the computational complexity, we build upon [24] and study robust estimators from a non-asymptotic point of view under standard assumptions, such as the system being stable and the attacks being sub-Gaussian.

**Contributions:** We study discrete-time linear time-invariant systems of the form  $x_{i+1} = \bar{A}x_i + \bar{B}u_i + \bar{d}_i$ , where  $\bar{A} \in \mathbb{R}^{n \times n}$  and  $\bar{B} \in \mathbb{R}^{n \times m}$  are unknown matrices of the model. We aim to learn these matrices from the samples  $\{x_i, u_i\}_{i=0}^{T-1}$  of a single initialization of the system when the disturbance vectors  $\bar{d}_i$  are adversarial. Here, the adversarial noise refers to a vector that is designed to deteriorate the performance of the estimator. Thus, the adversarial vectors  $\{\bar{d}_i\}_{i=0}^{T-1}$  can take arbitrarily large finite values, be dependent over time, and can have any undesirable structures. We say that an adversarial attack occurs whenever  $\bar{d}_i$  is non-zero, and we have no information on the value of  $\bar{d}_i$ . If  $\bar{d}_i$  is zero, there is no attack or adversary at time *i*. In our setting, we study systems that are not subject to ordinary minor measurement or modeling errors, and instead the non-zero noise or disturbance stems from an adversarial event.

We study two convex estimators based on the minimization of the  $\ell_2$  and  $\ell_1$  norms of the estimated disturbance vectors,  $\sum_{i=0}^{T-1} ||d_i||_2$  and  $\sum_{i=0}^{T-1} ||d_i||_1$ , with the decision variables *A*, *B*, and  $\{d_i\}_{i=0}^{T-1}$  subject to  $x_{i+1} = Ax_i + Bu_i + d_i$ , given the samples  $\{x_i, u_i\}_{i=0}^{T-1}$ :

$$\min_{A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times m}} \sum_{t=0}^{T-1} \|x_{t+1} - Ax_t - Bu_t\|_{\circ}, \quad \circ \in \{1, 2\}.$$

This is equivalent to an empirical risk minimization problem for which the loss function is the  $\ell_1$  and  $\ell_2$  norms depending on the choice of  $\circ$ . We employ a non-smooth objective function to obtain a robust estimator. The arbitrary injection of adversaries may happen infrequently in time. In that case, the attacks occur sparsely in time. Conversely, the vector  $d_i$  at each attack time *i* could be dense, and there is no limitation on how sparse the vector is. The  $\ell_2$  norm estimator is the most effective in this case. In contrast, the  $\ell_1$  norm estimator is preferable if the vector  $d_i$  at each attack time is structured and known to be sparse. We summarize our contributions below.

i) We first consider the case when the adversarial noise injections, i.e., adversarial attacks, happen periodically over time with the period  $\Delta$ . We show that both of our estimators exactly recover the true system matrices  $\overline{A}$  and  $\overline{B}$  when the system is stable and the number of samples, i.e., T, is larger than  $n + \Delta$ .

**ii**) We then consider a probabilistic model for the occurrence of attacks, in which there is an arbitrary noise injection at each time instance i with probability p, independent of previous time periods. Nevertheless, we allow these noise injections, or attack vectors, to be dependent. We study the required

sample complexity of our estimators for exact recovery when the attack vectors are stealthy. Suppose that the adversarial noise and the input sequence are sub-Gaussian random vectors and possibly dependent. Then, the estimators achieve exact recovery with probability at least  $1 - \delta$  if the time horizon *T* satisfies the inequality  $T \ge \Theta(\max\{T_{\text{sample}}^1, T_{\text{sample}}^2\})$ , where  $T_{\text{sample}}^1$  and  $T_{\text{sample}}^2$  are defined as

and

$$nmR_2\log\left(\frac{nR_2}{\delta}\right),$$

 $n^2 R_1 \log\left(\frac{nR_1}{\delta}\right),$ 

with the constants  $R_1$  and  $R_2$  defined in Theorem 4.

**iii)** As a corollary to the previous result, we show that the estimators converge to true system matrices almost surely when the attack vectors are stealthy. Otherwise, if the attack vectors are not stealthy, the system operator could detect the abnormalities and stop the system, which is not a desired outcome for the adversarial agent or attacker. This is the first paper that studies the adversarial attack structure for the system identification problem to obtain sample complexity using non-asymptotic analysis techniques.

This paper is organized as follows. In Sections 2 and 3, we introduce the notations used in the paper and formulate the problem, respectively. In Section 4, we study the convergence and sample complexity properties of our estimators in the case when the system is autonomous. In Section 5, we generalize the results to non-autonomous systems. In Section 6, we demonstrate the results on a biomedical system that models blood sugar levels with the injection of bolus insulin. This work provides the first bound in the literature on sample complexity for dynamical systems under adversaries, and its techniques can be adopted to study other robust online learning problems.

## **II. NOTATION AND PRELIMINARIES**

For a matrix Z,  $||Z||_F$  denotes the Frobenius norm of a matrix. For a vector z,  $||z||_1$ ,  $||z||_2$ , and  $||z||_{\infty}$  denote its  $\ell_1$ ,  $\ell_2$ , and  $\ell_{\infty}$  norms, respectively. Given two functions f and g, the notation  $f(x) = \Theta[g(x)]$  means that there exist universal positive constants  $c_1$  and  $c_2$  such that  $c_1g(x) \le f(x) \le c_2g(x)$ . The relation  $f(x) \leq g(x)$  holds if there exists a universal positive constant  $c_3$  such that  $f(x) \le c_3 g(x)$  holds with high probability when T is large. The relation  $f(x) \gtrsim g(x)$  holds if  $g(x) \leq f(x)$ . |S| shows the cardinality of a given set S. For two vectors v and w,  $\langle v, w \rangle$  is the inner product between those vectors in their respective vector space. Furthermore, we use the notation  $v \otimes w = vw^T$  to denote the outer product.  $\mathbb{P}(\cdot)$  and  $\mathbb{E}[\cdot]$  denote the probability of an event and the expectation of a random variable. A Gaussian random variable X with mean  $\mu$  and covariance matrix  $\Sigma$  is written as  $X \sim N(\mu, \Sigma)$ . Since we restrict the disturbance vectors to be sub-Gaussian, we formally define them below.

Definition 1 (Sub-Gaussian Random Variable [25]): A random variable  $X \in \mathbb{R}$  with mean  $\mu = \mathbb{E}[X]$  is sub-Gaussian

with parameter  $\sigma$  if

$$\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\lambda^2 \sigma^2/2}, \quad \forall \lambda \in \mathbb{R}.$$

Moreover, a random vector  $X \in \mathbb{R}^n$  with mean  $\mu = \mathbb{E}[X]$  is sub-Gaussian with parameter  $\sigma$  if

$$\mathbb{E}[e^{\lambda \langle \mathbf{v}, X - \mu \rangle}] \le e^{\lambda^2 \sigma^2/2}, \quad \forall \lambda \in \mathbb{R}, \mathbf{v} \in \mathbb{R}^n, \|\mathbf{v}\|_2 = 1$$

Informally, a sub-Gaussian random variable with parameter  $\sigma$  has the property that its tails are less dense than those of a Gaussian random variable with variance  $\sigma^2$ . We will utilize concentration bounds for sub-Gaussian random variables to verify that the optimality conditions for our proposed estimators are satisfied with high probability. The main concentration inequality for sub-Gaussian random variables is Hoeffding's bound.

Lemma 1: (Hoeffding's Bound [25]) Suppose that the variable X has mean  $\mu$  and sub-Gaussian parameter  $\sigma$ . Then, for all t > 0, we have

$$\mathbb{P}(|X-\mu|>t) \leq 2\exp\left(-\frac{t^2}{2\sigma^2}\right).$$

We use the union bound over the set of coordinates and other sets with finite cardinality. Let *S* be a set with finite cardinality,  $|S| < \infty$ , and  $E_i$  be the event related to element *i* in the set *S*. Then, we can write the union bound as

$$\mathbb{P}\left(\cup_{i\in S}E_i\right)\leq \sum_{i\in S}\mathbb{P}(E_i)$$

Since we use non-smooth objective functions with  $\ell_1$  and  $\ell_2$  norms, we introduce the subdifferentials of the  $\ell_1$  and  $\ell_2$  norms.

Definition 2 (Subdifferential of  $\ell_2$  Norm): Given a vector  $z \in \mathbb{R}^n$ , the subdifferential of  $||z||_2$  is denoted as  $\partial ||z||_2$  and is given as

$$\partial \|z\|_2 = \begin{cases} \frac{z}{\|z\|_2}, & \text{if } z \neq 0, \\ \mathbb{B}_2(1), & \text{otherwise.} \end{cases}$$

where  $\mathbb{B}_2(1) = \{x \in \mathbb{R}^n : ||x||_2 \le 1\}$  is the  $\ell_2$  norm unit ball.

Definition 3 (Subdifferential of  $\ell_1$  Norm): Given a vector  $z \in \mathbb{R}^n$  with entries  $z_i, i = 1, ..., n$ , the subdifferential of the  $||z||_1$  is denoted as  $\partial ||z||_1$  and is given as

$$\partial \|z\|_{1}^{i} = \begin{cases} 1, & \text{if } z_{i} > 0, \\ -1, & \text{if } z_{i} < 0, \\ [-1,1], & \text{otherwise} \end{cases}$$

where  $\partial ||z||_1^i$  is the *i*-th coordinate of the subdifferential of  $||z||_1$ .

Note that while the subdifferential of the  $\ell_1$  norm is coordinate-wise separable, the subdifferential of the  $\ell_2$  norm is not coordinate-wise separable. Whenever the vector z is equal to 0, the subdifferential of the  $\ell_2$  norm is the  $\ell_2$  norm unit ball, whereas the subdifferential of the  $\ell_1$  norm is the  $\ell_{\infty}$  norm unit ball, which is

$$\mathbb{B}_{\infty}(1) = \{ x \in \mathbb{R}^n : ||x||_{\infty} \le 1 \}.$$

We also define the unit ball  $S_2(1)$  as

$$\mathbb{S}_2(1) = \{ x \in \mathbb{R}^n : ||x||_2 = 1 \}$$

that is the set of all the points on the sphere with radius 1.

The asymptotic analysis of the system identification problem concerns the convergence rate to the true parameter at an infinite-time horizon. However, historically, asymptotic analysis has not provided the required sample complexity to obtain a solution within a given error tolerance. In contrast, non-asymptotic analysis deals with the finite-time behavior of the estimators using learning theory and high-dimensional statistics. It provides the required sample complexity to bound the estimation error within the specified tolerance with high probability. Consequently, non-asymptotic analysis is more challenging than asymptotic analysis. Our goal is to provide the minimum required number of samples to recover the true parameters of the system with a high probability using techniques designed for non-asymptotic analysis.

### **III. PROBLEM FORMULATION**

We consider a linear time-invariant dynamical system over the time horizon [0, T],  $x_{i+1} = \bar{A}x_i + \bar{B}u_i + d_i$ , i = 0, 1, ..., T-1, where  $\bar{A} \in \mathbb{R}^{n \times n}$  and  $\bar{B} \in \mathbb{R}^{n \times m}$  are unknown system matrices, and  $\bar{d_i} \in \mathbb{R}^n$  are unknown system disturbances. Given the set of state measurements  $\{x_i\}_{i=0}^T$  and the set of inputs  $\{u_i\}_{i=0}^{T-1}$ , the goal is to estimate the unknown system matrices  $\bar{A}$  and  $\bar{B}$ . In this paper, the disturbance vectors  $\{\bar{d}_i\}_{i=0}^{T-1}$  can be engineered to be large if there is an outside attack on the system from an agent or there is a sensor/actuation fault that leads to major corruption in the system dynamics. Throughout the paper, the disturbance vectors  $\{\bar{d}_i\}_{i=0}^{T-1}$  are also called (adversarial) attack vectors. Moreover, the agent who engineers the disturbance vectors is called an attacker. As opposed the majority of the literature, we assume that the disturbance vectors  $\{\bar{d}_i\}_{i=0}^{T-1}$  can be dependent on the disturbance vectors from the previous time instances and there is no specific distribution assumption for these vectors except the sub-Gaussian assumption. We represent the time indices of the attacks or large disturbance vectors with the set  $\mathscr{K}$ , that is  $\mathscr{K} = \{i : \overline{d_i} \neq 0, i \in [0, 1, \dots, T-1]\}$ 1}. These time instances are called the attack times and  $\mathcal{K}$ is the set of attack times. Similarly, the set of time instances without attack or corrupted data is shown with  $\mathscr{K}^c = \{i : \overline{d_i} =$  $0, i \in \{0, 1, \dots, T-1\}$ . These time instances are called the noattack times, and  $\mathcal{K}$  is the set of no-attack times. The data corresponding to attack times are corrupted, whereas the data corresponding to no-attack times are uncorrupted.

We establish the exact recovery of the proposed estimators when there are large disturbances in the system. In such cases, the least-squares method cannot achieve exact recovery, a fact that can be easily verified from its closed-form solution. Define the matrices  $X := [x_0, ..., x_{T-1}]$  and  $\overline{D} := [\overline{d_0}, ..., \overline{d_{T-1}}]$ . The solution for the least-squares problem is  $\widehat{A} = (\overline{A}X + \overline{D})^T X (X^T X)^{-1}$  in the absence of the input sequence  $\{u_i\}_{i=0}^{T-1}$ . Thus, the estimation error is  $\|\overline{D}^T X (X^T X)^{-1}\|$ , which is nonzero and arbitrarily large in the presence of arbitrarily large disturbance vectors. A similar calculation can be made in the presence of an input sequence. Consequently, the leastsquares estimator cannot achieve a zero estimation error, leading to a plateau in the estimation error of the leastsquares estimator in our numerical experiments in Section 6. We define the matrix  $D := [d_0, ..., d_{T-1}]$  with its columns being estimated disturbances, as well as the norms of matrices  $||D||_{1,1} := \sum_i ||d_i||_1$ , and  $||D||_{2,1} := \sum_i ||d_i||_2$ . To exactly recover the system matrices  $\overline{A}$  and  $\overline{B}$ , we analyze the following convex optimization problems with non-smooth objective functions:

$$\min_{\substack{A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times m}, \\ D \in \mathbb{R}^{n \times T}}} \|D\|_{2,1}$$
(CO-L2)
  
s.t.  $x_{i+1} = Ax_i + Bu_i + d_i, \quad i = 0, \dots, T-1,$ 

and

$$\min_{\substack{A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times m}, \\ D \in \mathbb{R}^{n \times T}}} \|D\|_{1,1}$$
(CO-L1)  
s.t.  $x_{i+1} = Ax_i + Bu_i + d_i, \quad i = 0, \dots, T-1,$ 

where the states  $\{x_i\}_{i=0}^T$  are generated according to  $x_{i+1} = \bar{A}x_i + \bar{B}u_i + \bar{d}_i$ , i = 0, ..., T - 1. The difference between problems (CO-L2) and (CO-L1) is their objective functions. Note that these two problems are equivalent when we have a firstorder system with  $x_i \in \mathbb{R}, i \in 0, ..., T-1$ . In problem (CO-L2), the sum of the  $\ell_2$  norm columns is analogous to the  $\ell_1$  norm minimization in the lasso problem. In other words, the  $\ell_1$  norm is applied at the group level to  $\{d_i\}_{i=0}^{T-1}$  because the occurrence of large injections of disturbances is rare and not frequent. We highlight that the vectors  $\{\bar{d}_i\}_{i=0}^{T-1}$  are not necessarily sparse. On the other hand, the  $\ell_1$  norm is applied both at the group level and the in-group levels to  $\{d_i\}_{i=0}^{\hat{T}-1}$  for problem (CO-L1). For those applications that the disturbance vectors can be assumed to be sparse, (CO-L1) is more suitable than (CO-L2). Furthermore, the states  $x_i$  are correlated to each other due to the system dynamics, which makes the nonasymptotic analysis of the problem more challenging than the robust regression literature for which the samples are assumed to be independently generated. One can write the optimization problems (CO-L2) and (CO-L1) as follows using the  $\ell_2$  and  $\ell_1$  norms, respectively:

$$\min_{\mathbf{A}\in\mathbb{R}^{n\times n}, \mathbf{B}\in\mathbb{R}^{n\times m}} \sum_{t=0}^{T-1} \|x_{t+1} - Ax_t - Bu_t\|_{\circ}.$$

This is equivalent to an empirical risk minimization problem for which the loss function is the  $\ell_1$  or  $\ell_2$  norm, depending on the choice of  $\circ$ . Although these types of sum-of-norm minimization non-smooth loss functions are utilized in other applications, this paper marks the first non-asymptotic analysis of these loss functions in the context of control and system identification with serially correlated data.

We remark that classical statistical theory on empirical risk minimization is not applicable to the problem under study in this paper due to the correlated data at each time instance. By representing the data points  $X_t$  as tuples  $(x_{t+1}, x_t, \mu_t)$ , it is impossible to claim that  $X_t$  and  $X_{t+1}$  are independent, which is a key assumption in the empirical risk minimization literature. As the first step of our proof technique, the Karush-Kuhn-Tucker (KKT) conditions will be used to analyze the properties of these estimators. Since (CO-L2) and (CO-L1) are convex optimization problems with linear equalities, the KKT conditions are necessary and sufficient to guarantee optimality, as stated below.

Theorem 1: Consider the convex optimization problems (CO-L2) and (CO-L1) and let  $o \in \{1,2\}$ . Given a pair of matrices  $(\hat{A}, \hat{B})$ , if the following conditions hold simultaneously

$$0 \in \sum_{i \notin \mathscr{K}} x_i \otimes \partial \| (\bar{A} - \hat{A}) x_i + (\bar{B} - \hat{B}) u_i \|_{\circ} + \sum_{i \in \mathscr{K}} x_i \otimes \partial \| (\bar{A} - \hat{A}) x_i + (\bar{B} - \hat{B}) u_i + \bar{d}_i \|_{\circ}, \quad (1)$$

$$0 \in \sum_{i \notin \mathscr{K}} u_i \otimes \partial \| (\bar{A} - \hat{A}) x_i + (\bar{B} - \hat{B}) u_i \|_{\circ} + \sum_{i \in \mathscr{K}} u_i \otimes \partial \| (\bar{A} - \hat{A}) x_i + (\bar{B} - \hat{B}) u_i + \bar{d}_i \|_{\circ}, \quad (2)$$

then  $(\hat{A}, \hat{B})$  is a solution to (CO-L1) when  $\circ = 1$  and a solution to (CO-L2) when  $\circ = 2$ .

The proof for the KKT conditions when  $\circ = 2$  is provided in [26], and the proof for the case  $\circ = 1$  can be done similarly. We will utilize the conditions above to study in what scenarios the exact recovery is achievable. As a simple corollary to Theorem 1, we can state that  $(\bar{A}, \bar{B})$  is a solution to our estimator(s) if the following conditions hold:

$$0 \in \sum_{i \notin \mathscr{K}} x_i \otimes \partial ||0||_{\circ} + \sum_{i \in \mathscr{K}} x_i \otimes \partial ||\bar{d}_i||_{\circ},$$
$$0 \in \sum_{i \notin \mathscr{K}} u_i \otimes \partial ||0||_{\circ} + \sum_{i \in \mathscr{K}} u_i \otimes \partial ||\bar{d}_i||_{\circ}.$$

# IV. AUTONOMOUS SYSTEMS

In this section, we consider autonomous systems, meaning that  $u_0 = \cdots = u_{T-1} = 0$ . Therefore, the system dynamics could be written as  $x_{i+1} = \bar{A}x_i + \bar{d}_i$  for  $i = 0, \dots, T-1$ . Throughout this section, we assume that the system is stable and that it is initialized at the origin.

Assumption 1: Given an autonomous system  $x_{i+1} = \bar{A}x_i + \bar{d}_i$ for i = 0, ..., T - 1 with dimension *n*, assume that  $x_0 = 0$  and all eigenvalues of  $\bar{A}$  are inside the unit circle.

The stability assumption is standard in system identification problems to avoid an unbounded growth of the states during the learning process. Without loss of generality, we initialize the trajectories at the origin since an initialization at other points affects the results only with a constant factor. We study noiseless systems under an adversary to obtain exact recovery results, meaning that if there is no attack at time  $i, i \in \mathscr{K}^c$ , then  $\bar{d_i} = 0$ .

In the noisy case, one can consider the following setup. If there is no attack at time  $i, i \in \mathcal{K}^c$ , then  $\overline{d_i}$  is likely non-zero with a small variance and its value is independent of those for other time periods. If there is an attack at time  $i \in \mathcal{K}$ , then  $\overline{d_i}$  is a combination of two terms: a small noise vector that is Gaussian and independent of past time periods, and a large noise vector that could have an arbitrary distribution and possibly be dependent on past time instances. The noisy case, where the system is subjected to small independent and identically distributed Gaussian errors due to measurements and modeling errors, in addition to the adversarial vectors, can be easily addressed using our framework. The perturbation analysis allows us to bound how far the recovered solution is from the true solution in terms of the values of small noise vectors.

Therefore, we only study the noiseless case as described above. Thus, we are interested in recovering the system matrix  $\overline{A}$  using the following convex optimization problems for autonomous systems:

$$\min_{\substack{A \in \mathbb{R}^{n \times n} \\ D \in \mathbb{R}^{n \times T'}}} \sum_{i=0}^{T-1} \|d_i\|_2 \tag{CO-L2-Aut}$$
s.t.  $x_{i+1} = Ax_i + d_i$ ,

and

$$\min_{\substack{A \in \mathbb{R}^{n \times n} \\ D \in \mathbb{R}^{n \times T'}}} \sum_{i=0}^{T-1} \|d_i\|_1 \tag{CO-L1-Aut}$$
s.t.  $x_{i+1} = Ax_i + d_i$ .

The optimality conditions for problem (CO-L2-Aut) with  $\circ = 2$  and problem (CO-L1-Aut) with  $\circ = 1$  can be written as follows using Theorem 1:

$$0 \in \sum_{i \notin \mathscr{K}} x_i \otimes \partial \| (\bar{A} - A) x_i \|_{\circ} + \sum_{i \in \mathscr{K}} x_i \otimes \partial \| ((\bar{A} - A) x_i + \bar{d}_i) \|_{\circ}.$$
(3)

As a remark, although the set of attack times  $\mathcal{K}$  appears in the optimality conditions, this set is not known a priori to the system operator. The set is only used during the analysis of the proposed estimators to derive sufficient conditions for exact recovery.

We first consider first-order systems where  $x_i, \bar{d_i} \in \mathbb{R}, i = 0, 1, ..., T - 1$  and  $\bar{A} \in \mathbb{R}$ . We examine the first-order case to gain some insight into the ideas behind the proof techniques for general systems. When n = 1, the problems (CO-L1-Aut) and (CO-L2-Aut) are equivalent, and therefore, we only focus on (CO-L2-Aut). After establishing the optimality conditions for these problems, we will examine two types of attack structures. An attack structure refers to the pattern of attack occurrences. In other words, it involves the distribution of each time instance at which a large disturbance vector is injected into the system. Namely, we inspect the structure of the set  $\mathcal{K}$ .

The first attack structure is a deterministic attack model for which the attacks occur at every  $\Delta$  time period. For instance, if  $\Delta = 2$ , the set  $\mathscr{K}$  could be  $\{1,3,5,\ldots,2k+1\}$ , meaning that an agent injects a disturbance vector into the system at every odd time instance. Later, we investigate a probabilistic attack structure where each attack may occur with probability p at each time instance *i*, independent of the past periods. We first define the deterministic attack model, borrowed from [26].

Definition 4 ( $\Delta$ -spaced Attack Structure): Given a positive integer  $\Delta > 2$ , the disturbance sequence  $\{\bar{d}_i\}_{i=0}^{T-1}$  is said to be  $\Delta$ -spaced if for every  $i \in \{0, 1, \dots, T - \Delta - 1\}$  such that  $\bar{d}_i \neq 0$ , we have  $\bar{d}_j = 0$ , for all  $j \in \{i+1,\dots,i+\Delta-1\}$  and  $\bar{d}_{i+\Delta} \neq 0$ . In addition, for  $i \in \{0, 1, \dots, \Delta - 1\}$ , we must have at least one non-zero disturbance vector, i.e.  $\bar{d}_i \neq 0$ .

We will show that the convex formulation (CO-L2-Aut) exactly recovers  $\overline{A}$  in the case of  $\Delta$ -spaced disturbance sequence with  $\Delta \geq 2$ . Proposition 1: Consider a first-order autonomous system with  $\Delta$ -spaced disturbance sequence with  $\Delta \geq 2$ . Then, the convex formulation (CO-L2-Aut) (or equivalently (CO-L1-Aut)) has the unique solution  $\overline{A}$  as long as the sample complexity satisfies the inequality  $T \geq \Delta + 1$ .

This proposition implies that whenever there are more than  $\Delta + 1$  data samples, the exact recovery is guaranteed to be achieved. Note that Proposition 1 does not make any assumption on the vector set  $\{\overline{d}_i : i \in \mathcal{K}\}$  and each element of the set could be arbitrarily large and correlated as long as they are finite. As a result, regardless of the severity of the attack, an exact recovery is guaranteed for (CO-L1-Aut) and (CO-L2-Aut). One important implication of Proposition 1 is for the case where there is a  $\Delta$ -spaced disturbance sequence with  $\Delta = 2$ , meaning that half of the observations are corrupted. In the robust regression estimation literature, exact recovery is possible only if the number of attacked observations is less than half of the total observations. The main difference between robust regression and system identification problems is that the observations are correlated with each other in the latter. This enables exact recovery for the convex formulation even if half of the data is corrupted via an adversarial agent. The proof of Proposition 1 is based on the following lemma.

*Lemma 2:* (Theorem 1 in [26]) Consider the convex optimization problem (CO-L2-Aut). If  $\sum_{i \notin \mathcal{H}} |x_i| > \sum_{i \in \mathcal{H}} |x_i|$ , then  $\overline{A}$  is the unique solution to the problem.

The proof of Lemma 2 is based on the KKT conditions of the problem provided earlier. A natural question arises as to whether one can generalize the above result to higherorder systems. The next proposition extends Proposition 1 to autonomous dynamical systems with an arbitrary order *n* under a  $\Delta$ -spaced disturbance sequence with  $\Delta \ge n+1$ .

Proposition 2: Consider an autonomous system of order n under a  $\Delta$ -spaced disturbance sequence with  $\Delta \ge n+1$ . Suppose that  $\bar{A}$  is diagonalizable with eigenvalues,  $\bar{\lambda}_l, l = 1, 2, ..., n$ , and that the condition

$$\bar{d}_{i+\Delta} \in \operatorname{span}\{\bar{d}_i, \bar{A}\bar{d}_i, \dots, \bar{A}^{\Delta-2}\bar{d}_i\}, \quad \forall i = 0, 1, \dots, T-1 \quad (4)$$

is satisfied. Then,  $\overline{A}$  is a solution to the convex formulation (CO-L2-Aut) if  $T \ge n + \Delta$ , provided that

$$\left|\sum_{k_1+\dots+k_n=\Delta-n}\bar{\lambda}(k_1,\dots,k_n)\right| \leq \sum_{t=0}^{\Delta-n-1} \left|\sum_{k_1+\dots+k_n=t}\bar{\lambda}(k_1,\dots,k_n)\right|,\tag{5}$$

where the notation  $\bar{\lambda}(k_1, \ldots, k_n)$  denotes  $\bar{\lambda}_1^{k_1} \times \bar{\lambda}_2^{k_2} \times \cdots \times \bar{\lambda}_n^{k_n}$ .

This result is a generalization of Proposition 1, and we do not require all the eigenvalues of  $\overline{A}$  to lie inside the unit circle (i.e., it allows the violation of Assumption 1). The condition (4) is necessary to ensure that the KKT condition is satisfied, which eliminates the alignment of the attack vectors with eigenspaces of the matrix  $\overline{A}$ . In real-life applications, this circumstance can be avoided by injecting a small perturbation to the system. To gain insight into equation (5), which involves the product of eigenvalues, consider a special case where  $\overline{A}$  has the eigenvalue  $\lambda$  with multiplicity *n* and *n* distinct

$C_{n,k}$	k=1	k=2	k=3	k=5	k=7	k=10
n=1	1.0000	1.6180	1.8393	1.9659	1.9920	1.9990
n=2	0.5000	1.0000	1.2886	1.5725	1.7010	1.7951
n=3	0.3300	0.7287	1.0000	1.3181	1.4892	1.6310
n=5	0.2000	0.4740	0.6938	1.0000	1.1956	1.3087
n=7	0.1429	0.3516	0.5320	0.8069	1.0000	1.1979
n=10	0.1000	0.2535	0.3944	0.6263	0.8036	1.0000

Fig. 1. Upper-Bound Value  $C_{n,k}$  for Different Values of *n* and *k*.

eigenvectors. In this case, we can simplify (5) as follows. Define  $k := \Delta - n$ . Then, (5) is equivalent to

$$\binom{n+k-1}{k}|\lambda|^k-\sum_{i=0}^{k-1}\binom{n+i-1}{i}|\lambda|^i<0$$

This condition is satisfied if  $|\lambda| \leq C_{n,k}$ , where  $C_{n,k}$  denotes the upper bound on the eigenvalue magnitudes given the parameters *n* and *k*. Figure 1 summarizes the values of  $C_{n,k}$  for different choices of *n* and *k*. Note that  $C_{n,k} \leq C_{m,k}$  if n > m and  $C_{n,k} \leq C_{n,l}$  if k < l, due to the definition of  $C_{n,k}$ . It can be shown that  $C_{1,k} \rightarrow 2$  as  $k \rightarrow \infty$ . As a result,  $|\lambda| \leq C_{n,k} \leq C_{1,k} \rightarrow 2$ . This shows that the stability of the system is not necessary for exact recovery when the attack vectors are injected less frequently. In addition, whenever k = n or  $\Delta = 2n$ ,  $|\lambda| < 1$  is sufficient for exact recovery as suggested by Proposition 2. This conclusion is analogous to the stability of the system. Proposition 2 can still be applied to problem (CO-L1-Aut). However, the KKT conditions will differ due to the subdifferential of the  $\ell_2$  and  $\ell_1$  norms. In fact, they both have a similar shape. Therefore, one can show that this proposition still holds with the same condition even if convex formulation (CO-L1-Aut) with the  $\ell_1$ norm of the disturbance vectors is used.

It is natural to ask whether it is possible to learn the system when there is more corrupted data than clean data. We cannot use a  $\Delta$ -spaced disturbance sequence model because the minimum value of  $\Delta$  is 2, which does not allow the size of corrupted data to exceed the size of clean data. Thus, we investigate a probabilistic attack structure. In this structure, a non-zero disturbance vector  $d_i$  is injected into the system at time instance *i* with probability p > 0, which is independent of the past and future time periods. To address this, we consider a probabilistic attack model where there is a parameter p specifying the probability of an attack at each time instance. Specifically, given a time instance *i*,  $d_i$ is non-zero with probability p, and this is independent of all previous and future time instances. As a result, the event of having an attack at each time instance is identically and independently distributed with a Bernoulli distribution with parameter p. Nevertheless, the attack vectors are still allowed to be correlated with each other. Our goal is to discover the properties of (CO-L1-Aut) and (CO-L2-Aut) for an arbitrary value of p, especially p > 0.5. We make the following stealth attack assumption.

Assumption 2: For each  $k \in \mathcal{K}$ , the attack vector is defined by

$$\bar{d}_k := \bar{\ell}_k \bar{f}_k$$
, where  $\bar{\ell}_k \in \mathbb{R}$  and  $\bar{f}_k \in \mathbb{S}_2(1)$ 

where  $\bar{f}_k$  plays the role of the direction of the attack while  $\bar{\ell}_k$  plays the role of the length (that is allowed to take negative

values too). Define the filtration

$$\mathscr{F}_k := \sigma\{x_1,\ldots,x_k\}, \quad \forall k \in \{0,\ldots,T-1\}.$$

For all  $k \in \mathcal{K}$ , conditioning on  $\mathcal{F}_k$ , the following statements hold:

- 1)  $\bar{\ell}_k$  is independent from the direction  $\bar{f}_k$ ;
- 2) The direction  $\bar{f}_k$  obeys the uniform distribution on  $\mathbb{S}_2(1)$ ;
- 3)  $\bar{\ell}_k$  is mean-zero and sub-Gaussian with parameter  $\sigma$ ;
- 4) The variance of  $\bar{\ell}_k$  is  $\sigma_k^2 \in [c^2 \sigma^2, \sigma^2]$  for some constant c > 0.

Under the stealth assumption, the length  $\bar{\ell}_k$  can depend on the previous attacks  $\bar{d}_{k'}$ , and in particular  $\bar{\ell}_{k'}$  and  $\bar{f}_{k'}$  for k' < k. In addition, we note that the above assumption of symmetry of the disturbance vectors reflected in  $f_k$  is not restrictive and corresponds to stealth attacks. If this assumption does not hold, the attacks may be detectable, and their effects could be nullified, or the system could be stopped to investigate the possible influence from outside agents. For an attack to be stealthy, its value should be zero on expectation, and our assumption has a similar flavor. If the symmetric assumption does not hold, it has been shown that there is a bias in estimation, and there is no way to avoid this bias [27]. In the special cases when the length distribution is Gaussian or bounded, the constant c is equal to 1. Furthermore, we mention that the uniform distribution assumption of  $\bar{f}_k$  can be relaxed to an arbitrary distribution on the sphere with zero mean and fullrank covariance matrix. In that more general case, the sample complexity in Theorems 2-5 will depend on the conditional number of the covariance matrix, which is equal to 1 under Assumption 2.

Since the KKT conditions include random variables and random sets due to the randomness in the attack structure, it is not possible to obtain deterministic sample complexity for exact recovery as in Proposition 2. Therefore, it is essential to quantify the required number of samples for exact recovery with high probability using non-asymptotic analysis. Under Assumption 2, the attack vector at time *i*,  $\overline{d_i}$ , has a sub-Gaussian distribution with parameter  $\sigma$  given  $\mathscr{F}_i$ , as described in Assumption 2. The sub-Gaussianity assumption does not specify the distribution of the disturbance vector but assures that the disturbance vectors have light tails. For instance, any distribution over a bounded space is sub-Gaussian, making this assumption extremely mild. As a result, the sub-Gaussian assumption is not restrictive.

The KKT conditions for exact recovery, which are necessary and sufficient, can be restated as

$$\exists \gamma_i \in \partial \|0\|_{\circ}, \quad \forall i \notin \mathscr{K} \quad \text{s.t.} \ \sum_{i \notin \mathscr{K}} x_i \otimes \gamma_i = \sum_{i \in \mathscr{K}} x_i \otimes \partial \|\bar{d}_i\|_{\circ}.$$

because of the properties of the subdifferentials at the origin. In order to simplify the analysis, we use the relationship between the unit balls of the  $\ell_{\infty}$  and  $\ell_2$  norms, that is  $\frac{1}{\sqrt{n}}\mathbb{B}_{\infty}(1) \subseteq \mathbb{B}_2(1)$ . Additionally, we examine the results for each coordinate of the subdifferentials since they are separable due to the properties of the  $\ell_{\infty}$  norm. Therefore, the following propositions provide sufficient conditions to satisfy the KKT conditions.

Proposition 3: The KKT conditions for the problem (CO-L2-Aut) and (CO-L1-Aut) are satisfied if there exist

scalars  $\gamma_i^l \in [-1,1], i \notin \mathcal{K}, l = 1, \dots, n$  such that

$$\sum_{i \notin \mathscr{K}} \gamma_i^l x_i / \sqrt{n} = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_2^l x_i, \quad \forall l = 1, \dots, n$$
 (6)

and

$$\sum_{i \notin \mathscr{K}} \gamma_i^l x_i = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_1^l x_i, \quad \forall l = 1, \dots, n,$$
(7)

respectively. Here,  $\partial \|\bar{d}_i\|_{\circ}^l$  is the *l*-th element of the subgradient.

Because analyzing the conditions (6) and (7) directly is cumbersome, we investigate the equivalent condition provided in the lemma below, derived using Farkas' lemma [28] and the duality of linear programs.

*Lemma 3:* Given a matrix  $\mathbf{F} \in \mathbb{R}^{n \times m}$  and the vector  $g \in \mathbb{R}^n$ , the following statements are equivalent:

- i) There exists a vector  $w \in \mathbb{R}^m$  with  $||w||_{\infty} \leq 1$  satisfying  $\mathbf{F}w = g$ .
- ii) For every  $z \in \mathbb{R}^n$  with  $||z||_2 = 1$ , it holds that  $f(z) := z^T g + ||z^T \mathbf{F}||_1 \ge 0$ .

It is important to notice that the conditions (6) and (7) amount to finding a vector for the set of equations in the form of  $\mathbf{F}w = g$  where w is restricted as  $||w||_{\infty} \leq 1$ . Given a coordinate l, the matrix  $\mathbf{F} \in \mathbb{R}^{n \times (T - |\mathcal{K}|)}$  associated with the conditions (6) and (7) is a matrix with columns  $\frac{x_i}{\sqrt{n}}$  and  $x_i$ , and the vector  $g \in \mathbb{R}^n$  is  $\sum_{i \in \mathcal{K}} \partial ||\bar{d}_i||_2^l x_i$  and  $\sum_{i \in \mathcal{K}} \partial ||\bar{d}_i||_1^l x_i$ , respectively. Moreover, the vector  $w \in \mathbb{R}^{T - |\mathcal{K}|}$  has the elements  $\gamma_i^l, i \notin \mathcal{K}$  for both conditions. Hence, we study the second statement in Lemma 3. We use the union bound to study the satisfaction of this condition. However, there are infinitely many points inside the  $\ell_2$  unit ball  $\mathbb{B}_2(1)$ . In order to show that the function  $f(z) = z^T g + ||z^T \mathbf{F}||_1$  is non-negative at every point inside the  $\ell_2$  unit ball, we employ the discretization technique that uses a finite set of points. The set of such points is called the cover of the unit ball.

Definition 5 (Covering Number [25]): Let  $(\mathbb{T}, \rho)$  be a compact metric space with a set  $\mathbb{T}$  and a norm operator  $\rho$ .  $\varepsilon$ -cover of the set  $\mathbb{T}$  with respect to the norm  $\rho$  is a set  $\{\theta^1, \theta^2, \ldots, \theta^N\} \subset \mathbb{T}$  such that for each  $\theta \in \mathbb{T}$ , there exists some  $i \in \{1, \ldots, N\}$  such that  $\rho(\theta, \theta^i) \leq \varepsilon$ . The  $\varepsilon$ -covering number  $\mathscr{N}(\varepsilon, \mathbb{T}, \rho)$  is the cardinality of the smallest  $\varepsilon$ -cover.

Given a  $\varepsilon > 0$ , the logarithm of the covering number of the unit ball or the metric entropy of the unit ball can be upper bounded using the volumetric arguments of the balls. Indeed, the number of  $\varepsilon$  balls exceeding  $\exp\{n\log(1+2/\varepsilon)\}$ is sufficient to cover the unit ball with balls of radius  $\varepsilon$ .

Lemma 4 (Covering Number of the Unit Ball [25]): Given an *n*-dimensional unit ball  $\mathbb{B}(1)$  with the norm  $\|\cdot\|$ ,

$$\mathbb{B}(1) = \{ x \in \mathbb{R}^n : ||x|| \le 1 \},\$$

the logarithm of the  $\varepsilon$ -covering number, i.e., the metric entropy of the unit ball, can be upper bounded by

$$\log \mathscr{N}(\varepsilon, \mathbb{B}(1), \|\cdot\|) \le n \log \left(1 + \frac{2}{\varepsilon}\right).$$

We show that the function f(z) can be lower bounded by some positive number  $\theta > 0$  at every point in the  $\varepsilon$ -cover of the unit circle with high probability, and that the function value inside the  $\varepsilon$ -ball does not change more than this positive number  $\theta$  with high probability. Thus, f(z) must be nonnegative at every point of the unit circle with high probability. Utilizing this idea, the next theorem shows that the required number of samples for the exact recovery grows with  $n^2$  and  $(1-p)^{-2}$  for the general systems of order *n*.

Theorem 2: Consider an autonomous system of order *n* under a probabilistic attack model with frequency *p*. Suppose that Assumptions 1 and 2 hold. Then, for all  $\delta \in (0,1]$ , if the time horizon satisfies  $T \ge \Theta(T_{\text{sample}})$ , where  $T_{\text{sample}}$  is defined as

 $nR\left[n\log(nR) + \log\left(\frac{1}{\delta}\right)\right],$ 

and

$$R := \max\left\{\frac{\log(1/c)}{nc^4 p(1-p)\log(1/\rho)}, \frac{\log^2(1/c)}{c^{10}(1-p)^2(1-\rho)^3\log^2(1/\rho)}, \frac{1}{np(1-p)}\right\},$$

with  $\rho$  denoting the largest magnitude of the eigenvalues of  $\bar{A}$ , then  $\bar{A}$  is a solution to the convex optimization (CO-L2-Aut) with probability at least  $1 - \delta$ .

An implication of the above theorem is that even when p is large (e.g., p > 0.5) corresponding to the system being under attack frequently, exact recovery of the system dynamics is still possible as long as the time horizon is above the threshold. Similar results can be obtained if one prefers to use problem (CO-L1-Aut) to recover the system matrix  $\overline{A}$ .

Theorem 3: Under the same assumptions as in Theorem 2, if the time horizon T satisfies  $T \ge \Theta(T_{\text{sample}})$ , where  $T_{\text{sample}}$  is defined as

$$R\left[n\log(nR) + \log\left(\frac{1}{\delta}\right)\right],$$

and *R* is defined in Theorem 2, then  $\overline{A}$  is a solution to the convex optimization (CO-L1-Aut) with probability at least  $1 - \delta$ .

The proof of Theorem 3 is highly similar to that of Theorem 2 and therefore, it is omitted. Because the conditions (6) and (7) differ by a factor of  $\sqrt{n}$ , the sample complexity results in those theorems differ by a factor of n.

The required amount of data increases with the value  $(1 - p)^{-2}$  and the order of the system *n*. Hence, as *p* and *n* increase, the number of samples for exact recovery with high probability grows. The results on sample complexity are intuitive: as the probability of having an attack increases, a larger time horizon is required for exact recovery. We note that the dependence on  $p^{-1}(1-p)^{-1}$  is an artifact of the high probability bound. More specifically, this dependence guarantees that the number of attacks is bounded by  $\Theta(pT)$  with high probability. In addition, if the system is at the verge of instability with eigenvalues close to the unit circle, the sample complexity increases significantly. Even in the case when the probability *p* is close to 1, resulting in significantly more corrupt data than clean data, this result guarantees asymptotic exact recovery as long as there are a sufficient number of clean samples.

8

Last but not at least, due to the logarithmic probability bound and the Borel-Cantelli lemma, Theorems 2 and 3 imply almost sure asymptotic convergence as a corollary. Almost sure convergence of random variables implies the convergence in probability and convergence in distribution for a sequence of random variables. Almost sure convergence of random variables is defined as below for completeness.

Definition 6 (Almost Sure Convergence): A sequence of random variables  $X_1, X_2, X_3, \ldots$  converges to X almost surely if

$$\mathbb{P}\left(\lim_{n\to\infty}X_n=X\right)\to 1.$$

The following corollary states that the sequence of estimators over time converges to the true system matrices almost surely.

Corollary 1: Under the same assumptions as in Theorem (2),  $\overline{A}$  is almost surely a solution of convex formulations (CO-L2-Aut) and (CO-L1-Aut) when T goes to infinity.

# V. SYSTEMS WITH INPUT SEQUENCE

It is desirable to understand the role of an input sequence in exact recovery because the majority of dynamical systems are controlled by an external input. Since the input sequence is generated by a controller, one can design it in such a way that it accelerates the exact recovery. In the non-autonomous case, the system dynamics is given as  $x_{i+1} = \bar{A}x_i + \bar{B}u_i + \bar{d}_i, i = 0, ..., T - 1$ , where  $\bar{A} \in \mathbb{R}^{n \times n}$  and  $\bar{B} \in \mathbb{R}^{n \times m}$ . Similar to the autonomous case, the true system matrices  $\bar{A}$  and  $\bar{B}$  are not known and the goal is to obtain these matrices using the state trajectories and the sequence of inputs. Unlike the disturbance vectors  $\bar{d}_i, i \in \{0, ..., T - 1\}$ , the sequence of system states  $x_i, i \in \{0, ..., T\}$ , and the sequence of the inputs  $u_i, i \in \{0, ..., T - 1\}$  are known. We will investigate the estimators (CO-L2) and (CO-L1) defined earlier.

We choose the input vectors  $u_i$  to be Gaussian given  $\mathscr{F}_i$ . This allows us to obtain a high-probability bound for the exact recovery of the matrices  $\overline{A}$  and  $\overline{B}$ . A random input sequence is commonly used in system identification and online learning because it enables the exploration of the system to learn the system dynamics faster. The Gaussian input assumption may seem restrictive. Nevertheless, it is satisfied when  $u_i$ is designed in the linear feedback form as  $u_i = Kx_i + \omega$ . Conditioning on  $\mathscr{F}_i$ , if the input is excited with Gaussian noise  $\omega$ , the input vector  $u_i$  is also Gaussian. Therefore, the most common input sequence used in optimal control satisfies this assumption. Note that the closed loop system could be written as  $x_{i+1} = (\overline{A} + \overline{B}K)x_i + \overline{B}\omega + \overline{d_i}$ . Thus, the problem is equivalent to estimating the matrices  $(\overline{A} + \overline{B}K)$  and  $\overline{B}$  when the linear feedback control is used.

The KKT conditions for the exact recovery that are both necessary and sufficient can be restated as

$$\exists \gamma_i \in \partial \|0\|_{\circ} \quad \forall \ i \notin \mathscr{K} \text{ s.t. } \sum_{i \notin \mathscr{K}} x_i \otimes \gamma_i = \sum_{i \in \mathscr{K}} x_i \otimes \partial \|\bar{d}_i\|_{\circ}$$

and

$$\exists \mu_i \in \partial \|0\|_{\circ} \quad \forall \ i \notin \mathscr{K} \ \text{s.t.} \ \sum_{i \notin \mathscr{K}} u_i \otimes \mu_i = \sum_{i \in \mathscr{K}} u_i \otimes \partial \|\bar{d}_i\|_{\circ}$$

The first set of conditions corresponds to the KKT conditions for the system states while the second set is for the KKT conditions for the input sequence. Similar to Proposition 3, the sufficient conditions can be tightened so that the equations become coordinate-wise separable.

*Proposition 4:* The KKT conditions for problem (CO-L2) are satisfied if there exist scalars  $\gamma_i^l, \mu_i^l \in [-1,1]$  for all  $i \notin \mathcal{K}, l \in \{1, ..., n\}$  such that

$$\sum_{\substack{\notin \mathscr{K}}} \gamma_i^l x_i / \sqrt{n} = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_2^l x_i, \quad \forall l = 1, \dots, n,$$
(8)

and

$$\sum_{i \notin \mathscr{K}} \mu_i^l u_i / \sqrt{n} = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_2^l u_i, \quad \forall l = 1, \dots, n,$$
(9)

where  $\partial \|\bar{d}_i\|_2^l$  denotes the *l*-th element of the subgradient.

The proof of Proposition 4 is omitted because it relies on the same technique as in Proposition 3. As in the case of autonomous systems, two sets of equations that guarantee the satisfaction of the KKT conditions can be written for problem (CO-L1) by omitting the factor  $\sqrt{n}$ . To establish the exact recovery guarantees, we require the following controllability assumption.

Assumption 3: The ground truth  $(\bar{A}, \bar{B})$  satisfies

rank {  $\begin{bmatrix} \bar{B} & \bar{A}\bar{B} & \cdots & \bar{A}^{n-1}\bar{B} \end{bmatrix}$  } = n.

Intuitively, the controllability of a non-autonomous system denotes the ability to move a system around in its entire state space using the admissible manipulations, namely, the input sequence  $\{u_t\}_{t=0}^{T-1}$ . Controllability is an important property of a control system and plays a crucial role in many control problems, such as stabilization of unstable systems by feedback. Under the above assumption, we implement the non-asymptotic analysis of the general non-autonomous system in a similar fashion to Theorem 2 using the covering arguments and Farkas' lemma.

Theorem 4: Consider an autonomous system of order *n* under a probabilistic attack model with frequency *p*. Suppose that Assumptions 1, 3 and the first three conditions in Assumption 2 hold. Assume also that the input vectors  $u_i | \mathscr{F}_i$  are selected to be independent from the attack vectors and obey the Gaussian distribution  $\mathscr{N}(0, \frac{\xi^2}{m}I_m)$ . For all  $\delta \in (0, 1]$ , let

$$T_{\text{sample}}^1 := nR_1 \left[ n \log(nR_1) + \log\left(\frac{1}{\delta}\right) \right]$$

and

$$T_{\text{sample}}^2 := nR_2 \left[ m \log(nR_2) + \log\left(\frac{1}{\delta}\right) \right],$$

where

$$R_{1} := \max\left\{\frac{\log(\kappa/c)}{nc^{4}\log(1/\rho)}, \frac{p\kappa^{2}}{c^{10}(1-p)^{2}(1-\rho)^{2}} \\ \frac{p\kappa^{2}\log^{2}(\kappa/c)}{c^{10}(1-\rho)^{2}\log^{2}(1/\rho)}, \frac{1}{np}\right\}$$
$$R_{2} := \max\left\{\frac{1}{np}, \frac{p}{(1-p)^{2}}, \frac{m}{n}\right\}.$$

Here, constants  $c \in (0,1]$  and  $\kappa \ge (1-\rho)^{-1}$  depend on *m*, *n*,  $\sigma$ ,  $\xi$  and  $\bar{B}$ . If the time horizon satisfies the inequality  $T \ge \Theta[\max\{T_{\text{sample}}^1, T_{\text{sample}}^2\}]$ , then  $(\bar{A}, \bar{B})$  is a solution to (CO-L2) with probability at least  $1-\delta$ .

We have obtained a high probability bound for the exact recovery of the system matrices A and  $\overline{B}$ . The first term in the sample complexity corresponds to the satisfaction of the KKT conditions for the state measurements  $\{x_i\}_{i=0}^T$ , whereas the second term corresponds to the satisfaction of the KKT conditions for the input sequence  $\{u_i\}_{i=0}^{T-1}$ . Similar to the case of autonomous systems, the sample complexity increases as the probability of disturbances increases. Because there is a logarithmic dependence on the satisfaction of the probability bound, Theorem 4 and the application of the Borel-Cantelli lemma imply almost sure asymptotic convergence to the correct matrices  $\overline{A}$  and  $\overline{B}$ . The sample complexity  $T_{\text{sample}}^2$  is needed to satisfy the KKT conditions associated with on the input sequence. Compared with the previous theorems for the autonomous case, we require a sample complexity that scales with  $p/(1-p)^2$  and terms depending on the spectral norm of  $\overline{A}$ . The introduction of the input sequence removes the requirement on the variance of the attack vectors. In addition, the dependence of the sample complexity on p is improved from  $1/(1-p)^2$  to  $p/(1-p)^2$ . Moreover, the dependence on the spectrum of  $\overline{A}$  is reduced from  $1/[(1-\rho)^3 \log^2(1/\rho)]$  to  $1/[(1-\rho)^2 \log^2(1/\rho)]$ . Finally, we mention that the dependence on 1/(np) is also to guarantee that the number of attacks is bounded by  $\Theta(pT)$  with high probability.

The following theorem studies problem (CO-L1).

Theorem 5: Under the assumptions of Theorem 4, for all  $\delta \in (0,1]$ , let  $T_{\text{sample}}^1$  and  $T_{\text{sample}}^2$  be defined as

$$R_1\left[n\log(nR_1) + \log\left(\frac{1}{\delta}\right)\right]$$
 and  $R_2\left[m\log(nR_2) + \log\left(\frac{1}{\delta}\right)\right]$ ,

where  $R_1$  and  $R_2$  are given in Theorem 4. If the time horizon satisfies the inequality  $T \ge \Theta[\max\{T_{\text{sample}}^1, T_{\text{sample}}^2\}]$ , then  $(\bar{A}, \bar{B})$  is a solution to (CO-L1) with probability at least  $1 - \delta$ .

As expected, even if more than half of the data are corrupted, that is p > 1/2, the exact recovery is still attainable with high probability. We note that when the input sequence  $u_i = Kx_i$  is used to control the system, this input sequence satisfies the assumptions in the above theorems if  $x_i$  are sub-Gaussian. The closed-loop system with the matrix  $(\bar{A} + \bar{B}K)$ results in a second solution  $\hat{A} = \bar{A} + \bar{B}K$  and  $\hat{B} = 0$ . Nevertheless, the ground-truth system matrix pair  $(\bar{A}, \bar{B})$  is also a solution to our estimators. This phenomenon occurs due to the existence of multiple optimal solutions and it could be avoided if the input is excited with a small noise in the form of  $u_i = Kx_i + \omega$ . Moreover, if all the input vectors  $u_i$  are set to zero, it is not possible to uniquely recover the system matrix  $\overline{B}$ . Nevertheless, because the input sequence  $\{u_i\}_{i=0}^{T-1}$ is zero, the KKT conditions are trivially satisfied. Therefore, the estimators have multiple optimum solutions where  $\bar{B}$  and 0 matrices are possible solutions among all optimum solutions.

### VI. NUMERICAL EXPERIMENT

We conduct a numerical experiment inspired by biomedical applications to demonstrate the results of this paper. We consider a compartmental model of blood sugar and insulin dynamics in the human body, as described in [29]. Accurately estimating the parameters of the dynamics is crucial when regulating the blood sugar level through the injection of a bolus of insulin into the system. Due to the complex structure of the human body, the dynamics vary among individuals. We consider a linear system based on Hovarka's model as follows [30]:

$$\begin{aligned} \dot{x}_1 &= -k_{a1}x_1 - k_{b1}I + d_1, \\ \dot{x}_2 &= -k_{a2}x_1 - k_{b2}I + d_2, \\ \dot{x}_3 &= -k_{a3}x_1 - k_{b3}I + d_3, \\ \dot{S}_1 &= -S_1/t_{max,I} + d_4, \\ \dot{S}_2 &= S_1/t_{max,I} - S_2/t_{max,I} + d_5, \\ \dot{I} &= S_2/(t_{max,I}V_I) - k_eI + d_6, \end{aligned}$$

where given a time-dependent variable z(t),  $\dot{z}(t)$  represents its derivative with respect to time t. The states  $x_1, x_2, x_3$  represent the influence of insulin on the system of the body.  $S_1$  and  $S_2$  represent the absorption rate of insulin in the, directly and indirectly, accessible compartment models, respectively. Lastly, the state I represents the blood sugar level in the body. The disturbance  $d_4$  corresponds to the bolus injection into the body, while the remaining disturbance vectors model sudden changes in the body due to diseases such as diabetes. Although the injected insulin amount could be known, the exact amount of insulin and its timing reaching the effective body parts are unknown. Hence, the  $d_i$  values are treated as unknown. Even though the disturbance in this application is not a malicious attack, it exhibits similar characteristics for identification purposes: the arrival time of the bolus is unknown, and once it arrives, it has a large magnitude.

In this experiment, we discretize the continuous-time system to obtain an LTI system using  $\Delta_t = 0.5$ . The resulting matrix  $\bar{A}$  is stable. Our objective is to estimate the parameters  $(k_{ai}, k_{bi}, t_{\max,I}, V_I, k_e)$ , where the true values are obtained from Table 1 in [31]. We model the attack vectors given the historical data as zero-mean Gaussian random vectors with an identity covariance matrix with variance 10. Thus, the attack vectors are conditionally independent, although they are dependent. We run our model with the probability of an attack being p = 0.2, p = 0.4, and p = 0.6. We report the estimation error  $|\hat{A} - \bar{A}|_F$  for the least-squares estimator, problem (CO-L2), and problem (CO-L1).

Figure 2 suggests that our proposed estimators attain exact recovery while the least-squares estimator fails to do so. As the probability of having an attack p increases, the number of required time periods for exact recovery grows proportionally to  $p/(1-p)^2$ . Note that there are more corrupted data than clean data in the case of p = 0.6. Additionally, because there is no sparsity assumption on the attack vectors, (CO-L2) performs slightly better than (CO-L1).

We compare the performance of (CO-L2) and (CO-L1) by running a similar experiment with and without sparse



Fig. 2. Estimation errors for Least-Squares, (CO-L2), and (CO-L1) with attack probability of p = 0.2, 0.4, 0.6 (left-to-right).



Fig. 3. Estimation errors for Least-Squares, (CO-L2), and (CO-L1) with attack probability p = 0.6 not Sparse *d* (top) Sparse *d* (bottom).

disturbances. When the disturbances are sparse,  $d_1$ ,  $d_2$ ,  $d_3$ ,  $d_5$  are set to zero while  $d_4$  and  $d_6$  have the same Gaussian distribution as before. Figure 3 shows that the two methods perform similarly when the attack vectors are also sparse.

# VII. DISCUSSION AND CONCLUSION

We investigated the problem of learning LTI systems under adversarial attacks by studying two lasso-type estimators. We considered both deterministic and probabilistic attack models regarding the time occurrence of the attack and developed strong conditions for the exact recovery of the system dynamics. When the attacks occur deterministically every  $\Delta$  period, exact recovery is possible after  $n + \Delta$  time steps. Moreover, if the system is attacked at each time instance with probability p, the system matrices are recovered with high probability when *T* is on the order of  $\Theta((1-p)^{-2})$  and a polynomial in the dimension of the problem. Similar results were obtained when the system is controlled by an input sequence. These findings were supported by a numerical experiment in biology that to validate the non-asymptotic analytic results. This work provides the first set of mathematical guarantees for the robust non-asymptotic analysis of dynamic systems.

Since our estimators have non-smooth objective functions, closed-form solutions to the optimization problem are not obtainable. We did not provide any specific numerical algorithm to solve the provided estimation problems. However, both (CO-L2) and (CO-L1) are convex optimization problems, allowing the use of the subgradient descent algorithm to obtain these estimators. It is a well-established result that the subgradient algorithm has a convergence rate on the order of  $\frac{1}{\sqrt{k}}$ , where k is the iteration update number. Although the algorithm is considered fast, one possible extension of this work would be to design an algorithm to predict and update  $(\hat{A}_{t+1}, \hat{B}_{t+1})$  using the latest estimation  $(\hat{A}_t, \hat{B}_t)$  and the new data  $(x_{t+1}, u_t)$ , instead of solving the problem from scratch at each time period. Initial experiments hinted that a single subgradient update at each iteration using the new information,  $(x_{t+1}, u_t)$ , asymptotically converges to the true system matrices. We leave the analysis of this algorithm and online control of dynamic systems under adversaries as future work.

#### **APPENDIX**

# A. Proofs for Results in Main Part

1) Proof of Proposition 1: The proof of Proposition 1 is established based on Lemma 2 defined in the paper. Let  $i_1, i_2, ...$  be the set of attack times over time horizon *T*. Therefore,  $\mathscr{K} = \{i_1, i_2, ...\}$ . Due to  $\Delta$ -spaced attack model, the first attack time must be smaller than  $\Delta$ , i.e.,  $i_1 \leq \Delta$ . Since  $x_0 = 0$ , we have  $x_t = 0$  for  $t = 0, 1, ..., i_1$ . Define  $\mathbb{N}$  as the set of natural numbers. We can utilize Lemma 2 to show that  $\overline{A}$ is the unique solution. Using these facts, we can decompose the sum of the magnitudes of the states at non-attack times as

$$\sum_{i \notin \mathscr{K}} |x_i| = \sum_{i \notin \mathscr{K}, i > i_1} |x_i| = \sum_{i \in \mathscr{K}'} |x_i| + \sum_{i \in \mathscr{K}''} |x_i|,$$

where  $\mathscr{K}^{c+} = \mathbb{N} \setminus (\mathscr{K} \cup \{0, 1, \dots, i_1 - 1\}), \ \mathscr{K}' = \mathscr{K}^{c+} \setminus \mathscr{K}'',$ and  $\mathscr{K}'' = \{i_2 - 1, i_3 - 1, \dots\}$ . The second term on the righthand side is the sum of magnitudes at the time step just before the attack while the first term covers the rest of the magnitudes of the states. In addition, the magnitudes of the states at attack times can be written as

$$\sum_{i \in \mathscr{K}} |x_i| = \sum_{i \in \mathscr{K}, i \ge i_2} |x_i| = \sum_{i \in \mathscr{K}''} |\bar{A}x_i| = \sum_{i \in \mathscr{K}''} |\bar{A}||x_i|$$

The second equality follows from the fact that  $x_{i_k} = \bar{A}x_{i_k-1}$  due to lack of attack. We compare the sum of the magnitudes of the states at attack times for the non-attack times to check if the condition in Lemma 2 holds:

$$\sum_{i \notin \mathscr{K}} |x_i| - \sum_{i \in \mathscr{K}} |x_i| = \sum_{i \in \mathscr{K}'} |x_i| + \sum_{i \in \mathscr{K}''} |x_i| - \sum_{i \in \mathscr{K}''} |\bar{A}| |x_i|$$
$$= \sum_{i \in \mathscr{K}'} |x_i| + (1 - |\bar{A}|) \sum_{i \in \mathscr{K}''} |x_i| > 0. \quad (10)$$

Note that the term  $\sum_{i \notin \mathscr{K}} |x_i|$  becomes positive at time period  $i_1 + 1$  while  $\sum_{i \in \mathscr{K}} |x_i|$  is positive first time at time step  $i_2$ . Consequently, the strict inequality for (10) holds for every time step after  $i_1$  because  $(1 - |\bar{A}|) > 0$  by assumption. As a result, we have a unique and exact recovery for every time period  $T \ge \Delta + 1 \ge i_1 + 1$ .

2) Proof of Proposition 2: By using (3), the necessary and sufficient condition for this problem is

$$0 \in \sum_{i \notin \mathscr{K}} x_i \otimes \partial \| (\bar{A} - A) x_i \|_2 + \sum_{i \in \mathscr{K}} x_i \otimes \partial \| (\bar{A} - A) x_i + \bar{d}_i \|_2.$$

Then,  $\overline{A}$  is a solution to the problem if and only if

$$0 \in \sum_{i \notin \mathscr{K}} x_i \otimes \partial \|0\|_2 + \sum_{i \in \mathscr{K}} x_i \otimes \partial \|\bar{d}_i\|_2.$$
(11)

Let  $i_1$  be the time stamp of the first attack time. Then, we have  $i_1 \in \{1, ..., \Delta\}$ . The set of attack times is  $\mathcal{H} = \{i_1, i_1 + \Delta, i_1 + 2\Delta, i_1 + 3\Delta, ...\}$ . Since  $x_0 = 0$ , we have  $x_t = 0$  whenever  $t = 0, 1, ..., i_1$  and  $x_{i_1+1} = \overline{d}_{i_1}$ . Let  $T = \Delta + i_1$ , i.e., the time step at which a cycle of disturbance is completed. In this case, the sufficient condition (3) can be written as

$$0 \in \sum_{t=1}^{\Delta-1} x_{i_1+t} \otimes \partial \|0\|_2 + x_{i_1+\Delta} \otimes \partial \|\bar{d}_{i_1+\Delta}\|_2$$
$$= \sum_{t=0}^{\Delta-2} \bar{A}^t \bar{d}_{i_1} \otimes \partial \|0\|_2 + \bar{A}^{\Delta-1} \bar{d}_{i_1} \otimes \frac{\bar{d}_{i_1+\Delta}}{\|\bar{d}_{i_1+\Delta}\|}.$$

The matrix 0 may belong to the right-hand side term for arbitrary  $\bar{d}_{i_1+\Delta}$  if  $\bar{d}_{i_1+\Delta} \in span\{\bar{d}_{i_1}, \bar{A}\bar{d}_{i_1}, \dots, \bar{A}^{\Delta-2}\bar{d}_{i_1}\}$ . This is satisfied by the assumption in the proposition statement.

However, this is not sufficient to ensure that KKT condition (3) holds. The reason is that  $\partial ||0||_2 = \{x \in \mathbb{R}^n : ||x||_2 \le 1\}$ . The vectors chosen for  $\partial ||0||_2$  have a bounded norm. Therefore, we need a condition that bounds the norm of the columns of  $\bar{A}^{\Delta-1}\bar{d}_{i_1} \otimes \frac{\bar{d}_{i_1+\Delta}}{||\bar{d}_{i_1+\Delta}||_2}$ , so it can be expressed as a linear combination of the vectors  $\{\bar{d}_{i_1}, \bar{A}\bar{d}_{i_1}, \dots, \bar{A}^{\Delta-2}\bar{d}_{i_1}\}$ . Let  $(\lambda_j, v_j)$  be eigenvalue-eigenvector pairs for the matrix  $\bar{A}^T$ . Let  $e_1, \dots, e_{\Delta-1} \in \partial ||0||_2$ . Then, the KKT condition can be written as follows after dropping the sub-index  $i_1$ :

$$0 \in e_1 \bar{d}^T + e_2 \bar{d}^T \bar{A}^T + \dots + e_{\Delta - 1} \bar{d}^T (\bar{A}^T)^{\Delta - 2} + f \bar{d}^T (\bar{A}^T)^{\Delta - 1},$$

where  $f = \frac{\bar{d}_{i_1+\Delta}}{\|\bar{d}_{i_1+\Delta}\|_2}$  and  $\|f\|_2 = 1$ . If we multiply the equation above by the eigenvector  $v_i$  of  $\bar{A}^T$ , we obtain

$$D \in e_{1}\bar{d}^{T}v_{j} + \dots + e_{\Delta-1}\bar{d}^{T}(\bar{A}^{T})^{\Delta-2}v_{j} + f\bar{d}^{T}(\bar{A}^{T})^{\Delta-1}v_{j}$$
  
$$\in (e_{1} + \lambda_{j}e_{2} + \dots + \lambda_{j}^{\Delta-2}e_{\Delta-1} + \lambda_{j}^{\Delta-1}f)\bar{d}^{T}v_{j}.$$

Note that because  $\overline{A}$  is diagonalizable, we only need to satisfy this condition along the direction of each eigenvector, since all eigenvectors span the whole space. Therefore, the KKT condition holds if

$$0 \in e_1 + \lambda_j e_2 + \dots + \lambda_j^{\Delta-2} e_{\Delta-1} + \lambda_j^{\Delta-1} f, \quad \forall j = 1, \dots, n.$$

There are  $(\Delta - 1)n$  free variables and  $n^2$  equations. One can use the substitution to eliminate  $n^2$  variables, which leads to

$$\sum_{1+\dots+k_n=\Delta-n} \lambda(k_1,\dots,k_n)f = \sum_{t=0}^{\Delta-n-2} \sum_{k_1+\dots+k_n=t} \lambda(k_1,\dots,k_n)e_{t+n+1}.$$

Taking the norm of both sides and using the triangle inequality yields that

$$\left| \sum_{\substack{k_1 + \dots + k_n = \Delta - n}} \lambda(k_1, \dots, k_n) \right| \|f\|_2$$
  
$$\leq \sum_{t=0}^{\Delta - n - 1} \left| \sum_{\substack{k_1 + \dots + k_n = t}} \lambda(k_1, \dots, k_n) \right| \|e_{t+n+1}\|_2$$

Using the fact that  $||e_j||_2 = 1$  for all *j* and  $||f||_2 = 1$ , we obtain

$$\left|\sum_{k_1+\dots+k_n=\Delta-n}\lambda(k_1,\dots,k_n)\right|\leq \sum_{t=0}^{\Delta-n-1}\left|\sum_{k_1+\dots+k_n=t}\lambda(k_1,\dots,k_n)\right|.$$

This completes the proof for the proposition.

k

*3) Proof of Proposition 3:* The KKT condition for the exact recovery that is the necessary and sufficient condition can be restated as

$$\exists \gamma_i \in \partial \|0\|_{\circ}, i \notin \mathscr{K} \text{ s.t. } \sum_{i \notin \mathscr{K}} x_i \otimes \gamma_i = \sum_{i \in \mathscr{K}} x_i \otimes \partial \|\bar{d}_i\|_{\circ}.$$
(12)

For the problem (CO-L2-Aut) with  $\circ = 2$ . the condition (12) becomes

$$\exists \gamma_i \in \partial \|0\|_2, i \notin \mathscr{K} \text{ s.t. } \sum_{i \notin \mathscr{K}} x_i \otimes \gamma_i = \sum_{i \in \mathscr{K}} x_i \otimes \partial \|\bar{d}_i\|_2.$$

Since  $\frac{1}{\sqrt{n}}\partial \|0\|_1 = \frac{1}{\sqrt{n}}\mathbb{B}_{\infty}(1) \subseteq \mathbb{B}_2(1) = \partial \|0\|_2$ , we can rewrite it as

$$\exists \gamma_i \in \partial \|0\|_1, i \notin \mathscr{K} \text{ s.t. } \sum_{i \notin \mathscr{K}} \frac{x_i}{\sqrt{n}} \otimes \gamma_i = \sum_{i \in \mathscr{K}} x_i \otimes \partial \|\bar{d}_i\|_2.$$

We can check the condition at each coordinate because the set  $\mathbb{B}_{\infty}(1)$  is coordinate wise separable. Thus, the condition becomes that KKT condition holds for (CO-L2-Aut) if there exist scalars  $\gamma_i^l \in [-1,1], i \notin \mathcal{K}, l = 1, ..., n$  such that

$$\sum_{i \notin \mathscr{K}} \gamma_i^l x_i / \sqrt{n} = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_2^l x_i, \quad \forall l = 1, \dots, n,$$

where  $\partial \|\bar{d}_i\|_{\circ}^l$  is the *l*-th element of the subgradient. Similar algebraic manipulation can be done for (CO-L1-Aut) except for the transforming subdifferential of the  $\ell_2$  norm to subdifferential of the  $\ell_1$  norm to obtain the second part of the result.

4) Proof of Lemma 3: The condition "Given a matrix  $\mathbf{F} \in \mathbb{R}^{n \times m}$  and the vector  $g \in \mathbb{R}^n$ , there exists a vector  $w \in \mathbb{R}^m$  with  $||w||_{\infty} \leq 1$  satisfying  $\mathbf{F}w = g$ ." is equivalent to the feasibility of the linear programming (LP) below with objective function equal to 0:

$$\max_{w \in \mathbb{R}^m} \quad 0$$
  
s.t.  $\mathbf{F}w = g,$   
 $\|w\|_{\infty} \le 1.$ 

Due to the strong duality, the dual problem of the LP above must have the optimum objective value equal to 0. The dual problem can be formulated as

$$\min_{\substack{y \in \mathbb{R}^m, z \in \mathbb{R}^n}} z^T g + \|y^T\|_1$$
  
s.t.  $z^T \mathbf{F} + y^T = 0$ ,

or equivalently,

$$\min_{z\in\mathbb{R}^n} \quad f(z):=z^Tg+\|z^T\mathbf{F}\|_1.$$

Thus, for any  $z \in \mathbb{R}^n$ , f(z) must be nonnegative, i.e.,  $f(z) \ge 0$ . Because f(cz) = cf(z) for all c > 0, the condition  $f(z) \ge 0$ for all  $z \in \mathbb{R}^n$  is satisfied if  $f(z) \ge 0$  for all  $z \in \mathbb{R}^n$  such that  $||z||_2 = 1$ . This completes the proof.

5) Proof of Theorem 2: Due to the system dynamics and given  $x_0 = 0$ ,  $x_i$  can be expressed as

$$x_i = \sum_{k \in \mathcal{K}} \bar{A}^{(i-k-1)_+} \bar{d}_k$$

where  $A^{(i)_+}$  is defined as

$$A^{(i)_{+}} := \begin{cases} 0, & \text{if } i < 0\\ I, & \text{if } i = 0 \\ A^{i}, & \text{if } i > 0 \end{cases}$$

By Lemma 3, given a coordinate  $l \in \{1, ..., n\}$ , the optimality condition for the recovery of  $\overline{A}$  is equivalent to

$$f(z) := z^T g + \|z^T \mathbf{F}\|_1 \ge 0, \quad \forall z \in \mathbb{S}_2(1),$$
(13)

where the unit sphere  $\mathbb{S}_2(1)$  is  $\{z \in \mathbb{R}^n : ||z||_2 = 1\}$ , the matrix  $\mathbf{F} \in \mathbb{R}^{n \times (T-|\mathcal{K}|)}$  has the columns

$$\mathbf{F}^{i} := \sum_{k \in \mathscr{K}} \frac{\bar{A}^{(i-k-1)_{+}} \bar{d}_{k}}{\sqrt{n}}, \quad \forall i \notin \mathscr{K},$$

and the vector  $g \in \mathbb{R}^n$  is

$$g := \sum_{i \in \mathscr{K}} \sum_{k \in \mathscr{K}} \bar{A}^{(i-k-1)_+} \bar{d}_k \cdot \bar{f}_i^l$$

We prove that condition (13) holds with high probability in two steps.

Step 1: We first prove that condition (13) holds with high probability for a fixed  $z \in S_2(1)$ .

a) Step 1-1: We first analyze the term  $||z^T \mathbf{F}||_1$ , namely,

$$\mathbb{E} \| z^T \mathbf{F} \|_1 = \frac{1}{\sqrt{n}} \sum_{i \notin \mathscr{K}} \mathbb{E} \left| \sum_{k \in \mathscr{K}} z^T \bar{A}^{(i-k-1)_+} \bar{d}_k \right|.$$
(14)

We construct the index set

Let

$$\mathscr{I}_1 := \{ i \mid i \notin \mathscr{K}, \ i-1 \in \mathscr{K} \}.$$

$$S := \left[ \log_{\rho} \Theta \left[ \frac{c^5}{\log(|\mathscr{I}_1|/\delta)} \right] \right]$$
$$= \Theta \left[ \frac{\log\log(|\mathscr{I}_1|/\delta) + \log(1/c)}{\log(1/\rho)} \right]$$

where  $\lceil x \rceil$  is the minimal integer that is not smaller than x and  $\delta \in (0, 1)$  is the specified probability. We construct a subset of  $\mathscr{I}_1$  in the following way:

$$\mathscr{I} := \{i_1, \ldots, i_I \mid i_j \in \mathscr{I}_1, \ i_j - i_{j-1} \ge S, \ \forall j\}$$

It is straightforward to construct  $\mathscr{I}$  such that

$$I = |\mathscr{I}| \ge \frac{1}{S}|\mathscr{I}_1|.$$

In addition, due to the probabilistic attack model, it holds with probability at least  $1 - \exp[-\Theta[p(1-p)T]]$  that

$$|\mathscr{I}_1| \ge \frac{p(1-p)T}{2}$$

Therefore, we have an estimate on the size of  $\mathscr{I}$ :

$$\mathbb{P}\left(I \ge \frac{p(1-p)T}{2S}\right) \ge 1 - \exp[-\Theta[p(1-p)T]].$$
(15)

For each  $j \in \{1, \ldots, I\}$ , we define

$$\mathscr{K}_j := \{k \in \mathscr{K} \mid i_{j-1} < k < i_j\},\$$

where we denote  $i_0 := -1$ . Moreover, we define

$$X_{j,\ell} := \sum_{k \in \mathscr{K}_j} z^T ar{A}^{i_\ell - k - 1} ar{d}_k, \quad orall j, \ell \in \{1, \dots, I\} \quad ext{s.t.} \ j \leq \ell.$$

Using equation (14), we can calculate that

$$||z^{T}\mathbf{F}||_{1} \geq \frac{1}{\sqrt{n}} \sum_{\ell=1}^{I} \left| \sum_{j=1}^{\ell} X_{j,\ell} \right|$$

$$\geq \frac{1}{\sqrt{n}} \sum_{j=1}^{I} \left( |X_{j,j}| - \sum_{\ell=j+1}^{I} |X_{j,\ell}| \right).$$
(16)

We utilize the following lemma to bound  $|X_{i,\ell}|$ .

*Lemma 5:* Suppose that a random variable X is sub-Gaussian with parameter  $\sigma_X$ , where the mean and the variance of X are 0 and  $\tilde{\sigma}_X^2$ , respectively. Then, we have

$$\mathbb{P}\left(|X| \geq ilde{\sigma}_X\right) \geq rac{ ilde{\sigma}_X^4}{64\sigma_Y^4}.$$

For all  $j \in \{1, ..., I\}$ , the stealthy assumption (Assumption 2) implies that the standard deviation and the sub-Gaussian parameter of  $X_{j,\ell}$  are

$$\begin{split} \tilde{\sigma}_{j,\ell} &:= \sqrt{\frac{1}{n} \sum_{k \in \mathscr{K}_j} \| z^T \bar{A}^{i_\ell - k - 1} \|_2^2 \sigma_k^2}, \\ \sigma_{j,\ell} &:= \sqrt{\frac{1}{n} \sum_{k \in \mathscr{K}_j} \| z^T \bar{A}^{i_\ell - k - 1} \|_2^2 \sigma^2}, \end{split}$$

respectively. It follows from Lemma 5 that

$$\mathbb{P}(|X_{j,j}| \ge \tilde{\sigma}_{j,j}) \ge \frac{\tilde{\sigma}_{j,j}^4}{64\sigma_{j,j}^4}$$

which further leads to

$$\mathbb{P}(|X_{j,j}| \ge c\sigma_{j,j}) \ge \frac{c^4}{64}.$$
(17)

On the other hand, the sub-Gaussian parameter of  $\sum_{\ell=j+1}^{I} |X_{j,\ell}|$  is at most

$$\sum_{\ell=j+1}^{I} \sigma_{j,\ell} \leq \sum_{\ell=j+1}^{I} \rho^{(\ell-j)S} \sigma_{j,j} \leq \frac{\rho^S}{1-\rho^S} \sigma_{j,j}.$$

Therefore, it holds with probability at least  $1 - \delta/(4I)$  that

$$-\sum_{\ell=j+1}^{I} |X_{j,\ell}| \ge -\frac{\rho^{S}}{1-\rho^{S}} \sigma_{j,j} \cdot \sqrt{2\log(4I/\delta)}$$
(18)  
$$\ge -\frac{\rho^{S}}{1-\rho^{S}} \sigma_{j,j} \cdot \sqrt{2\log(4|\mathscr{I}_{1}|/\delta)}$$
$$\ge -\frac{c^{4}}{512} \cdot c \sigma_{j,j},$$

where the last step is by the choice of S. Using the bound in (15), if we choose

$$T \geq \Theta\left(\frac{\log\log(1/\delta) + \log(1/c)}{p(1-p)c^4\log(1/\rho)}\right),$$

it holds with high probability that

$$\frac{c^4}{64} - \frac{\delta}{4I} \ge \frac{c^4}{128}$$

Note that we have dropped the  $|\mathscr{I}_1|$  term in the definition of *S* since  $\log \log(|\mathscr{I}_1|)$  is bounded by  $\log \log(T)$  and will not change the order of the above bound. Let  $q_j$  be the  $(1-c^4/128)$ -quantile of  $|X_{j,j}| - \sum_{\ell=j+1}^{I} |X_{j,\ell}|$ . We define the indicator function

$$\mathbf{1}_j := \begin{cases} 1, \text{ if } |X_{j,j}| - \sum_{\ell=j+1}^I |X_{j,\ell}| \ge q_j, \\ 0, \text{ otherwise,} \end{cases} \quad \forall j \in \{1, \dots, I\}.$$

Since the value of the Bernoulli random variable  $\mathbf{1}_j$  only depends on attacks in  $\mathcal{K}_j$ , which are disjoint from each other, the random variables

$$\mathbf{1}_1 - c^4/128, \ \ldots, \ \mathbf{1}_I - c^4/128$$

form a martingale sequence with respect to filtration  $\mathscr{F}_{i_1}, \ldots, \mathscr{F}_{i_l}$ . For all  $j \in \{1, \ldots, I\}$ , we can calculate that

$$\mathbb{E}\left[\exp\left(s\mathbf{1}_{j}\right)\right] \leq \exp\left[\frac{c^{4}}{128}\left(e^{s}-1\right)\right], \quad \forall s \in \mathbb{R}.$$

By the tower property of expectation, we have

$$\mathbb{E}\left[\exp\left(s\sum_{j=1}^{I}\mathbf{1}_{j}\right)\right] \leq \exp\left[\frac{c^{4}I}{128}\left(e^{s}-1\right)\right], \quad \forall s \in \mathbb{R}.$$

Therefore, by applying Chernoff's bound and choosing  $s := -\log(2)$ , it follows that

$$\mathbb{P}\left(\sum_{j=1}^{I} \mathbf{1}_{j} \leq \frac{c^{4}}{256} \cdot I\right) \leq \exp\left[-\frac{c^{4}I}{256} \cdot s + \frac{c^{4}I}{128} \left(e^{s} - 1\right)\right]$$
$$\leq \exp\left[-\frac{c^{4}I}{256} \cdot \log\left(\frac{1}{2}\right) - \frac{c^{4}I}{128} \cdot \frac{1}{2}\right]$$
$$= \exp\left[-\Theta\left(\frac{c^{4}}{128} \cdot I\right)\right].$$

Equivalently, we know

$$\mathbb{P}\left(\sum_{j=1}^{I} \mathbf{1}_{j} \ge \frac{c^{4}}{256} \cdot I\right) \ge 1 - \exp\left[-\Theta\left(\frac{c^{4}}{128} \cdot I\right)\right].$$
(19)

Furthermore, since  $i_j - 1 \in \mathscr{K}_j$ , we can estimate that

$$\sigma_{j,j} \geq \sqrt{\frac{1}{n}} \|z\|_2^2 \sigma^2 = \frac{1}{\sqrt{n}} \sigma.$$

By the definition of  $q_j$  and  $\mathbf{1}_j$ , when the event in inequality (19) happens, inequalities (17) and (18) imply that

$$\begin{aligned} \|z^{T}\mathbf{F}\|_{1} &\geq \frac{1}{\sqrt{n}} \sum_{j=1}^{I} \left( |X_{j,j}| - \sum_{\ell=j+1}^{I} |X_{j,\ell}| \right) \\ &\geq \frac{1}{\sqrt{n}} \sum_{j=1}^{I} \left[ \frac{c^{4}}{256} \cdot c \sigma_{j,j} - \frac{c^{4}}{512} \cdot c \sigma_{j,j} \right] \geq \frac{c^{5} \sigma}{512n} \cdot I \end{aligned}$$

holds with probability at least  $1 - \delta/4$ . Hence, we obtain

$$\mathbb{P}\left[\|\boldsymbol{z}^{T}\mathbf{F}\|_{1} \geq \frac{c^{5}\boldsymbol{\sigma}}{512n} \cdot \boldsymbol{I}\right] \geq 1 - \exp\left[-\Theta\left(c^{4}\boldsymbol{I}\right)\right] - \frac{\delta}{4}.$$
 (20)

b) Step 1-2: For the term  $z^T g$ , we can establish an upper bound on

$$\mathbb{E}\left[\exp\left(\lambda \cdot z^{T}g\right)\right] = \mathbb{E}\left[\exp\left(\lambda \sum_{k \in \mathscr{K}} \sum_{i \in \mathscr{K}} z^{T} \bar{A}^{(i-k-1)_{+}} \bar{d}_{k} \cdot \bar{f}_{i}^{l}\right)\right].$$

Define the filtration

$$\mathscr{F}^f := \mathbf{\sigma}\{\overline{f}_t, t \in \mathscr{K}\}.$$

By the stealth assumption, for each  $k \in \mathcal{K}$ , conditional on  $\mathcal{F}_k$ and  $\mathcal{F}^f$ , we have

# $\bar{\ell}_k$ is sub-Gaussian with parameter $\sigma$ .

Let T' be the second last time instance in  $\mathcal{K}$ . We have

$$\mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}\sum_{k\in\mathscr{K}}z^{T}\bar{A}^{(i-k-1)_{+}}\bar{d}_{k}\cdot\bar{f}_{i}^{l}\right)\right]$$
(21)  
$$=\mathbb{E}\left[\exp\left(\lambda\sum_{k\in\mathscr{K},k< T'}\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-k-1)_{+}}\bar{d}_{k}\cdot\bar{f}_{i}^{l}\right) \times \mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-1-T')_{+}}\bar{d}_{T}'\cdot\bar{f}_{i}^{l}\right)\middle|\mathscr{F}_{T'},\mathscr{F}^{f}\right]\right].$$

Using the decomposition in Assumption 2, we have

$$\mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-1-T')_{+}}\bar{d}_{T'}\cdot\bar{f}_{i}^{l}\right)\middle|\mathscr{F}_{T'},\mathscr{F}^{f}\right]$$
$$=\mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-1-T')_{+}}\bar{f}_{T'}\bar{f}_{i}^{l}\cdot\bar{\ell}_{T'}\right)\middle|\mathscr{F}_{T'},\mathscr{F}^{f}\right]$$
$$\leq \exp\left[\frac{\lambda^{2}\sigma^{2}}{2}\left(\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-1-T')_{+}}\bar{f}_{T'}\bar{f}_{i}^{l}\right)^{2}\right].$$

Substituting back into (21), it follows that

$$\mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}\sum_{k\in\mathscr{K}}z^{T}\bar{A}^{(i-k-1)_{+}}\bar{d}_{k}\cdot\bar{f}_{i}^{l}\right)\right]$$
  
$$\leq \mathbb{E}\left[\exp\left(\lambda\sum_{k\in\mathscr{K},k< T'}\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-k-1)_{+}}\bar{d}_{k}\cdot\bar{f}_{i}^{l}\right)\right.$$
  
$$\times \exp\left[\frac{\lambda^{2}\sigma^{2}}{2}\left(\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-1-T')_{+}}\bar{f}_{T'}\bar{f}_{i}^{l}\right)^{2}\right]\right]$$

Continuing the process for all  $k \in \mathcal{K}$ , we obtain

$$\mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}\sum_{k\in\mathscr{K}}z^{T}\bar{A}^{(i-k-1)_{+}}\bar{d}_{k}\cdot\bar{f}_{i}^{l}\right)\right]$$
(22)  
$$\leq \mathbb{E}\left[\exp\left[\frac{\lambda^{2}\sigma^{2}}{2}\sum_{k\in\mathscr{K}}\left(\sum_{i\in\mathscr{K}}z^{T}\bar{A}^{(i-1-k)_{+}}\bar{f}_{k}\bar{f}_{i}^{l}\right)^{2}\right]\right]$$
$$\leq \mathbb{E}\left[\exp\left[\frac{\lambda^{2}\sigma^{2}}{2}\sum_{k\in\mathscr{K}}\left(\sum_{i\in\mathscr{K}}\left|z^{T}\bar{A}^{(i-1-k)_{+}}\bar{f}_{k}\right|\right)^{2}\right]\right],$$

where the last inequality holds because  $\bar{f}_i^l$  is bounded in [-1,1]. For each  $i,k \in \mathcal{K}$ , the value of  $\left(z^T \bar{A}^{(i-1-k)_+} \bar{f}_k\right)^2$  concentrates around its expectation  $||z^T \bar{A}^{(i-1-k)_+}||_2^2/n$ . Therefore, inequality (22) leads to

$$\mathbb{E}\left[\exp\left(\lambda\sum_{i\in\mathscr{K}}\sum_{k\in\mathscr{K}}z^{T}\bar{A}^{(i-k-1)_{+}}\bar{d}_{k}\cdot\bar{f}_{i}^{l}\right)\right]$$
(23)  
$$\leq \exp\left[\Theta\left[\frac{\lambda^{2}\sigma^{2}}{2n}\sum_{k\in\mathscr{K}}\left(\sum_{i\in\mathscr{K}}\left\|z^{T}\bar{A}^{(i-1-k)_{+}}\right\|_{2}\right)^{2}\right]\right]$$
$$\leq \exp\left[\Theta\left[\frac{\lambda^{2}\sigma^{2}}{2n}\sum_{k\in\mathscr{K}}\left(\sum_{i\in\mathscr{K}}\rho^{(i-k-1)_{+}}\right)^{2}\right]\right].$$

Suppose the elements in  $\mathcal{K}$  are

$$j_1 < j_2 < \cdots < j_{|\mathscr{K}|}$$

Define

$$\Delta_k := j_k - j_{k-1} - 1, \quad \forall k \in \{2, \dots, |\mathscr{K}|\}.$$

We can calculate that

$$\sum_{i\in\mathscr{K}}\rho^{(i-1-j_k)_+}\leq \frac{\rho^{\Delta_k}}{1-\rho}$$

Since  $\rho^{\Delta_k} \in [0,1]$  are bounded random variables, they are sub-Gaussian and concentrate around the mean with high probability. The expectation of  $\rho^{2\Delta_k}$  is

$$\sum_{\Delta=0}^{\infty} p(1-p)^{\Delta} \rho^{2\Delta} = \frac{p}{1-(1-p)\rho^2}$$

Therefore, with probability at least  $1 - \exp[-\Theta(pT)]$ , we have

$$\sum_{k=2}^{|\mathscr{K}|} \rho^{2\Delta_k} \lesssim \frac{|\mathscr{K}|p}{1-(1-p)\rho^2} \leq \frac{|\mathscr{K}|p}{1-\rho}.$$

Hence, inequality (23) implies that with the same probability,  $z^T g$  is sub-Gaussian with parameter

$$\Theta\left[\sqrt{\frac{\sigma^2}{n}\sum_{k=2}^{|\mathscr{K}|}\frac{\rho^{2\Delta_k}}{(1-\rho)^2}}\right] \le \Theta\left[\sqrt{\frac{|\mathscr{K}|p\sigma^2}{n(1-\rho)^3}}\right]$$

Therefore, Hoeffding's inequality leads to

$$\mathbb{P}\left[z^{T}g \leq -\Theta\left(\sqrt{\frac{|\mathscr{K}|p\sigma^{2}}{n(1-\rho)^{3}}\log\left(\frac{4}{\delta}\right)}\right)\right] \leq \frac{\delta}{4}.$$
 (24)

By combining inequalities (20) and (24), it holds with probability at least

$$1 - \exp\left[-\Theta(c^4 I)\right] - \frac{\delta}{2}$$

that

$$f(z) \ge \Theta\left[\frac{c^{5}\sigma I}{n} - \sqrt{\frac{|\mathscr{K}|p\sigma^{2}}{n(1-\rho)^{3}}\log\left(\frac{1}{\delta}\right)}\right]$$

Similar to the bound in (15), it holds with probability at least  $1 - \exp[-\Theta(pT)]$  that

$$|\mathscr{K}| \leq 2pT$$

As a result, if we choose

$$T \ge \Theta \left[ \max \left\{ \frac{\log \log(1/\delta) + \log(1/c)}{c^4 p(1-p) \log(1/\rho)} \log\left(\frac{1}{\delta}\right), \quad (25) \right. \\ \left. \frac{1}{p(1-p)} \log\left(\frac{1}{\delta}\right), \\ \left. \frac{n \log(1/c)^2}{c^{10}(1-p)^2(1-\rho)^3 \log^2(1/\rho)} \log\left(\frac{1}{\delta}\right) \right\} \right] \\ = \Theta \left[ nR \log\left(\frac{1}{\delta}\right) \right],$$

where

$$\begin{split} R &:= \max \left\{ \frac{\log(1/c)}{c^4 p(1-p)\log(1/\rho)} \log\left(\frac{1}{\delta}\right), \\ & \frac{\log^2(1/c)}{c^{10}(1-p)^2(1-\rho)^3 \log^2(1/\rho)}, \frac{1}{np(1-p)} \right\}, \end{split}$$

we have

$$\mathbb{P}\left[f(z) \ge \Theta\left(\frac{c^5 \sigma I}{n}\right)\right] \ge 1 - \delta.$$
(26)

Step 2: In the second step, we apply discretization techniques to prove that condition (13) holds for all  $z \in S_2(1)$  with high probability. Suppose that  $\varepsilon > 0$  is a small constant. We construct an  $\varepsilon$ -cover of the unit sphere  $S_2(1)$ , denoted as

$$\{z^1,\ldots,z^N\},\$$

Namely, for all  $z \in S_2(1)$ , we can find  $r \in \{1, 2, ..., N\}$  such that  $||z - z^r||_2 \le \varepsilon$ . The number of points N can be bounded by

$$\log(N) \leq \log[\mathscr{N}(\varepsilon, \mathbb{S}_2(1), \|\cdot\|_2)] \leq n \log\left(1 + \frac{2}{\varepsilon}\right).$$

Define a to be the lower bound of f(z) in inequality (26). Then, we have

$$a = \Theta\left(\frac{c^5 \sigma I}{n}\right).$$

Our goal is to prove that

$$f(z) - f(z') \ge -a, \quad \forall z, z' \in \mathbb{S}_2(1) \text{ s.t. } \|z - z'\|_2 \le \varepsilon$$

holds with high probability. Notice that

$$\begin{split} f(z) - f(z') &= (z - z')^T g + (\|z^T \mathbf{F}\|_1 - \|(z')^T \mathbf{F}\|_1) \\ &\geq (z - z')^T g - \|(z - z')^T \mathbf{F}\|_1 \\ &\geq -\|z - z'\|_2 \|g\|_2 - \|z - z'\|_2 \sum_{i \notin \mathscr{M}} \|\mathbf{F}^i\|_2 \\ &\geq -\varepsilon \left( \left\| \sum_{i \in \mathscr{K}} \sum_{k \in \mathscr{M}} \bar{A}^{(i-k-1)_+} \bar{d}_k \right\|_2 \right) \\ &+ \frac{1}{\sqrt{n}} \sum_{i \notin \mathscr{M}} \left\| \sum_{k \in \mathscr{M}} \bar{A}^{(i-k-1)_+} \bar{d}_k \right\|_2 \right) \\ &\geq -\varepsilon \sum_{k \in \mathscr{K}} \sum_{i > k} \rho^{(i-k-1)} |\bar{\ell}_k|. \end{split}$$

Using the property of exponential sequences, we have

$$\sum_{k\in\mathscr{K}}\sum_{i>k}\rho^{(i-k-1)}|\bar{\ell}_k|\leq \frac{1}{1-\rho}\sum_{k\in\mathscr{K}}|\bar{\ell}_k|.$$

Using a similar proof, we can show that  $\sum_{k \in \mathscr{K}} |\bar{\ell}_k|$  is sub-Gaussian with parameter  $|\mathscr{K}|\sigma$ . Therefore, Hoeffding's inequality implies that

$$\mathbb{P}\left(\frac{1}{1-\rho}\sum_{k\in\mathscr{K}}|\bar{\ell}_k|>\frac{a}{\varepsilon}\right)\leq 2\exp\left[-\frac{(1-\rho)^2a^2}{2\varepsilon^2|\mathscr{K}|^2\sigma^2}\right].$$

Letting

$$\varepsilon := \frac{(1-\rho)a}{|\mathscr{K}|\sigma\sqrt{2\log(4/\delta)}},$$

it holds that

$$\mathbb{P}\left[f(z) - f(z') \ge -a, \quad \forall z, z' \in \mathbb{S}_2(1) \text{ s.t. } \|z - z'\|_2 \le \varepsilon\right]$$
$$\ge \mathbb{P}\left(\frac{1}{1 - \rho} \sum_{k \in \mathscr{K}} |\bar{\ell}_k| \le \frac{a}{\varepsilon}\right) \ge 1 - \frac{\delta}{2}.$$

Now, after we replace  $\delta$  in (25) with  $\delta/(2N)$ , it holds with probability at least  $1 - \delta/2$  that

$$f(z^r) \ge a, \quad \forall r \in \{1, \dots, N\}$$

After combining the above two inequalities, we apply the union bound to obtain

$$\mathbb{P}[f(z) \ge 0, \quad \forall z \in \mathbb{S}_2(1)] \ge 1 - \delta.$$

The corresponding sample complexity is

$$T \ge \Theta\left[nR\log\left(\frac{2N}{\delta}\right)\right].$$

Since it holds with probability  $1 - \exp[-\Theta[p(1-p)T]]$  that

$$|\mathscr{I}_1| = \Theta[p(1-p)T], \quad |\mathscr{K}| = \Theta(pT),$$

we get the estimate

$$\log(N) \le n \log\left(1 + \frac{2}{\varepsilon}\right)$$
$$= n \log\left[1 + \Theta\left(\frac{n\sqrt{\log(1/\delta)}\log(1/c)}{(1-p)c^5(1-\rho)\log(1/\rho)}\right)\right]$$
$$= \Theta[n \log(nR)].$$

By omitting the constants in the expression, the final sample complexity can be written as

$$T \ge \Theta \left[ nR \left[ n\log\left(nR\right) + \log\left(\frac{1}{\delta}\right) \right] \right]$$

Finally, we replace  $\delta$  with  $\delta/n$  and apply the union bound to all coordinates  $\ell \in \{1, ..., n\}$ . The sample complexity remains on the same order as the above expression.

6) Proof of Theorem 4: Due to the system dynamics and given  $x_0 = 0$ ,  $x_i$  can be expressed as

$$x_i = \sum_{k \notin \mathscr{K}} \bar{A}^{(i-k-1)_+} \bar{B} u_k + \sum_{k \in \mathscr{K}} \bar{A}^{(i-k-1)_+} (\bar{B} u_k + \bar{d}_k).$$

From Proposition 4, we want to show that there exist scalars  $\gamma_i^l, \mu_i^l \in [-1, 1]$  for all  $i \notin \mathcal{K}, l \in \{1, ..., n\}$  such that

$$\sum_{i \notin \mathscr{K}} \gamma_i^l x_i / \sqrt{n} = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_2^l x_i, \quad \forall l = 1, \dots, n,$$
(27)

and

$$\sum_{i \notin \mathscr{K}} \mu_i^l u_i / \sqrt{n} = \sum_{i \in \mathscr{K}} \partial \|\bar{d}_i\|_2^l u_i, \quad \forall l = 1, \dots, n.$$
 (28)

We finish the proof in two steps.

Step 1: We first analyze condition (27) with a given coordinate  $l \in \{1, ..., n\}$ . From Lemma 3, condition (27) is equivalent to

$$f(z) := z^T g + ||z^T \mathbf{F}||_1 \ge 0, \quad \forall z \in \mathbb{S}_2(1),$$

where the matrix  $\mathbf{F} \in \mathbb{R}^{n \times (T - |\mathscr{K}|)}$  has the columns

$$\mathbf{F}^{i} := \sum_{k \notin \mathscr{K}} \frac{\bar{A}^{(i-k-1)_{+}} \bar{B} u_{k}}{\sqrt{n}} + \sum_{k \in \mathscr{K}} \frac{\bar{A}^{(i-k-1)_{+}} (\bar{B} u_{k} + \bar{d}_{k})}{\sqrt{n}}, \ \forall i \notin \mathscr{K},$$

and the vector  $g \in \mathbb{R}^n$  is

$$g := \sum_{i \in \mathscr{K}} \left[ \sum_{k \notin \mathscr{K}} \bar{A}^{(i-k-1)_+} \bar{B} u_k + \sum_{k \in \mathscr{K}} \bar{A}^{(i-k-1)_+} (\bar{B} u_k + \bar{d}_k) \right] \bar{f}_i^l.$$

Similar to the proof of Theorem 2, we first prove that  $f(z) \ge a$ holds with high probability for a fixed  $z \in S_2(1)$  and some positive constant *a*. For each  $k \notin \mathcal{H}$ , the standard deviation and sub-Gaussian parameter of  $z^T \bar{A}^{(i-k-1)_+} \bar{B}u_k$  are both

$$\frac{1}{\sqrt{m}} \| z^T \bar{A}^{(i-k-1)_+} \bar{B} \|_2 \xi$$

For each  $k \in \mathscr{K}$ , the standard deviation and sub-Gaussian parameter of  $z^T \bar{A}^{(i-k-1)_+}(\bar{B}u_k + \bar{d}_k)$  are, respectively,

$$\sqrt{\frac{1}{m}} \|z^T \bar{A}^{(i-k-1)_+} \bar{B}\|_2^2 \xi^2 + \frac{1}{n} \|z^T \bar{A}^{(i-k-1)_+}\|_2^2 \sigma_k^2, 
\sqrt{\frac{1}{m}} \|z^T \bar{A}^{(i-k-1)_+} \bar{B}\|_2^2 \xi^2 + \frac{1}{n} \|z^T \bar{A}^{(i-k-1)_+}\|_2^2 \sigma^2.$$

Note that we have utilized the independence between  $u_k$  and  $\bar{d}_k$  in the above calculation. Let

$$S := \left[ \log_{\rho} \Theta \left[ \frac{\sqrt{\frac{1}{m} \eta_B^2 \xi^2 + \frac{p}{n} \sigma^2} \cdot c^5}{\sqrt{\frac{1}{m} \rho_B^2 \xi^2 + \frac{p}{n} \sigma^2} \cdot \sqrt{\log(1/\delta)}} \right] \right],$$

where  $\rho_B$  is the maximal singular value of  $\overline{B}$  and  $\eta_B$  is the minimal singular value of the matrix

$$\frac{1}{(1-\rho)^2} \begin{bmatrix} \bar{B} & \bar{A}\bar{B} & \cdots & \bar{A}^{n-1}\bar{B} \end{bmatrix}$$

By the controllability assumption, the above matrix is rank*n* and thus, the parameter  $\eta_B$  is strictly positive. We define  $i_0 := -1$  and construct the index set

$$\mathscr{I} := \{i_1, \dots, i_I \mid i_j \notin \mathscr{K}, \ i_j - i_{j-1} \ge S, \quad \forall j\}$$

It is straightforward to construct  $\mathscr{I}$  such that  $I = |\mathscr{I}|$  is on the order of

$$\min\left\{(1-p)T,\frac{T}{S}\right\}$$

For each  $j \in \{1, \ldots, I\}$ , we define

$$\mathscr{K}_j := \{k \in \mathscr{K} \mid i_{j-1} \le k < i_j\}, \ \mathscr{K}_j^c := \{k \notin \mathscr{K} \mid i_{j-1} \le k < i_j\}.$$

Moreover, we define

$$\begin{aligned} X_{j,\ell} &:= \sum_{k \in \mathscr{K}_j} z^T \bar{A}^{i_\ell - k - 1} \left( \bar{B} u_k + \bar{d}_k \right) + \sum_{k \in \mathscr{K}_j^c} z^T \bar{A}^{i_\ell - k - 1} \bar{B} u_k, \\ \forall j, \ell \in \{1, \dots, I\} \quad \text{s.t. } j \le \ell. \end{aligned}$$

For all  $j \in \{1, ..., I\}$ , the stealthy assumption (Assumption 2) implies that the standard deviation and the sub-Gaussian parameter of  $X_{j,\ell}$  is

$$\begin{split} \tilde{\sigma}_{j,\ell} &:= \\ \sqrt{\frac{1}{m} \sum_{k \in \mathscr{K}_j \cup \mathscr{K}_j^c} \| z^T \bar{A}^{i_\ell - k - 1} \bar{B} \|_2^2 \xi^2} + \frac{1}{n} \sum_{k \in \mathscr{K}_j} \| z^T \bar{A}^{i_\ell - k - 1} \|_2^2 \sigma_k^2, \\ \sigma_{j,\ell} &:= \\ \sqrt{\frac{1}{m} \sum_{k \in \mathscr{K}_j \cup \mathscr{K}_j^c} \| z^T \bar{A}^{i_\ell - k - 1} \bar{B} \|_2^2 \xi^2} + \frac{1}{n} \sum_{k \in \mathscr{K}_j} \| z^T \bar{A}^{i_\ell - k - 1} \|_2^2 \sigma^2, \end{split}$$

respectively. Define

$$c_{j,\ell} := rac{ ilde{\sigma}_{j,\ell}}{\sigma_{j,\ell}}, \quad \forall j,\ell \in \{1,\ldots,I\} \quad \text{s.t. } j \leq \ell$$

Similar to the proof of Theorem 2, we have the bound

$$||z^T \mathbf{F}||_1 \ge \frac{1}{\sqrt{n}} \sum_{j=1}^{I} \left( |X_{j,j}| - \sum_{\ell=j+1}^{I} |X_{j,\ell}| \right).$$

By Lemma 5, we have

$$\mathbb{P}(|X_{j,j}| \ge c_{j,j}\sigma_{j,j}) \ge \frac{c_{j,j}^4}{64}.$$
(29)

For all vector  $y \in \mathbb{R}^n$ , the controllability assumption leads to

$$\sum_{k=0}^{n-1} \|y^{T} \bar{A}^{k} \bar{B}\|_{2}^{2} \geq \frac{\eta_{B}^{2}}{(1-\rho)^{2}} \cdot \|y\|_{2}^{2}$$
(30)  
$$\geq \frac{\eta_{B}^{2}}{(1-\rho)^{2}} \cdot (1-\rho)^{2} \sum_{k=0}^{n} \rho^{2k} \|y\|_{2}^{2} \geq \eta_{B}^{2} \sum_{k=0}^{n} \|y^{T} \bar{A}^{k}\|_{2}^{2}.$$

Therefore, we can divide the set  $\mathscr{K}_j \cup \mathscr{K}_j^c$  into segments with *n* consecutive time instances and apply inequality (30) to each segment. When *T* is large enough such that  $I \ge \Theta(n)$ , we obtain the estimation

$$\sum_{k \in \mathscr{K}_j \cup \mathscr{K}_j^c} \| z^T \bar{A}^{i_\ell - k - 1} \bar{B} \|_2^2 \gtrsim \sum_{k = i_{j-1}}^{i_j - 1} \| z^T \bar{A}^{i_\ell - k - 1} \|_2^2 \eta_B^2.$$

Applying concentration inequalities to set  $\mathcal{K}_j$ , the distribution of its elements will surround their expected values. Therefore, for the simplicity of presentation, we use the following approximation:

$$\sigma_{j,j}^2 \gtrsim rac{1}{m} \eta_B^2 \xi^2 + rac{p}{n} \sigma^2 := ar\sigma^2.$$

In addition, the parameter  $c_{j,j}$  can be bounded by

For the numerator, we can estimate that

$$\sum_{k \in \mathscr{K}_j \cup \mathscr{K}_j^c} \| z^T \bar{A}^{i_\ell - k - 1} \bar{B} \|_2^2 \gtrsim \sum_{k = i_{j-1}}^{i_j - 1} \| z^T \bar{A}^{i_\ell - k - 1} \|_2^2 \eta_B^2.$$

On the other hand, since

$$\sum_{\substack{k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c} \\ k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c}}} \| z^{T} \bar{A}^{i_{\ell}-k-1} \bar{B} \|_{2}^{2} \xi^{2} \leq \sum_{\substack{k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c} \\ k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c}}} \| z^{T} \bar{A}^{i_{\ell}-k-1} \|_{2}^{2} \sigma^{2} \lesssim p \sum_{\substack{k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c} \\ k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c}}} \| z^{T} \bar{A}^{i_{\ell}-k-1} \|_{2}^{2} \sigma^{2},$$

we get

$$c_{j,j}^2 \gtrsim rac{rac{1}{m}\eta_B^2\xi^2}{rac{1}{m}
ho_B^2\xi^2 + rac{p}{n}\sigma^2} := c.$$

17

Therefore, inequality (29) implies

$$\mathbb{P}(|X_{j,j}| \ge c\bar{\sigma}) \ge \frac{c^4}{64}.$$
(31)

Since the sub-Gaussian parameter of  $\sum_{\ell=j+1}^{I} |X_{j,\ell}|$  is  $\sum_{\ell=j+1}^{I} \sigma_{j,\ell}$ , Hoeffding's inequality implies that

$$\mathbb{P}\left(\sum_{\ell=j+1}^{I} |X_{j,\ell}| \le \sum_{\ell=j+1}^{I} \sigma_{j,\ell} \cdot \sqrt{2\log\left(\frac{4I}{\delta}\right)}\right) \ge 1 - \frac{\delta}{4I}.$$
 (32)

We can bound the sub-Gaussian parameter by

$$\begin{split} &\sum_{\ell=j+1}^{I} \sigma_{j,\ell} \\ &\leq \sum_{\ell=j+1}^{I} \sqrt{\frac{1}{m} \sum_{k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c}} \rho^{2(i_{\ell}-k-1)} \rho_{B}^{2} \xi^{2} + \frac{1}{n} \sum_{k \in \mathscr{K}_{j}} \rho^{2(i_{\ell}-k-1)} \sigma^{2}} \\ &\leq \frac{\rho^{S}}{1-\rho^{S}} \sqrt{\frac{1}{m} \sum_{k \in \mathscr{K}_{j} \cup \mathscr{K}_{j}^{c}} \rho^{2(i_{j}-k-1)} \rho_{B}^{2} \xi^{2} + \frac{1}{n} \sum_{k \in \mathscr{K}_{j}} \rho^{2(i_{j}-k-1)} \sigma^{2}} \\ &\leq \frac{\rho^{S}}{1-\rho^{S}} \sqrt{\frac{1}{m(1-\rho)} \rho_{B}^{2} \xi^{2} + \frac{1}{n} \sum_{k \in \mathscr{K}_{j}} \rho^{2(i_{j}-k-1)} \sigma^{2}}. \end{split}$$

In the same way, we have the following bound with high probability:

$$\sum_{k \in \mathscr{K}_j} \rho^{2(i_j - k - 1)} \lesssim p \sum_{k \in \mathscr{K}_j \cup \mathscr{K}_j^c} \rho^{2(i_j - k - 1)} \leq \frac{p}{1 - \rho},$$

which holds with high probability when T is large. Therefore, we have the bound

$$\sum_{\ell=j+1}^{I} \sigma_{j,\ell} \lesssim \frac{\rho^{S}}{1-\rho^{S}} \sqrt{\frac{1}{m(1-\rho)}\rho_{B}^{2}\xi^{2} + \frac{p}{n(1-\rho)}\sigma^{2}}$$
$$:= \frac{\rho^{S}}{1-\rho^{S}}\tilde{\sigma}.$$

By the choice of S, we get

$$\sum_{\ell=j+1}^{I} \sigma_{j,\ell} \lesssim \frac{c^4}{256} \cdot c\bar{\sigma} \cdot \left(\sqrt{2\log\left(\frac{4I}{\delta}\right)}\right)^{-1},$$

Therefore, inequality (32) leads to

$$\mathbb{P}\left(\sum_{\ell=j+1}^{I} |X_{j,\ell}| \le \frac{c^4}{256} \cdot c\bar{\sigma}\right) \ge 1 - \frac{\delta}{4I}.$$
 (33)

Choosing

$$T \ge \Theta\left(\frac{\log\log(1/\delta)}{c^4 \min\{1-p, 1/S\}}\right),$$

we have

$$\frac{c^4}{64} - \frac{\delta}{4I} \ge \frac{c^4}{128}.$$

By the same construction of the martingale sequence and the application of Azuma-Hoeffding's inequality, inequalities (29)

and (32) imply that

$$\|\boldsymbol{z}^{T}\mathbf{F}\|_{1} \geq \frac{1}{\sqrt{n}} \sum_{j=1}^{I} \left( |\boldsymbol{X}_{j,j}| - \sum_{\ell=j+1}^{I} |\boldsymbol{X}_{j,\ell}| \right)$$

$$\geq \frac{1}{\sqrt{n}} \left( \frac{c^{4}I}{256} \cdot c\bar{\boldsymbol{\sigma}} - \frac{c^{4}I}{512} c\bar{\boldsymbol{\sigma}} \right) = \frac{c^{5}\bar{\boldsymbol{\sigma}}}{512\sqrt{n}} \cdot I$$
(34)

holds with probability at least

$$1 - \exp[-\Theta(c^4 I)] - \delta/4.$$

On the other hand, for the term  $z^T g$ , we can bound its sub-Gaussian parameter by

$$\sqrt{\frac{|\mathscr{K}|\rho_B^2\xi^2}{m(1-\rho)^3} + \frac{|\mathscr{K}|p\sigma^2}{n(1-\rho)^3}} = \sqrt{\frac{|\mathscr{K}|}{(1-\rho)^2}}\tilde{\sigma}.$$

Then, Hoeffding's inequality leads to

$$\mathbb{P}\left[z^{T}g \leq -\Theta\left(\sqrt{\frac{|\mathscr{K}|\tilde{\sigma}^{2}}{(1-\rho)^{2}}\log\left(\frac{4}{\delta}\right)}\right)\right] \leq \frac{\delta}{4}.$$
 (35)

Combining inequalities (34) and (35), it holds with probability at least

$$1 - \exp[-\Theta(c^4 I)] - \frac{\delta}{2}$$

that

$$f(z) \ge \Theta\left[\frac{c^5 \bar{\sigma} I}{\sqrt{n}} - \sqrt{\frac{\tilde{\sigma}^2 |\mathscr{K}|}{(1-\rho)^2} \log\left(\frac{1}{\delta}\right)}\right]$$

Similar to the bound in (15), it holds with probability at least  $1 - \exp[-\Theta(pT)]$  that

$$|\mathscr{K}| \leq 2pT.$$

As a result, if we choose

$$T \ge \Theta \left[ \max \left\{ \frac{1}{c^4 \min\{1-p, 1/S\}} \log\left(\frac{1}{\delta}\right), \\ \frac{1}{p} \log\left(\frac{1}{\delta}\right), \\ \frac{np\kappa^2}{c^{10}(1-\rho)^2 \min\{(1-p)^2, 1/S^2\}} \log\left(\frac{1}{\delta}\right) \right\} \right]$$
$$= \Theta \left[ nR_1 \log\left(\frac{1}{\delta}\right) \right],$$

where  $\kappa := \tilde{\sigma}/\bar{\sigma} \ge (1-\rho)^{-1}$  and

$$R_{1} := \max\left\{\frac{1}{c^{4}(1-p)}, \frac{\log(\kappa/c)}{c^{4}\log(1/\rho)}, \frac{p\kappa^{2}}{c^{10}(1-p)^{2}(1-\rho)^{2}}, \frac{p\kappa^{2}\log^{2}(\kappa/c)}{c^{10}(1-\rho)^{2}\log^{2}(1/\rho)}, \frac{1}{np}\right\}$$

we have

$$\mathbb{P}\left[f(z) \ge \Theta\left(\frac{c^5 \bar{\sigma}I}{\sqrt{n}}\right)\right] \ge 1 - \delta.$$
(36)

Next, we apply the discretization techniques and estimate the size of  $\varepsilon$ -net, which we denote as *N*. Similar to the proof of Theorem 2, it is sufficient to choose

$$\log(N) \le n \log\left(1 + \frac{2}{\varepsilon}\right),$$

and

$$arepsilon := \Theta\left(rac{a}{\|g\|_2 + \sum_{i \notin \mathscr{K}} \|\mathbf{F}^i\|_2}
ight).$$

where a > 0 is the lower bound of f(z) in (36). We can estimate that

$$\begin{split} \|g\|_{2} &\leq \sum_{i \in \mathscr{H}} \left\| \sum_{k \notin \mathscr{H}} \bar{A}^{(i-k-1)_{+}} \bar{B}u_{k} + \sum_{k \in \mathscr{H}} \bar{A}^{(i-k-1)_{+}} (\bar{B}u_{k} + \bar{d}_{k}) \right\|_{2} \\ &\leq \frac{\rho_{B}}{1-\rho} \sum_{k=0}^{T-1} \|u_{k}\|_{2} + \frac{1}{1-\rho} \sum_{j \in \mathscr{H}} \|\bar{d}_{k}\|_{2}. \end{split}$$

Therefore, the sub-Gaussian parameter of  $||g||_2$  is bounded by

$$\frac{1}{1-\rho}\sqrt{\rho_B^2 T\xi^2 + |\mathscr{K}|\sigma^2} \lesssim \frac{1}{1-\rho}\sqrt{\rho_B^2 T\xi^2 + pT\sigma^2} := \sigma'\sqrt{T}$$

Similarly, the sub-Gaussian parameter of  $\sum_{i \notin \mathscr{K}} \|\mathbf{F}^i\|_2$  is bounded by

$$\frac{1}{(1-\rho)\sqrt{n}}\sqrt{\rho_B^2 T\xi^2 + pT\sigma^2} = \frac{\sigma'\sqrt{T}}{\sqrt{n}}$$

with high probability. Hoeffding's inequality implies

$$\|g\|_2 + \sum_{i \notin \mathscr{K}} \|\mathbf{F}^i\|_2 \le \Theta \left[\sigma' \sqrt{T} \sqrt{\log\left(\frac{1}{\delta}\right)}\right]$$

with probability at least  $1 - \delta/2$ . With the same probability, we have

$$\log\left(1+\frac{2}{\varepsilon}\right) \leq \log\left[1+\Theta\left(\frac{\sqrt{n}\sigma'\sqrt{\log(1/\delta)}}{c^{5}\bar{\sigma}\min\{1-p,1/S\}\sqrt{T}}\right)\right]$$
$$\leq \Theta[\log(nR_{1})],$$

which further leads to

$$\log(N) \lesssim \Theta[n\log(nR_1)].$$

Replacing  $\delta$  with  $\delta/N$  in (36), the final sample complexity bound is

$$T \ge \Theta\left[nR_1\left[n\log(nR_1) + \log\left(\frac{1}{\delta}\right)\right]\right].$$

Step 2: In the second step, we consider condition (28). From Lemma 3, given a coordinate  $l \in \{1, ..., n\}$ , (28) is equivalent to

$$f(z) := z^T g + \| z^T \mathbf{F} \|_1 \ge 0, \quad \forall z \in \mathbb{S}_2(1),$$

where the matrix  $\mathbf{F} \in \mathbb{R}^{m \times (T - |\mathcal{K}|)}$  has the columns

$$\mathbf{F}^i := \frac{u_i}{\sqrt{n}}, \quad \forall i \notin \mathscr{K},$$

and the vector  $g \in \mathbb{R}^m$  is

$$g := \sum_{i \in \mathscr{K}} u_i \bar{f}_i^l.$$

For a given  $z \in S_2(1)$ , we have

$$\mathbb{E}[f(z)] = \mathbb{E} \| z^T \mathbf{F} \|_1 = \sum_{i \notin \mathscr{K}} \mathbb{E} | z^T \mathbf{F}^i |$$
  
=  $\Theta\left(\frac{(T - |\mathscr{K}|)\xi}{\sqrt{mn}}\right) \gtrsim \Theta\left(\frac{(1 - p)T\xi}{\sqrt{mn}}\right).$ 

The sub-Gaussian parameter of  $||z^T \mathbf{F}||_1 + z^T g$  is

$$\sqrt{\frac{(T-|\mathscr{K}|)\xi^2}{mn}+\frac{|\mathscr{K}|\xi^2}{m}} \lesssim \sqrt{\left(\frac{1-p}{mn}+\frac{p}{m}\right)} \cdot T\xi^2.$$

Therefore, Hoeffding's inequality implies that

$$f(z) \ge \Theta\left(\frac{(1-p)T\xi}{\sqrt{mn}}\right)$$

holds with probability at least

$$1 - \exp\left[-\Theta\left(\frac{(1-p)^2T}{1-p+np}\right)\right].$$

Choosing

$$T \ge \Theta \left[ \max\left\{ \frac{1}{p} \log\left(\frac{1}{\delta}\right), \frac{1-p+np}{(1-p)^2} \log\left(\frac{1}{\delta}\right) \right\} \right]$$
$$= \Theta \left[ \max\left\{ \frac{1}{p} \log\left(\frac{1}{\delta}\right), \frac{np}{(1-p)^2} \log\left(\frac{1}{\delta}\right) \right\} \right],$$

we have

$$\mathbb{P}\left[f(z) \geq \Theta\left(\frac{(1-p)T\xi}{\sqrt{mn}}\right)\right] \geq 1 - \frac{\delta}{2}$$

Similarly, applying the discretization techniques, it is sufficient to choose N points, where

$$\log(N) = m \log \left[ 1 + \Theta \left( \frac{\sum_{i \notin \mathscr{K}} ||u_i||_2 / \sqrt{n} + \sum_{i \in \mathscr{K}} ||u_i||_2}{(1 - p)T\xi / \sqrt{mn}} \right) \right]$$
  
$$\lesssim m \log \left[ 1 + \Theta \left( \frac{\sqrt{(1 - p)T/n + pT\xi} \cdot \sqrt{\log(1/\delta)}}{(1 - p)T\xi / \sqrt{mn}} \right) \right]$$
  
$$= m \log \left[ 1 + \Theta \left( \frac{\sqrt{1 - p + np} \cdot \sqrt{\log(1/\delta)}}{(1 - p)\sqrt{T/m}} \right) \right]$$
  
$$\leq m \log \left[ m \log \left( \frac{1}{\delta} \right) \right].$$

Hence, the overall sample complexity is

$$T \ge \Theta\left[nR_2\left[m\log(nR_2) + \log\left(\frac{1}{\delta}\right)\right]\right],$$

where we define

$$R_2 := \max\left\{\frac{1}{np}, \frac{p}{(1-p)^2}, \frac{m}{n}\right\}$$

Combining the two steps, we get the conclusion of this theorem.

#### B. Proofs for Results in Appendix

1) Proof of Lemma 5: For the notational simplicity, we omit the X in subscripts. Let

$$\eta:=rac{ ilde{\sigma}}{\sigma}, \quad \delta:=1-rac{ ilde{\sigma}^4}{64\sigma^4}.$$

Assume conversely that

$$\mathbb{P}(|X| \ge \eta \sigma) < 1 - \delta.$$

Then, we can calculate that

$$\begin{split} \mathbb{E}(X^2) &= \int_0^\infty \theta^2 \ d\left[-\mathbb{P}(|X| \ge \theta)\right] = \int_0^\infty 2\theta \mathbb{P}(|X| \ge \theta) \ d\theta \\ &\leq (\eta\sigma)^2 + \int_{\eta\sigma}^\infty 2\theta \min\left\{1 - \delta, 2\exp\left(-\frac{\theta^2}{2\sigma^2}\right)\right\} \ d\theta \\ &= (\eta\sigma)^2 + (1 - \delta)\left[(\sigma')^2 - (\eta\sigma)^2\right] \\ &+ \int_{\sigma'}^\infty 2\theta \cdot 2\exp\left(-\frac{\theta^2}{2\sigma^2}\right) \ d\theta \\ &= (\eta\sigma)^2 + (1 - \delta)\left[(\sigma')^2 - (\eta\sigma)^2\right] \\ &+ 4\sigma^2\exp\left(-\frac{(\sigma')^2}{2\sigma^2}\right), \end{split}$$

where we define

$$\sigma' := \sqrt{2\sigma^2 \log\left(rac{2}{1-\delta}
ight)}$$

Rearranging the above inequality, we get

$$\eta^2 \cdot \delta + 2(1-\delta) \log\left(\frac{2}{1-\delta}\right) + 2(1-\delta) \ge \frac{\tilde{\sigma}^2}{\sigma^2}.$$

Hence, it holds that

$$\begin{split} \eta^2 &\geq \frac{1}{\delta} \left[ \frac{\tilde{\sigma}^2}{\sigma^2} - 2(1-\delta) \log\left(\frac{2}{1-\delta}\right) - 2(1-\delta) \right] \\ &> 2 \left[ \frac{\tilde{\sigma}^2}{\sigma^2} + 4(1-\delta) \log\left(\frac{1-\delta}{2}\right) \right], \end{split}$$

where the second inequality holds because  $\delta > 1/2$  and  $\log[2/(1-\delta)] > 1$ . Using the fact that

$$(1-\delta)\log\left(\frac{1-\delta}{2}\right) \ge -\sqrt{1-\delta} \ge -\frac{\tilde{\sigma}^2}{8\sigma^2},$$

we get

$$\eta^2 > rac{ ilde{\sigma}^2}{\sigma^2},$$

which contradicts with the definition of  $\eta$ . Therefore, we have proved that

$$\mathbb{P}(|X| \ge \tilde{\sigma}) \ge \frac{\tilde{\sigma}^4}{64\sigma^4}$$

#### REFERENCES

- H.-F. Chen and L. Guo, *Identification and stochastic adaptive control*. Springer Science & Business Media, 2012.
- [2] E. Hazan, H. Lee, K. Singh, C. Zhang, and Y. Zhang, "Spectral filtering for general linear dynamical systems," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [3] H. Mania, S. Tu, and B. Recht, "Certainty equivalence is efficient for linear quadratic control," Advances in Neural Information Processing Systems, vol. 32, 2019.

- [4] T. Sarkar, A. Rakhlin, and M. A. Dahleh, "Nonparametric finite time LTI system identification," arXiv preprint arXiv:1902.01848, 2019.
- [5] A. Tsiamis, I. Ziemann, N. Matni, and G. J. Pappas, "Statistical learning theory for control: A finite sample perspective," *arXiv preprint arXiv:2209.05423*, 2022.
- [6] A. Alan, A. J. Taylor, C. R. He, A. Ames, and G. Orosz, "Control barrier functions and input-to-state safety with application to automated vehicles," *IEEE Transactions on Control Systems Technology*, vol. 31, pp. 2744–2759, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:249461776
- [7] L. Wang, E. A. Theodorou, and M. Egerstedt, "Safe learning of quadrotor dynamics using barrier certificates," 2018 IEEE International Conference on Robotics and Automation (ICRA), pp. 2460–2465, 2017. [Online]. Available: https://api.semanticscholar.org/CorpusID:35948052
- [8] S. M. Khansari-Zadeh and A. Billard, "Learning control lyapunov function to ensure stability of dynamical system-based robot reaching motions," *Robotics Auton. Syst.*, vol. 62, pp. 752–765, 2014. [Online]. Available: https://api.semanticscholar.org/CorpusID:14374268
- [9] 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*. PMLR, 2018, pp. 439–473.
- [10] M. Simchowitz and D. Foster, "Naive exploration is optimal for online LQR," in *International Conference on Machine Learning*. PMLR, 2020, pp. 8937–8948.
- [11] R. Zhang, Y. Li, and N. Li, "On the regret analysis of online LQR control with predictions," in 2021 American Control Conference (ACC). IEEE, 2021, pp. 697–703.
- [12] I. Ziemann, A. Tsiamis, B. Lee, Y. Jedra, N. Matni, and G. J. Pappas, "A tutorial on the non-asymptotic theory of system identification," 2023.
- [13] L. Bako and H. Ohlsson, "Analysis of a nonsmooth optimization approach to robust estimation," *Automatica*, vol. 66, pp. 132–145, Apr. 2016.
- [14] H. Xu, C. Caramanis, and S. Mannor, "Robustness and Regularization of Support Vector Machines." *Journal of machine learning research*, vol. 10, no. 7, 2009.
- [15] L. Bako, "On a Class of Optimization-Based Robust Estimators," *IEEE Transactions on Automatic Control*, vol. 62, no. 11, pp. 5990–5997, Nov. 2017.
- [16] D. Bertsimas and M. S. Copenhaver, "Characterization of the equivalence of robustification and regularization in linear and matrix regression," *European Journal of Operational Research*, vol. 270, no. 3, pp. 931–942, Nov. 2018.
- [17] S. Pesme and N. Flammarion, "Online robust regression via SGD on the 1-1 loss," *Advances in Neural Information Processing Systems*, vol. 33, pp. 2540–2552, 2020.
- [18] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu, "On the sample complexity of the linear quadratic regulator," *Foundations of Computational Mathematics*, vol. 20, no. 4, pp. 633–679, 2020.
- [19] S. Mendelson, "Learning without concentration," *Journal of the ACM (JACM)*, vol. 62, no. 3, pp. 1–25, 2015.
- [20] Y. Li, S. Das, J. Shamma, and N. Li, "Safe adaptive learning-based control for constrained linear quadratic regulators with regret guarantees," *arXiv preprint arXiv:2111.00411*, 2021.
- [21] S. Fattahi, N. Matni, and S. Sojoudi, "Learning sparse dynamical systems from a single sample trajectory," in 2019 IEEE 58th Conference on Decision and Control (CDC). IEEE, 2019, pp. 2682–2689.
- [22] Y. Jedra and A. Proutiere, "Finite-time identification of stable linear systems optimality of the least-squares estimator," in 2020 59th IEEE Conference on Decision and Control (CDC). IEEE, 2020, pp. 996– 1001.
- [23] A. Wagenmaker and K. Jamieson, "Active learning for identification of linear dynamical systems," in *Conference on Learning Theory*. PMLR, 2020, pp. 3487–3582.
- [24] H. Feng, B. Yalcin, and J. Lavaei, "Learning of dynamical systems under adversarial attacks - null space property perspective," 2023 American Control Conference (ACC), pp. 4179–4184, 2023.
- [25] M. J. Wainwright, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, ser. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- [26] H. Feng and J. Lavaei, "Learning of dynamical systems under adversarial attacks," in 2021 60th IEEE Conference on Decision and Control (CDC), 2021, pp. 3010–3017.
- [27] Y. Chen, J. Fan, C. Ma, and Y. Yan, "Bridging convex and nonconvex optimization in robust pca: Noise, outliers, and missing data," *Annals of statistics*, vol. 49, no. 5, p. 2948, 2021.

- [28] J. Farkas, "Theorie der einfachen ungleichungen." Journal für die reine und angewandte Mathematik, vol. 124, pp. 1–27, 1902. [Online]. Available: http://eudml.org/doc/149129
- [29] R. Hovorka, F. Shojaee-Moradie, P. V. Carroll, L. J. Chassin, I. J. Gowrie, N. C. Jackson, R. S. Tudor, A. M. Umpleby, and R. H. Jones, "Partitioning glucose distribution/transport, disposal, and endogenous production during ivgtt," *American Journal of Physiology-Endocrinology and Metabolism*, vol. 282, no. 5, pp. E992–E1007, 2002.
- [30] I. Hajizadeh, M. Rashid, and A. Cinar, "Integrating compartment models with recursive system identification," in 2018 Annual American Control Conference (ACC), 2018, pp. 3583–3588.
- [31] R. Hovorka, V. Canonico, L. J. Chassin, U. Haueter, M. Massi-Benedetti, M. O. Federici, T. R. Pieber, H. C. Schaller, L. Schaupp, T. Vering *et al.*, "Nonlinear model predictive control of glucose concentration in subjects with type 1 diabetes," *Physiological measurement*, vol. 25, no. 4, p. 905, 2004.



Baturalp Yalcin, is a Ph.D.candidate at the University of California, Berkeley in Industrial Engineering and Operations Research. His research interests include the landscape of optimization and adversarial learning problems. He holds a B.S. in Industrial Engineering from Bogazici University and an M.S. in Industrial Engineering and Operations Research from the University of California, Berkeley.



Haixiang Zhang, is a Ph.D.candidate at the University of California, Berkeley in Applied Mathematics. His research interests include nonconvex optimization, especially low-rank matrix optimization and optimization via simulation. He holds B.S. in Computer Science and Technology and B.S. in Computational Mathematics from Peking University. He is the receipent of Two Sigma Ph.D. Fellowship.



Javad Lavaei, received the Ph.D. degree in control and dynamical systems from the California Institute of Technology, Pasadena, CA, in 2011. He is an Associate Professor in the Department of Industrial Engineering and Operations Research, the University of California, Berkeley, Berkeley, CA. Dr. Lavaei received multiple awards, including the NSF CAREER Award, the Office of Naval Research Young Investigator Award, and the Donald P. Eckman Award.



Murat Arcak, is currently a Professor with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA, USA. His research interests include dynamical systems and control theory with applications to synthetic biology, multi-agent systems, and transportation. Prof. Arcak was a recipient of the NSF CAREER Award and the Donald P. Eckman Award.