
Improved rates for prediction and identification for partially observed linear dynamical systems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Identification of a linear time-invariant dynamical system from partial observa-
2 tions is a fundamental problem in control theory. Particularly challenging are sys-
3 tems exhibiting long-term memory. A natural question is how learn such systems
4 with non-asymptotic statistical rates depending on the inherent dimensionality (or-
5 der) d of the system, rather than on the possibly much larger memory length. We
6 propose an algorithm that given a single trajectory of length T with gaussian ob-
7 servation noise, learns the system with a near-optimal rate of $\tilde{O}\left(\sqrt{\frac{d}{T}}\right)$ in \mathcal{H}_2 er-
8 ror, with only logarithmic, rather than polynomial dependence on memory length.
9 We also give bounds under process noise and improved bounds for learning a
10 realization of the system. Our algorithm is based on multi-scale low-rank approx-
11 imation: SVD applied to Hankel matrices of geometrically increasing sizes. Our
12 analysis relies on careful application of concentration bounds on the Fourier do-
13 main – we give sharper concentration bounds for sample covariance of correlated
14 inputs and for \mathcal{H}_∞ norm estimation, which may be of independent interest.

15 1 Introduction

16 We consider the problem of prediction and identification of an *unknown* partially-observed linear
17 time-invariant (LTI) dynamical system with stochastic noise,

$$x(t) = Ax(t-1) + Bu(t-1) + \xi(t) \quad (1)$$

$$y(t) = Cx(t) + Du(t) + \eta(t), \quad (2)$$

18 with a single trajectory of length T , given access only to input and output data. Here, $u(t) \in \mathbb{R}^{d_u}$ are
19 inputs, $x(t) \in \mathbb{R}^d$ are the hidden states, $y(t) \in \mathbb{R}^{d_y}$ are observations (or outputs), $\xi(t) \sim N(0, \Sigma_x)$
20 and $\eta(t) \sim N(0, \Sigma_y)$ are iid gaussian noise, and $A \in \mathbb{R}^{d \times d}$, $B \in \mathbb{R}^{d \times d_u}$, $C \in \mathbb{R}^{d_y \times d}$, $D \in \mathbb{R}^{d_y \times d_u}$
21 are matrices. Partial observability refers to the fact that we do not observe the state $x(t)$, but rather
22 a noisy linear observation $y(t)$.

23 As a simple and tractable family of dynamical systems, LTI systems are a central object of study for
24 control theory and time series analysis. The problem of prediction and filtering for a known system
25 dates back to [Kal60]. However, in many machine learning applications, the system is *unknown* and
26 must be learned from input and output data. Identification of an unknown system is often a necessary
27 first step for robust control [DMM⁺19, BMR18]. In a long line of recent work, the interplay between
28 machine learning and control theory has borne fruit in an improved understanding of the statistical
29 and online learning guarantees for prediction, identification, and control for unknown systems. In
30 machine learning, LTI systems also serve as a simple model problem for learning from correlated

31 data in stateful environments, and can give insight into understanding the successes of reinforcement
32 learning [Rec19, TR19] and recurrent neural networks [HMR18].

33 Partial observability poses a significant challenge to system identification: In the fully observed
34 setting, given access to $x(t)$, there is no obstacle to learning the matrices directly through linear
35 regression. However, in the partially observed setting, the most natural form of the optimization
36 problem is non-convex.

37 Systems exhibiting *long-term memory* are particularly challenging to learn. Restricting to marginally
38 stable systems, this occurs when the spectral radius of A , $\rho(A)$, is close to 1, and it implies that the
39 output at a particular time cannot be accurately estimated without taking into account inputs over
40 many previous time-steps—on the order of $O\left(\frac{1}{1-\rho(A)}\right)$ times steps. Such systems often arise in
41 practice. A particular class of such systems are those exhibiting *multiscale* behavior, with different
42 state variables that change on vastly different timescales [CR10]. For example, the body’s pH level is
43 affected both by long-term changes on a timescale of days or weeks, as well as breathing rate which
44 changes over a timescale of seconds. For such systems, it makes sense to discretize at the scale of the
45 fastest changing variable, which leads to a long memory for the slowest-changing variable. With few
46 exceptions, existing guarantees for learning partially observed LTI systems degrade as the memory
47 length increases. However, counting the number of parameters in the model (1)–(2) suggests that the
48 right measure of statistical complexity is the intrinsic dimensionality of the system, not the memory
49 length. This leads to the following natural question.

50 **Question:** How can we learn partially observed LTI systems with (non-asymptotic) statistical rates
51 that depend on the *intrinsic dimensionality* of the system, rather than the memory length?

52 Despite the simplicity of the question, little in the way of theoretical results are known. We focus
53 on the particular problem of learning the *impulse response (IR) function* of the system—which fully
54 determines its input-output behavior—in \mathcal{H}_2 norm. This is a natural norm for prediction problems
55 as it measures the expected prediction error under random input. Known guarantees for learning the
56 IR depend on the memory length. One particularly undesirable consequence is that for a continuous
57 system with time discretization Δ going to 0, the memory scales as $1/\Delta$ (while the system order
58 stays constant), leading to suboptimal estimation by an arbitrarily large factor.

59 Our key contribution is an algorithm and analysis that gives statistical rates that are optimal up to
60 logarithmic factors. Unlike previous works, our rates depend on the system order d —the natural
61 dimensionality of the problem—and only *logarithmically* on the memory length of the system. Our
62 algorithm is based on taking a low-rank approximation (SVD) of the Hankel matrix, which is a
63 widely used technique in system identification. We consider a *multiscale* version of this algorithm,
64 where we repeat this process for a geometric sequence of sizes of the Hankel matrix. This is essential
65 for obtaining a stronger theoretical guarantee. In the setting of zero process noise, we prove that our
66 algorithm achieves near-optimal $\tilde{O}\left(\sqrt{\frac{d(d_u+d_y)}{T}}\right)$ rates in \mathcal{H}_2 error for the learned system.

67 Our analysis relies on careful application of concentration bounds on the Fourier domain to give
68 sharper concentration bounds for sample covariance and \mathcal{H}_∞ norm estimation, which may be of
69 independent interest. While we consider our algorithm in a simple setting, we hope that this is a
70 first step to understanding and improving more complex subspace identification algorithms. Indeed,
71 SVD and related spectral methods are a standard step used in subspace identification algorithms
72 such as N4SID; our analysis suggests that SVD has an important “de-noising effect”.

73 We also give improved bounds for system identification, that is, learning the matrices A, B, C, D
74 using the Ho-Kalman algorithm [HK66], with $\tilde{O}\left(\sqrt{\frac{Ld(d_u+d_y)}{T}}\right)$ rates.

75 1.1 Notation

76 **Norms.** We use $\|\cdot\|$ to denote the 2-norm of a vector. For a matrix A , let $\|A\| = \|A\|_2$
77 denote its operator norm, $\rho(A)$ denote its spectral radius (maximum absolute value of eigen-

78 value), and $\|A\|_F$ denote its Frobenius norm. For a matrix-valued function $M(t) \in \mathbb{C}^{d_1 \times d_2}$,
 79 $\|M\|_F := \sqrt{\sum_t \|M(t)\|_F^2}$. Let $\sigma_r(A)$ denote the r th singular value of A .

80 **Fourier transform.** Given a matrix-valued function $F : \mathbb{Z} \rightarrow \mathbb{C}^{m \times n}$, define the (discrete-time)
 81 Fourier transform as the function $\widehat{F} : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{C}^{m \times n}$ given by $\widehat{F}(\omega) = \sum_{t=-\infty}^{\infty} F(t)e^{-2\pi i \omega t}$.

82 **Matrices.** Given a sequence $(F(t))_{t=1}^{a+b-1}$ where each $F(t) \in \mathbb{C}^{m \times n}$, define $\text{Hankel}_{a \times b}(F)$ as the
 83 $am \times bn$ block matrix such that the (i, j) th block is $[\text{Hankel}_{a \times b}(F)]_{ij} = F(i + j - 1)$. Given a
 84 sequence $(F(t))_{t=0}^{a-1}$ where each $F(t) \in \mathbb{C}^{m \times n}$, define the Toeplitz matrix as the block matrix such
 85 that the (i, j) th block is $[\text{Toep}_{a \times b}(F)]_{ij} = F(i - j) \mathbb{1}_{i \geq j}$. For a matrix A , let A^\top, A^H, A^\dagger denote
 86 its transpose, Hermitian (conjugate transpose), and pseudoinverse, respectively. For a vector-valued
 87 function $v : \{a, \dots, b\} \rightarrow \mathbb{R}^n$, let $v_{a:b} \in \mathbb{R}^{(|a-b|+1)n}$ denote the vertical concatenation of
 88 $v(a), \dots, v(b)$.

89 **Control theory.** For a matrix $A \in \mathbb{C}^{d \times d}$, define its resolvent as $\Phi_A(z) = (zI - A)^{-1}$. For a
 90 linear dynamical system \mathcal{D} given by (1)–(2), let $\Phi_{\mathcal{D}} = \Phi_{u \rightarrow y}$ denote the transfer function from
 91 u to y (response to input). Then $\Phi_{\mathcal{D}} = \Phi_{u \rightarrow y} = C\Phi_A B + D = C(zI - A)^{-1}B + D$. Let
 92 $\mathbb{T} := \{z \in \mathbb{C} : |z| = 1\}$ be the unit circle in the complex plane. For a matrix-valued function $F : \mathbb{T} \rightarrow \mathbb{C}^{d_1 \times d_2}$,
 93 define the \mathcal{H}_2 and \mathcal{H}_∞ norms by

$$\|\Phi\|_{\mathcal{H}_2} = \sqrt{\frac{1}{2\pi} \int_{\mathbb{T}} \|\Phi(z)\|_F^2 dz} \quad \|\Phi\|_{\mathcal{H}_\infty} = \sup_{z \in \mathbb{T}} \|\Phi(z)\|.$$

94 For a function $F : \mathbb{N}_0 \rightarrow \mathbb{C}^{d_1 \times d_2}$, define its Z-transform to be $\mathcal{Z}[F](z) = \sum_{n=0}^{\infty} F(n)z^{-n}$. Con-
 95 sidered as a function $\mathbb{T} \rightarrow \mathbb{C}$, we can take its \mathcal{H}_2 and \mathcal{H}_∞ norms. Overloading notation, we will let
 96 $\|F\|_{\mathcal{H}_p} := \|\mathcal{Z}F\|_{\mathcal{H}_p}$ for $p = 2, \infty$. The \mathcal{H}_2 and \mathcal{H}_∞ norms can be interpreted as the Frobenius and
 97 operator norms of the linear operator from input to output, i.e., they measure the average power of
 98 the output signal under random or worst-case input, respectively. For background on control theory,
 99 see e.g., [ZDG⁺96].

100 **Constants.** In proofs, C may represent different constants from line to line.

101 2 Main results

102 We consider the problem of prediction and identification for an unknown linear dynamical sys-
 103 tem (1)–(2). Our main goal is to obtain error guarantees in \mathcal{H}_2 norm, which determines prediction
 104 error under random input [OO19, Lemma 3.3].

105 **Problem 2.1.** Consider the partially-observed LTI system \mathcal{D} (1)–(2) with gaussian inputs $u(t) \sim$
 106 $N(0, I_{d_u})$ for $0 \leq t < T$. Suppose that the system is stable, that is, $\rho(A) < 1$, and that we observe
 107 a single trajectory of length T started with $x(0) = 0$, that is, we observe $u(t) \sim N(0, I_{d_u})$ and $y(t)$
 108 for $0 \leq t < T$.

109 The goal is to learn a LTI system $\widetilde{\mathcal{D}}$ with the aim of minimizing $\|\Phi_{\widetilde{\mathcal{D}}} - \Phi_{\mathcal{D}}\|_{\mathcal{H}_2}$. Equivalently, letting

$$F^*(t) = \begin{cases} D, & t = 0 \\ CA^{t-1}B, & t \geq 1 \end{cases}$$

110 denote the impulse response function (also called the Markov parameters) of the system, the goal is
 111 to learn an impulse response \widetilde{F} minimizing $\|F^* - \widetilde{F}\|_{\mathcal{H}_2} = \|F^* - \widetilde{F}\|_F$.

112 Note that learning F^* is sufficient to fully understand the input-output behavior of the system, but
 113 we may also ask to recover the system parameters A, B, C, D up to similarity transformation (see
 114 Theorem 2.3).

115 Previous results [OO19, SRD19] roughly depend polynomially on the “memory” $\frac{1}{1-\rho(A)}$, which
 116 blows up as the spectral norm of A approaches 1. In the setting of zero process noise, our goal is

117 to obtain rates that are $\tilde{O}\left(\frac{\text{poly}(d, d_u, d_y)}{\sqrt{T}}\right)$, with only poly-logarithmic dependence on $\frac{1}{1-\rho(A)}$. See
 118 Figure 1 for a comparison.

119 We assume that $\rho(A) < 1$ because if \mathcal{D} is not stable, it is in general impossible to learn $\tilde{\mathcal{D}}$ with
 120 finite \mathcal{H}_2 error, as a system with infinite response can have arbitrarily small response on any finite
 121 time interval. However, it may still be possible to learn the response up to time $L \ll T$ in this
 122 case [SBR19]. The marginally stable case ($\rho(A) = 1$) is an important case we leave to future work.

Method	Rollout type	Min # samples	IR error
Least squares (IR) [TBPR17]	Multi	L	$\sigma\sqrt{\frac{L}{T}}$
Least squares (IR) [OO19]	Single	L	$\sigma\sqrt{\frac{L}{T}}$
Nuclear norm minimization [SOF20]	Multi	$\min\{d^2, L\}$	$\sigma\sqrt{\frac{L}{T}}$
	Multi	d	$\sigma\sqrt{\frac{dL}{T}}$
rank- d SVD (Theorem 2.2)	Single	L	$\sigma\sqrt{\frac{d}{T}}$

Figure 1: Here, L is the memory length for the system, which is $\tilde{O}\left(\frac{1}{1-\rho(A)}\right)$ for well-conditioned systems. *Rollout type* refers to whether we have access to a single trajectory or multiple trajectories. *Min # samples* refers to the minimum number of samples (up to log factors) before the bounds are operational. *IR error* refers to the error in the impulse response in Frobenius/ \mathcal{H}_2 norm. Logarithmic factors are omitted.

123 In our Algorithm 1, we first use linear regression to obtain a noisy estimate F of the impulse re-
 124 sponse. Next, following standard system identification procedures, we form the Hankel matrix
 125 $\text{Hankel}_{L \times L}(F)$ with the entries of F on its diagonals. Because the true Hankel matrix

$$\text{Hankel}_{L \times L}(F^*) = \begin{pmatrix} CB & CAB & \dots & CA^{L-1}B \\ CAB & CA^2B & & \vdots \\ \vdots & & \ddots & \vdots \\ CA^{L-1}B & \dots & \dots & CA^{2L-1}B \end{pmatrix}$$

126 has rank d , we take a low-rank SVD R_L of the Hankel matrix to “de-noise” the impulse response.
 127 We can then read off the estimated impulse response by averaging over the corresponding diagonal
 128 of R_L . For technical reasons, we need to repeat this process for a geometric sequence of sizes
 129 of the Hankel matrix: $L \times L$, $L/2 \times L/2$, $L/4 \times L/4$, and so forth. This is because the low-
 130 rank approximation objective for a $\ell \times \ell$ Hankel matrix encourages the diagonals that are $\Theta(\ell)$
 131 to be close—as those are the diagonals with the most entries—and hence estimates $F^*(t)$ well when
 132 $t = \Theta(\ell)$. In other words, low-rank estimation for $\text{Hankel}_{\ell \times \ell}(F)$ is only sensitive to the portion of
 133 the signal that is at timescale ℓ . Repeating this process ensures that we cover all scales.

134 Our main theorem is the following.

135 **Theorem 2.2.** *There is a constant C_1 such that following holds. In the setting of Problem 2.1,*
 136 *suppose that F^* is the impulse response function, T is such that $T \geq C_1 L d_u \log\left(\frac{L d_u}{\delta}\right)$, $\varepsilon_{\text{trunc}} :=$
 137 $\|F^* \mathbb{1}_{[2L, \infty)}\|_{\mathcal{H}_\infty} \sqrt{d_u} + \|G^* \mathbb{1}_{[2L, \infty)}\|_{\mathcal{H}_\infty} \left\| \Sigma_x^{1/2} \right\|_{\text{F}}$, and $M_{x \rightarrow y} = (O, C, CA, \dots, CA^{L-1})^\top \in$
 138 $\mathbb{R}^{(L+1)d \times d_y}$. Let $0 < \delta \leq \frac{1}{2}$ and $\sigma = \sqrt{\|\Sigma_y\| + \|\Sigma_x\| L \log\left(\frac{L d_u}{\delta}\right) \|M_{x \rightarrow y}\|^2}$. Then with probabil-
 139 ity at least $1 - \delta$, Algorithm 1 learns an impulse response function \tilde{F} such that*

$$\left\| \tilde{F} - F^* \right\|_{\text{F}} = O\left(\sigma \sqrt{\frac{d(d_y + d_u + \log\left(\frac{L}{\delta}\right)) \log L}{T}} + \varepsilon_{\text{trunc}} \sqrt{d} + \|F^* \mathbb{1}_{(L, \infty)}\|_{\text{F}} \right).$$

140 In the absence of process noise (when $\Sigma_x = O$), when L and T are chosen large enough, the first
 141 term dominates, and ignoring log factors, the dependence is $O\left(\sqrt{\frac{d(d_y + d_u)}{T}}\right)$. We expect this to be

Algorithm 1 Learning impulse response through multi-scale low-rank Hankel SVD

- 1: **Input:** Length L (power of 2), time T .
- 2: Part 1: Linear regression to recover noisy impulse response
- 3: Let $u(t) \sim N(0, I_{d_u})$ for $0 \leq t < T$, and observe the outputs $y(t) \in \mathbb{R}^{d_y}$, $0 \leq t < T$.
- 4: Solve the least squares problem

$$\min_{F: \text{Supp}(F) \subseteq [0, 2L-1]} \sum_{t=0}^{T-1} \|y(t) - F * u(t)\|^2.$$

- to obtain the noisy impulse response $F : [0, 2L - 1] \cap \mathbb{Z} \rightarrow \mathbb{R}^{d_y \times d_u}$.
 - 5: Part 2: Low-rank Hankel SVD to de-noise impulse response
 - 6: Let $\widetilde{F}(0) = F(0)$.
 - 7: **for** $k = 0$ to $\log_2 L$ **do**
 - 8: Let $\ell = 2^k$.
 - 9: Let R_ℓ be the rank- d SVD of $\text{Hankel}_{\ell \times \ell}(F)$ (i.e., $\text{argmin}_{\text{rank}(R) \leq d} \|R - \text{Hankel}_{\ell \times \ell}(F)\|$).
 - 10: For $\frac{\ell}{2} < t \leq \ell$, let $\widetilde{F}(t)$ be the $d_y \times d_u$ matrix given by $\widetilde{F}(t) = \frac{1}{t} \sum_{i+j=t} (R_\ell)_{ij}$, where $(\cdot)_{ij}$ denotes the (i, j) th block of the matrix.
 - 11: **end for**
 - 12: **Output:** Estimate of impulse response \widetilde{F} .
-

142 the optimal sample complexity up to logarithmic factors. However, in the presence of process noise,
 143 there is an undesirable factor of $\sqrt{L} \|M_{x \rightarrow y}\|$, which (for well-conditioned matrices) is expected to
 144 be $O\left(\frac{1}{1-\rho(A)}\right)$ or $O(L)$. We leave it an open problem to improve the guarantees in this setting.

145 *Remark 1.* The L -factor dependence on the process noise is unavoidable with the current algorithm:
 146 when the process noise has covariance $\Sigma_y = I$ and decays after L steps, it can cause perturbations of
 147 size $O(\sqrt{L})$ compared to the noiseless system. Even in the case $d = 1$, when the impulse response
 148 function is $ae^{-kt/L}$ for a known k , the noise will cause the estimate of a to be off by $O(\sqrt{L})$, and
 149 hence the \mathcal{H}_2 norm of the impulse response to be off by $O(L)$. Our algorithm only regresses on
 150 previous inputs, but in the presence of process noise, a better approach is to regress on both the
 151 previous inputs $u(t)$ and outputs $y(t)$ and then take a (weighted) SVD, as in N4SID [Qin06].

152 *Remark 2.* A burn-in time of $\Omega(L)$ is information-theoretically required to get poly(d) rates. At-
 153 tempting to extrapolate an impulse response function from time $o(L)$ to time L can magnify errors
 154 by $\exp(d)$, because the finite impulse response of a system of order d can approximate a polynomial
 155 of degree $d - 1$ on $[0, L]$.

156 We also show the following improved rates for learning the system matrices, by combining \mathcal{H}_∞
 157 bounds for the learned impulse response with stability results for the Ho-Kalman algorithm [OO19].
 158 Because the input-output behavior is unchanged under a similarity transformation $(A, B, C) \leftrightarrow$
 159 $(W^{-1}AW, W^{-1}B, CW)$, we can only learn the parameters up to similarity transformation.

160 **Theorem 2.3.** *Keep the assumptions and notation of Theorem 2.2, suppose \mathcal{D} is observable and*
 161 *controllable, and let*

$$\varepsilon' = \sigma \sqrt{\frac{L(d_y + d_u + \log(\frac{L}{\delta}))}{T}} + \varepsilon_{\text{trunc}}.$$

162 Let $H^- = \text{Hankel}_{L \times (L-1)}(F^*)$. Suppose that $\varepsilon' = O(\sigma_{\min}(H^-))$. Then with probability at least
 163 $1 - \delta$, the Ho-Kalman algorithm (Algorithm 1 in [OO19] with $T_1 = L, T_2 = L - 1$) returns $\widehat{A}, \widehat{B}, \widehat{C}$
 164 such that there exists a unitary matrix W satisfying

$$\max \left\{ \|C - \widehat{C}W\|_{\text{F}}, \|B - W^{-1}\widehat{B}\|_{\text{F}} \right\} = O(\sqrt{d} \cdot \varepsilon')$$

$$\|A - W^{-1}\widehat{A}W\|_{\text{F}} = O\left(\frac{1}{\sigma_{\min}(H^-)} \cdot \sqrt{d} \cdot \varepsilon' \cdot \left(\frac{\|\Phi_{\mathcal{D}}\|_{\mathcal{H}_\infty}}{\sigma_{\min}(H^-)} + 1\right)\right).$$

165 As L can be chosen to make $\varepsilon_{\text{trunc}}$ negligible, this gives $\tilde{O}\left(\sqrt{\frac{Ld(d_u+d_y)}{T}}\right)$ rates, however,
 166 with factors depending on the minimum eigenvalue of H . This is an improvement over the
 167 $\tilde{O}\left(\sqrt{d}^4\sqrt{\frac{L(d_u+d_y)}{T}}\right)$ rates in [OO19].

168 We prove Theorem 2.2 in Section 4 and Theorem 2.3 in Appendix B.

169 3 Related work

170 We survey two classes of methods for learning partially observable LDS's, subspace identification
 171 and improper learning. With the exception of [RJR20], all guarantees have sample complexity
 172 depending on the memory length L , which we wish to avoid.

173 3.1 Subspace identification

174 The basic idea of subspace identification [Lju98, Qin06, VODM12] is to learn a certain structured
 175 matrix (such as a Hankel matrix), take a best rank- k approximation (using SVD or another linear di-
 176 mensionality reduction method), and learn the system matrices A, B, C, D up to similarity transfor-
 177 mation. Usage of spectral methods circumvents the fact that the most natural optimization problem
 178 for A, B, C, D is non-convex. However, classical guarantees for these methods are asymptotic.

179 Recently, various authors have given non-asymptotic guarantees for system identification algo-
 180 rithms. [OO19] analyzed the Ho-Kalman algorithm [HK66] in this setting. [SRD19] consider the
 181 setting where system order is unknown and give an end-to-end result for prediction, while [TMP20]
 182 consider the problem of online filtering, that is, recovering $x(t)$'s up to some linear transformation.

183 An alternate, empirically successful approach is that of nuclear norm minimization or regulariza-
 184 tion [FPST13]. [SOF20] (building on [CQXY16]) give explicit rates of convergence, and show that
 185 the algorithm has a lower minimum sample complexity and is easier to tune.

186 Our algorithm is based on the classical approach of taking a low-rank approximation of the Hankel
 187 matrix, but we repeat this process with Hankel matrices of sizes $L \times L, L/2 \times L/2, L/4 \times L/4$, and so
 188 forth; this is key modification that allows us to obtain better statistical rates. Our analysis builds on
 189 the analyses given in [OO19, SOF20]. As essential part of the analysis is analyzing linear regression
 190 for correlated inputs, where we extend the work of [DMR19] to MIMO systems, as explained below.

191 3.1.1 Linear regression with correlated inputs

192 An important step in obtaining non-asymptotic rates for system identification is analyzing linear
 193 regression for correlated inputs. The most challenging step is to lower-bound the sample covariance
 194 matrix of inputs to the linear regression. A lower bound, rather than a matrix concentration result, is
 195 sufficient [Men14, SMT⁺18, MT19]; however, a concentration result is obtainable in our setting.

196 [TBPR17] give non-asymptotic bounds for learning the finite impulse response for a SISO system
 197 in ℓ^∞ Fourier norm; however, they require L rollouts of size $O(L)$ and hence $\Omega(L^2)$ timesteps.
 198 Addressing the more challenging single-rollout setting, [OO19] obtain bounds for a single rollout
 199 of $\tilde{\Omega}(L)$ timesteps, by using concentration bounds for random circulant matrices [KMR14] to de-
 200 rive concentration inequalities for the covariance matrix. These concentration inequalities for the
 201 covariance matrix were improved (by logarithmic factors) by [DMR19]. Although [DMR19] give
 202 an analysis in the SISO case, as we show in Theorem A.2, the results can be extended to the MIMO
 203 case with an ε -net argument.

204 3.2 Improper learning using autoregressive methods

205 Instead of solving the statistical problem of identifying parameters, another line of work develops
 206 algorithms for regret minimization in online learning. The goal is simply to do well in predicting
 207 future observations, with small loss (regret) compared to the best predictor in hindsight; the learned
 208 predictor is allowed to be improper, that is, take a different functional form. In the stochastic case,

209 this allows prediction almost as well as if the actual system parameters were known; however, the
 210 framework also allows for adversarial noise.

211 One popular strategy for improperly learning the system is to learn a linear autoregressive filter over
 212 previous inputs and observations, or ARMA model. Naturally, because we are optimizing over a
 213 larger hypothesis class, the statistical rates depend on L rather than the system order d .

214 [GLS⁺20, Theorem 4.7] consider the problem of online prediction for a fully or partially observed
 215 LDS, and give a regret bound that depends polynomially on the memory length L . Their approach
 216 works even for marginally stable systems, that is, systems with $\rho(A) \leq 1$. See also [AHMS13,
 217 HSZ17, HLS⁺18, KMTM19, TP20, RJR20] for previous work using autoregressive methods.

218 Of particular interest to us is [RJR20], which gives rates independent of spectral radius. Building
 219 on [HSZ17], they observe that it suffices to regress on previous inputs and outputs projected to a
 220 lower-dimensional space. Their algorithm works in the setting of process noise and competes with
 221 the Kalman filter, but only when $A - KC$ has real eigenvalues, where K is the Kalman gain.

222 4 Proof of main theorem

223 In this section, we prove Theorem 2.2. The proof hinges on the following lemma, which shows
 224 that if we observe a low-rank matrix plus noise, then taking a low-rank SVD can have a de-noising
 225 effect, producing a matrix that is closer to the true matrix.

226 **Lemma 4.1** (De-noising effect of SVD). *There exists a constant C such that the following holds.*
 227 *Suppose that $A \in \mathbb{C}^{m \times n}$ is a rank- r matrix, $\hat{A} = A + E$, and \hat{A}_r is the rank- r SVD of \hat{A} . Then*

$$\left\| \hat{A}_r - A \right\|_{\text{F}} \leq C\sqrt{r} \|E\|. \quad (3)$$

228 Compare this with the original error $\left\| \hat{A} - A \right\|_{\text{F}} = \|E\|_{\text{F}}$, which can only be bounded by
 229 $\sqrt{\min\{m, n\}} \|E\|$. When applied to the d -SVD of the Hankel matrix, this gives a factor of \sqrt{d}
 230 rather than \sqrt{L} for the error.

231 *Proof.* We have

$$\left\| \hat{A}_r - A \right\|_{\text{F}} \leq \sqrt{2r} \left\| \hat{A}_r - A \right\|_2 \quad (4)$$

$$\leq \sqrt{2r} \left(\left\| \hat{A}_r - \hat{A} \right\|_2 + \left\| \hat{A} - A \right\|_2 \right) \quad (5)$$

$$\leq 2\sqrt{2r} \|E\| \quad (6)$$

232 where (4) follows from $\hat{A}_r - A$ having rank at most $2r$, (5) follows from the triangle inequality,
 233 and (6) follows from Weyl's Theorem: $\left\| \hat{A}_r - \hat{A} \right\|_2 \leq \sigma_{r+1}(\hat{A}) \leq \sigma_{r+1}(A) + \|E\| = \|E\|$. \square

234 To prove Theorem 2.2, we will need to obtain bounds for $F : \{0, 1, \dots, 2L - 1\} \rightarrow \mathbb{R}^{d_y \times d_u}$ learned
 235 from linear regression in \mathcal{H}_∞ norm. The following is our main technical result.

236 **Lemma 4.2.** *There are C_1, C_2 such that the following hold. Suppose $y = F^* * u + G^* * \xi + \eta$ where $u(t) \sim N(0, I_{d_u})$, $\xi(t) \sim N(0, \Sigma_x)$, $\eta(t) \sim N(0, \Sigma_y)$ for $0 \leq t < T$,
 237 and $\text{Supp}(F^*), \text{Supp}(G^*) \subseteq [0, \infty)$. Let $F = \text{argmin}_{F \in \{0, \dots, L\} \rightarrow \mathbb{R}^{d_y \times d_u}} \sum_{t=0}^{T-1} |y(t) - (F * u)(t)|^2$, $M_{G^*} = (G^*(0), \dots, G^*(L))^\top \in \mathbb{R}^{(L+1)d \times d_y}$, and $\varepsilon_{\text{trunc}} = \|F^* \mathbb{1}_{[L+1, \infty)}\|_{\mathcal{H}_\infty} \sqrt{d_u} +$
 239 $\|G^* \mathbb{1}_{[L+1, \infty)}\|_{\mathcal{H}_\infty} \left\| \Sigma_x^{1/2} \right\|_{\text{F}}$. For $0 < \delta \leq \frac{1}{2}$, $T \geq C_1 L d_u \log\left(\frac{L d_u}{\delta}\right)$, $-1 \leq a < L - L'$,*

$$\begin{aligned} & \left\| (F - F^*) \mathbb{1}_{[a+1, a+L']} \right\|_{\mathcal{H}_\infty} \\ & \leq C_2 \left[\sqrt{\frac{1}{T}} \left(\sqrt{\|\Sigma_y\| L' \left(d_u + d_y + \log\left(\frac{L'}{\delta}\right) \right)} + \sqrt{\|\Sigma_x\| L' L d_u \log\left(\frac{L d_u}{\delta}\right)} \|M_{G^*}\| \right) + \varepsilon_{\text{trunc}} \right] \end{aligned}$$

241 *with probability at least $1 - \delta$.*

242 In the case $\Sigma_x = O$, this roughly says that the error in the learned impulse response, $F - F^*$, over
 243 any interval of length L' , has all Fourier coefficients bounded in spectral norm by $\tilde{O}\left(\sqrt{\frac{L'(d_u+d_y)}{T}}\right)$
 244 – what we expect if the error from linear regression is uniformly distributed over all frequencies.

245 A complete proof is in Appendix A; we give a brief sketch. First, because the errors are Gaussian,
 246 the error from linear regression, $F - F^*$, follows a Gaussian distribution. To bound its covariance,
 247 we lower-bound the smallest singular value of the sample covariance of the inputs (Lemma A.1,
 248 Appendix A.1). Here, the difficulty is that the inputs are *correlated* – the input at time t is $u_{t:t-L}$.
 249 Fortunately, the translation structure means it is close to a submatrix of an infinite block Toeplitz
 250 matrix, which becomes block diagonal in the Fourier domain. This “decoupling” allows us to show
 251 concentration. Compared to the SISO setting in [DMR19], we require an extra ε -net argument. Once
 252 we have a bound on the covariance, we can bound any $\left\|\widehat{(F - F^*)}(\omega)\right\|$ by matrix concentration
 253 (Appendix A.2); to bound the \mathcal{H}_∞ norm it suffices to bound this over a grid of ω 's (Lemma A.4).

254 Bounding the error in \mathcal{H}_∞ norm of the impulse response allows us to bound the error in operator
 255 norm of the Hankel matrix, as the following lemma shows.

256 **Lemma 4.3.** *For any $F : \mathbb{Z} \rightarrow \mathbb{C}^{m \times n}$, we have $\|\text{Hankel}_{a \times b}(F)\| \leq \|F\|_{\mathcal{H}_\infty}$.*

257 *Proof.* Note that when $v : \mathbb{Z} \rightarrow \mathbb{C}^n$, $\text{Supp}(v) \subseteq [0, b-1]$, we have $(F * v)_{b:b+a-1} =$
 258 $\text{Hankel}_{a \times b}(F)v_{b-1:0}$. Hence for any v set to 0 outside of $[0, b-1]$, using Parseval's Theorem
 259 and the fact that the Fourier transform of a convolution is the product of the Fourier transforms, we
 260 have

$$\|\text{Hankel}_{a \times b}(F)v_{b-1:0}\|_2 \leq \|F * v\|_2 = \|\widehat{F}\widehat{v}\|_2 \leq \sup_{\omega \in [0,1]} \|\widehat{F}(\omega)\|_2 \|\widehat{v}\|_2 = \|F\|_{\mathcal{H}_\infty} \|v\|_2$$

261 This shows that $\|\text{Hankel}_{a \times b}(F)\| \leq \|F\|_{\mathcal{H}_\infty}$. □

262 Theorem 2.2 will follow from the following bound after an application of the triangle inequality.

263 **Lemma 4.4.** *There are C_1, C_2 such that the following holds for the setting of Problem 2.1.*
 264 *Suppose L is a power of 2, and $T \geq C_1 L d_u \log\left(\frac{L d_u}{\delta}\right)$. Let $\|F^* \mathbb{1}_{[L+1, \infty)}\|_{\mathcal{H}_\infty} \sqrt{d_u} +$
 265 $\|G^* \mathbb{1}_{[L+1, \infty)}\|_{\mathcal{H}_\infty} \left\|\Sigma_x^{1/2}\right\|_{\text{F}}$ and $M_{x \rightarrow y} = (O, C, CA, \dots, CA^{L-1})^\top \in \mathbb{R}^{(L+1)d \times d_y}$. Then with
 266 probability at least $1 - \delta$, the output \widetilde{F} given by Algorithm 1 satisfies*

$$\left\|\left(\widetilde{F} - F^*\right)\mathbb{1}_{[1,L]}\right\|_{\text{F}} \leq C_2 \left(\sqrt{\frac{\|\Sigma_y\| d (d_y + d_u + \log\left(\frac{L}{\delta}\right)) \log L}{T}} + \sqrt{\frac{\|\Sigma_x\| L d d_u \log\left(\frac{L d_u}{\delta}\right)}{T}} \|M_{G^*}\| + \varepsilon_{\text{trunc}} \sqrt{d} \right)$$

267 *Proof.* We are in the situation of Lemma 4.2 with $G^*(t) = CA^{t-1} \mathbb{1}_{t \geq 1}$. Let $\mathcal{H}_\ell = \text{Hankel}_{\ell \times \ell}(F)$
 268 and $\mathcal{H}_\ell^* = \text{Hankel}_{\ell \times \ell}(F^*)$. Suppose $\ell \leq L$ is even. Note that

$$\mathcal{H}_\ell = \underbrace{\text{Hankel}_{\ell \times \ell}(F^*)}_{\mathcal{H}_\ell^*} + \text{Hankel}_{\ell \times \ell}(F - F^*)$$

269 where $\mathcal{H}_\ell^* = \text{Hankel}_{\ell \times \ell}(F^*)$ is a rank- d matrix, with error term is bounded by

$$\begin{aligned} \|\text{Hankel}_{\ell \times \ell}(F - F^*)\| &\leq \|(F - F_{\text{trunc}}^*) \mathbb{1}_{[1, 2\ell-1]}\|_{\mathcal{H}_\infty} && \text{by Lemma 4.3} \\ &< C \left[\sqrt{\frac{1}{T}} \left(\sqrt{\|\Sigma_y\| \ell \left(d_u + d_y + \log\left(\frac{L'}{\delta}\right) \right)} \right. \right. \\ &\quad \left. \left. + \sqrt{\|\Sigma_x\| \ell L d_u \log\left(\frac{L d_u}{\delta}\right)} \|M_{G^*}\| + \varepsilon_{\text{trunc}} \right) \right] && \text{by Lemma 4.2} \end{aligned} \quad (7)$$

270 with probability at least $1 - \delta$. Let R_ℓ be the rank- d SVD of \mathcal{H}_ℓ . Then by Lemma 4.1,

$$\|R_\ell - \mathcal{H}_\ell^*\|_{\text{F}} = O\left(\sqrt{d} \|\text{Hankel}_{\ell \times \ell}(F - F^*)\|\right). \quad (8)$$

271 Now letting $\tilde{F}(t) = \frac{1}{t} \sum_{i+j=t} (R_\ell)_{ij}$ when $\frac{\ell}{2} < t \leq \ell$ we have (using the fact that the mean
 272 minimizes the sum of squared errors)

$$\begin{aligned} \|R_\ell - \mathcal{H}_\ell^*\|_F^2 &\geq \sum_{t=\frac{\ell}{2}+1}^{\ell} \sum_{i+j=t} \|(R_\ell)_{ij} - F^*(t)\|_F^2 \\ &\geq \sum_{t=\frac{\ell}{2}+1}^{\ell} \left(t \cdot \|\tilde{F}(t) - F^*(t)\|_F^2 \right) \geq \left(\frac{\ell}{2} + 1 \right) \sum_{t=\frac{\ell}{2}+1}^{\ell} \left(\|\tilde{F}(t) - F^*(t)\|_F^2 \right). \end{aligned}$$

273 Note that we only get a lower bound with a factor of ℓ if we restrict to t that is $\Theta(\ell)$, i.e., restrict to
 274 diagonals that have many entries. This is the reason we will have to repeat this process for multiple
 275 sizes. Hence

$$\left\| (\tilde{F} - F^*) \mathbb{1}_{[\frac{\ell}{2}+1, \ell]} \right\|_F^2 \leq \frac{1}{\ell/2} \|R_\ell - \mathcal{H}_\ell^*\|_F^2.$$

276 Together with (8) and (7) this gives with probability $\geq 1 - \delta$ that

$$\left\| (\tilde{F} - F^*) \mathbb{1}_{[\frac{\ell}{2}+1, \ell]} \right\|_F \leq C \left(\sqrt{\frac{\|\Sigma_y\| d (d_y + d_u + \log(\frac{L}{\delta}))}{T}} + \sqrt{\frac{\|\Sigma_x\| L d_u \log(\frac{L d d_u}{\delta})}{T}} \|M_{G^*}\| + \frac{\varepsilon_{\text{trunc}} \sqrt{d}}{\sqrt{\ell}} \right)$$

277 Replacing δ by $\frac{\delta}{\log_2 L}$, using a union bound over powers of 2, and summing gives

$$\left\| (\tilde{F} - F^*) \mathbb{1}_{[1, L]} \right\|_F = O \left(\sqrt{\frac{\|\Sigma_y\| d (d_y + d_u + \log(\frac{L}{\delta})) \log L}{T}} + \sqrt{\frac{\|\Sigma_x\| L d d_u \log(\frac{L d_u}{\delta}) \log L}{T}} \|M_{G^*}\| + \varepsilon_{\text{trunc}} \sqrt{d} \right).$$

278 □

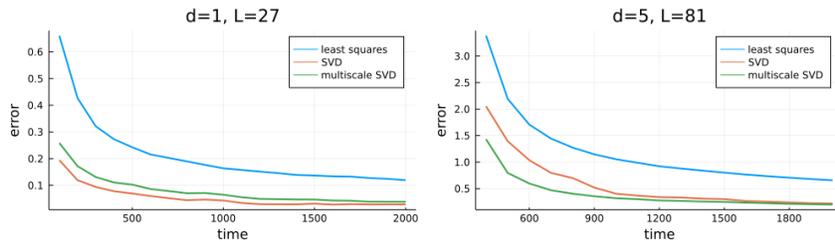
279 *Proof of Theorem 2.2.* We have the bound in Lemma 4.4, and also the same bound
 280 for $\left\| (\tilde{F} - F^*)(0) \right\|_F$ after applying Lemma 4.2 to $(F - F^*)\delta_0$. Finally, note that
 281 $\left\| (\tilde{F} - F^*) \mathbb{1}_{(L, \infty)} \right\|_F = \|F^* \mathbb{1}_{(L, \infty)}\|_F$ and use the triangle inequality. □

282 5 Experiments

283 We compared three algorithms for learning the impulse response function: least-squares, and low-
 284 rank Hankel SVD with and without the multi-scale repetition. We include details of the experimental
 285 setup in Appendix D. Note that to reduce the number of scales, we consider use a slight modification
 286 of our Algorithm 1 which triples the size at each iteration instead.

287 The plots show the error $\|F^* \mathbb{1}_{[1, L]} - F\|_2$, where F is the estimated impulse response on $[1, L]$,
 288 averaged over 10 randomly generated LDS's, as a function of the time T elapsed. We consider
 289 systems of order $d = 1, 3, 5, 10$, and memory lengths $L = 27, 81$.

290 Using SVD significantly reduces the error, supporting our theory which shows that SVD has a “de-
 291 noising” effect. Additionally, multiscale SVD has better performance than naive SVD when d is
 292 moderate, L is large, and data is limited, but the performance is similar in a data-rich setting.



293 **References**

- 294 [AHMS13] Oren Anava, Elad Hazan, Shie Mannor, and Ohad Shamir. Online learning for time
295 series prediction. In *COLT 2013 - The 26th Annual Conference on Learning Theory,*
296 *June 12-14, 2013, Princeton University, NJ, USA*, pages 172–184, 2013.
- 297 [BMR18] Ross Boczar, Nikolai Matni, and Benjamin Recht. Finite-data performance guarantees
298 for the output-feedback control of an unknown system. In *2018 IEEE Conference on*
299 *Decision and Control (CDC)*, pages 2994–2999. IEEE, 2018.
- 300 [BTR13] Badri Narayan Bhaskar, Gongguo Tang, and Benjamin Recht. Atomic norm denoising
301 with applications to line spectral estimation. *IEEE Transactions on Signal Processing*,
302 61(23):5987–5999, 2013.
- 303 [CQXY16] Jian-Feng Cai, Xiaobo Qu, Weiyu Xu, and Gui-Bo Ye. Robust recovery of complex
304 exponential signals from random gaussian projections via low rank hankel matrix re-
305 construction. *Applied and computational harmonic analysis*, 41(2):470–490, 2016.
- 306 [CR10] Shaunak Chatterjee and Stuart Russell. Why are dbns sparse? In *Proceedings of*
307 *the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages
308 81–88. JMLR Workshop and Conference Proceedings, 2010.
- 309 [DMM⁺19] Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the
310 sample complexity of the linear quadratic regulator. *Foundations of Computational*
311 *Mathematics*, pages 1–47, 2019.
- 312 [DMR19] Boualem Djehiche, Othmane Mazhar, and Cristian R Rojas. Finite impulse response
313 models: A non-asymptotic analysis of the least squares estimator. *arXiv preprint*
314 *arXiv:1911.12794*, 2019.
- 315 [FPST13] Maryam Fazel, Ting Kei Pong, Defeng Sun, and Paul Tseng. Hankel matrix rank
316 minimization with applications to system identification and realization. *SIAM Journal*
317 *on Matrix Analysis and Applications*, 34(3):946–977, 2013.
- 318 [GLS⁺20] Udaya Ghai, Holden Lee, Karan Singh, Cyril Zhang, and Yi Zhang. No-regret predic-
319 tion in marginally stable systems. In *COLT 2020 - The 33rd Annual Conference on*
320 *Learning Theory, July 9-12, 2020*, pages 1–44, 2020.
- 321 [HK66] BL Ho and Rudolph E Kalman. Effective construction of linear state-variable models
322 from input/output functions. *at-Automatisierungstechnik*, 14(1-12):545–548, 1966.
- 323 [HLS⁺18] Elad Hazan, Holden Lee, Karan Singh, Cyril Zhang, and Yi Zhang. Spectral filtering
324 for general linear dynamical systems. In *Advances in Neural Information Processing*
325 *Systems*, pages 4634–4643, 2018.
- 326 [HMR18] Moritz Hardt, Tengyu Ma, and Benjamin Recht. Gradient descent learns linear dynam-
327 ical systems. *The Journal of Machine Learning Research*, 19(1):1025–1068, 2018.
- 328 [HSZ17] Elad Hazan, Karan Singh, and Cyril Zhang. Learning linear dynamical systems via
329 spectral filtering. In *Advances in Neural Information Processing Systems*, pages 1–2,
330 2017.
- 331 [Kal60] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems.
332 *Journal of Basic Engineering*, 82.1:35–45, 1960.
- 333 [KMR14] Felix Krahmer, Shahar Mendelson, and Holger Rauhut. Suprema of chaos processes
334 and the restricted isometry property. *Communications on Pure and Applied Mathemat-*
335 *ics*, 67(11):1877–1904, 2014.

- 336 [KMTM19] Mark Kozdoba, Jakub Marecek, Tigran T. Tchrakian, and Shie Mannor. On-line learn-
337 ing of linear dynamical systems: Exponential forgetting in kalman filters. In *Proceed-*
338 *ings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4098–4105,
339 2019.
- 340 [Lju98] Lennart Ljung. *System identification: Theory for the User*. Prentice Hall, Upper Saddle
341 River, NJ, 2 edition, 1998.
- 342 [Men14] Shahar Mendelson. Learning without concentration. In *Conference on Learning The-*
343 *ory*, pages 25–39, 2014.
- 344 [MT19] Nikolai Matni and Stephen Tu. A tutorial on concentration bounds for system identifi-
345 cation. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 3741–
346 3749. IEEE, 2019.
- 347 [OO19] Samet Oymak and Necmiye Ozay. Non-asymptotic identification of lti systems from
348 a single trajectory. In *2019 American Control Conference (ACC)*, pages 5655–5661.
349 IEEE, 2019.
- 350 [Qin06] S Joe Qin. An overview of subspace identification. *Computers & chemical engineer-*
351 *ing*, 30(10-12):1502–1513, 2006.
- 352 [Rec19] Benjamin Recht. A tour of reinforcement learning: The view from continuous control.
353 *Annual Review of Control, Robotics, and Autonomous Systems*, 2:253–279, 2019.
- 354 [RJR20] Paria Rashidinejad, Jiantao Jiao, and Stuart Russell. Slip: Learning to predict in un-
355 known dynamical systems with long-term memory. *arXiv preprint arXiv:2010.05899*,
356 2020.
- 357 [RV⁺13] Mark Rudelson, Roman Vershynin, et al. Hanson-wright inequality and sub-gaussian
358 concentration. *Electronic Communications in Probability*, 18, 2013.
- 359 [SBR19] Max Simchowitz, Ross Boczar, and Benjamin Recht. Learning linear dynamical sys-
360 tems with semi-parametric least squares. *arXiv preprint arXiv:1902.00768*, 2019.
- 361 [SGA20] Arnab Sarker, Joseph E Gaudio, and Anuradha M Annaswamy. Parameter estimation
362 bounds based on the theory of spectral lines. *arXiv preprint arXiv:2006.12687*, 2020.
- 363 [SMT⁺18] Max Simchowitz, Horia Mania, Stephen Tu, Michael I Jordan, and Benjamin Recht.
364 Learning without mixing: Towards a sharp analysis of linear system identification. In
365 *Conference On Learning Theory*, pages 439–473, 2018.
- 366 [SOF20] Yue Sun, Samet Oymak, and Maryam Fazel. Finite sample system identification: Im-
367 proved rates and the role of regularization. In *2nd LADC 2020 - Learning for Dynamics*
368 *& Control, June 11-12, 2020, 2020*.
- 369 [SRD19] Tuhin Sarkar, Alexander Rakhlin, and Munther A Dahleh. Finite-time system
370 identification for partially observed lti systems of unknown order. *arXiv preprint*
371 *arXiv:1902.01848*, 2019.
- 372 [TBPR17] Stephen Tu, Ross Boczar, Andrew Packard, and Benjamin Recht. Non-asymptotic
373 analysis of robust control from coarse-grained identification. *arXiv preprint*
374 *arXiv:1707.04791*, 2017.
- 375 [TMP20] Anastasios Tsiamis, Nikolai Matni, and George Pappas. Sample complexity of kalman
376 filtering for unknown systems. In *Learning for Dynamics and Control*, pages 435–444.
377 PMLR, 2020.
- 378 [TP20] Anastasios Tsiamis and George Pappas. Online learning of the kalman filter with
379 logarithmic regret. *arXiv preprint arXiv:2002.05141*, 2020.

- 380 [TR19] Stephen Tu and Benjamin Recht. The gap between model-based and model-free meth-
381 ods on the linear quadratic regulator: An asymptotic viewpoint. In *Conference on*
382 *Learning Theory*, pages 3036–3083. PMLR, 2019.
- 383 [Ver18] Roman Vershynin. *High-dimensional probability: An introduction with applications*
384 *in data science*, volume 47. Cambridge University Press, 2018.
- 385 [VODM12] Peter Van Overschee and BL De Moor. *Subspace Identification for Linear Systems*.
386 Springer Science & Business Media, 2012.
- 387 [WJ20] Andrew Wagenmaker and Kevin Jamieson. Active learning for identification of linear
388 dynamical systems. In *COLT 2020 - The 33rd Annual Conference on Learning Theory,*
389 *July 9-12, 2020*, 2020.
- 390 [ZDG⁺96] Kemin Zhou, John Comstock Doyle, Keith Glover, et al. *Robust and optimal control*,
391 volume 40. Prentice hall New Jersey, 1996.

392 **Checklist**

- 393 1. For all authors...
- 394 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
395 contributions and scope? [Yes]
- 396 (b) Did you describe the limitations of your work? [Yes] The main limitations are the
397 sensitivity to process noise (Remark 1 after Theorem 2.2, elaborated in Appendix E),
398 the restriction to Gaussian noise, and the assumption $\rho(A) < 1$.
- 399 (c) Did you discuss any potential negative societal impacts of your work? [N/A] The
400 work is primarily a theory result.
- 401 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
402 them? [Yes]
- 403 2. If you are including theoretical results...
- 404 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 405 (b) Did you include complete proofs of all theoretical results? [Yes]
- 406 3. If you ran experiments...
- 407 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
408 mental results (either in the supplemental material or as a URL)? [Yes] Code is in the
409 supplement. Executing either the Julia file or notebook directly produces the graphs.
- 410 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
411 were chosen)? [Yes] Appendix D
- 412 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
413 iments multiple times)? [No] Graphs are averages over 10 runs; visual inspection of
414 raw data shows that the variability is small.
- 415 (d) Did you include the total amount of compute and the type of resources used (e.g., type
416 of GPUs, internal cluster, or cloud provider)? [No] Code runs in 5 minutes on laptop.
- 417 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 418 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 419 (b) Did you mention the license of the assets? [N/A]
- 420 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 421
- 422 (d) Did you discuss whether and how consent was obtained from people whose data
423 you're using/curating? [N/A]
- 424 (e) Did you discuss whether the data you are using/curating contains personally identifi-
425 able information or offensive content? [N/A]
- 426 5. If you used crowdsourcing or conducted research with human subjects...
- 427 (a) Did you include the full text of instructions given to participants and screenshots, if
428 applicable? [N/A]
- 429 (b) Did you describe any potential participant risks, with links to Institutional Review
430 Board (IRB) approvals, if applicable? [N/A]
- 431 (c) Did you include the estimated hourly wage paid to participants and the total amount
432 spent on participant compensation? [N/A]