
Finite Sample Analysis Of Dynamic Regression Parameter Learning

Mark Kozdoba

Technion, Israel Institute of Technology
markk@ef.technion.ac.il

Edward Moroshko

Technion, Israel Institute of Technology
edward.moroshko@gmail.com

Shie Mannor

Technion, Israel Institute of Technology and NVIDIA Research
shie@ee.technion.ac.il

Koby Crammer

Technion, Israel Institute of Technology
koby@ee.technion.ac.il

Abstract

We consider the dynamic linear regression problem, where the predictor vector may vary with time. This problem can be modeled as a linear dynamical system, with non-constant observation operator, where the parameters that need to be learned are the variance of both the process noise and the observation noise. While variance estimation for dynamic regression is a natural problem, with a variety of applications, existing approaches to this problem either lack guarantees altogether, or only have asymptotic guarantees without explicit rates. In particular, existing literature does not provide any clues to the following fundamental question: In terms of data characteristics, what does the convergence rate depend on? In this paper we study the global system operator – the operator that maps the noise vectors to the output. We obtain estimates on its spectrum, and as a result derive the first known variance estimators with finite sample complexity guarantees. The proposed bounds depend on the shape of a certain spectrum related to the system operator, and thus provide the first known explicit geometric parameter of the data that can be used to bound estimation errors. In addition, the results hold for arbitrary sub Gaussian distributions of noise terms. We evaluate the approach on synthetic and real-world benchmarks.

1 Introduction

A dynamic linear regression (West and Harrison, 1997, Chapter 3), or non-stationary regression, is a situation where we are given a sequence of scalar *observations* $\{Y_t\}_{t \leq T} \subset \mathbb{R}$, and *observation vectors* $\{u_t\}_{t \leq T} \subset \mathbb{R}^n$ such that $Y_t = \langle X_t, u_t \rangle + z_t$ where $X_t \in \mathbb{R}^n$ is a regressor vector, and z_t a random noise term. In contrast to a standard linear regression, the vector X_t may change with time. One common objective for this problem is at time T , to estimate the trajectory of X_t for $t \leq T$, given the observation vectors and observations, $\{u_t\}_{t \leq T}$, $\{Y_t\}_{t \leq T}$, and possibly to forecast Y_{T+1} if u_{T+1} is also known.

In this paper we model the problem as follows:

$$\begin{aligned} X_{t+1} &= X_t + h_t & (1) \\ Y_t &= \langle X_t, u_t \rangle + z_t, & (2) \end{aligned}$$

where $\langle \cdot, \cdot \rangle$ is the standard inner product on \mathbb{R}^n , z_t , the *observation noise*, are zero-mean sub Gaussian random variables, with variance η^2 , and the *process noise* variables h_t take values in \mathbb{R}^n , such that coordinates of h_t are zero-mean sub Gaussian, independent, and have variance σ^2 . All h_t and z_t variables are assumed to be mutually independent. The vectors u_t are an *arbitrary* sequence in \mathbb{R}^n , and the observed, known, quantities at time T are $\{Y_t\}_{t \leq T}$ and $\{u_t\}_{t \leq T}$.

The system (1)-(2) is a special case of a Linear Dynamical System (LDS). As is well known, when the parameters σ, η are given, the mean-squared loss optimal forecast for Y_{T+1} and estimate for X_T are obtained by the Kalman Filter (Anderson and Moore, 1979; Hamilton, 1994; Chui and Chen, 2017). In this paper we are concerned with estimators for σ, η , and *finite* sample complexity guarantees for these estimators.

Let us first make a few remarks about the particular system (1)-(2). First, as a natural model of time varying regression, this system is useful in a considerable variety of applications. We refer to West and Harrison (1997), Chapter 3, for numerous examples. In addition, an application to electricity consumption time series as a function of the temperature is provided in the experiments section of this paper. Second, one may regard the problem of estimating σ, η in (1)-(2) as a pure case of finding the optimal *learning rate* for X_t . Indeed, the Kalman filter equations for (1)-(2), are given by (3)-(4) below, where (3) describes the filtered covariance update and (4) the filtered state update. Here \bar{x}_t is the estimated state, given the observations Y_1, \dots, Y_t , see West and Harrison (1997).

$$C_{t+1} = \frac{\eta^2}{\langle (C_t + \sigma^2 I) u_{t+1}, u_{t+1} \rangle + \eta^2} C_t + \sigma^2 I \quad (3)$$

$$\bar{x}_{t+1} = \bar{x}_t + \frac{C_{t+1}}{\eta^2} u_{t+1} \cdot (Y_{t+1} - \langle \bar{x}_t, u_{t+1} \rangle). \quad (4)$$

In particular, following (4), the role of σ and η may be interpreted as regulating how much the estimate of \bar{x}_{t+1} is influenced, via the operator $\frac{C_{t+1}}{\eta^2}$, by the most recent observation and input Y_{t+1}, u_{t+1} . Roughly speaking, higher values of σ or lower values of η would imply that the past observations are given less weight, and result in an overfit of the forecast to the most recent observation. On the other hand, very low σ or high η would make the problem similar to the standard linear regression, where all observations are given equal weight, and result in a *lag* of the forecast. See Figure 3 in Supplementary Material Section A for an illustration.

Finally, it is worth mentioning that the system (1)-(2) is closely related to the study of *online gradient* (OG) methods (Zinkevich, 2003; Hazan, 2016). In this field, assuming quadratic cost, one considers the update

$$\bar{x}_{t+1} = \bar{x}_t + \alpha \cdot u_{t+1} \cdot (Y_{t+1} - \langle \bar{x}_t, u_{t+1} \rangle), \quad (5)$$

where α is the learning rate, and studies the performance guarantees of the forecaster $\langle \bar{x}_t, u_{t+1} \rangle$. Compared to (4), the update (5) is simpler, and uses a scalar rate α instead of the input-dependent operator rate $\frac{C_{t+1}}{\eta^2}$ of the Kalman filter. However, due to the similarity, every domain of applicability of the OG methods is also a natural candidate for the model (1)-(2) and vice-versa. As an illustration, we compare the OG to Kalman filter based methods with learned σ, η in the experiments section.

In this paper we introduce a new estimation algorithm for σ, η , termed STVE (Spectrum Thresholding Variance Estimator), and prove finite sample complexity bounds for it. In particular, our bounds are an explicit function of the parameters T and $\{u_t\}_{t=1}^T$ for any finite T , and indicate that the estimation error decays roughly as $T^{-\frac{1}{2}}$, with high probability. To the best of our knowledge, these are the first bounds of this kind. As we discuss in detail in Section 2, most existing estimation methods for LDSs, such as subspace identification (van Overschee and de Moor, 1996; Qin, 2006), or improper learning (Anava et al., 2013; Hazan et al., 2017; Kozdoba et al., 2019), do not apply to the system (1)-(2), due to non-stationarity. On the other hand, the methods that do apply to (1)-(2) either lack guarantees, or have only asymptotic analysis which in addition relies strongly on Gaussianity of the noises.

Moreover, our approach differs significantly from the existing methods. We show that the structure of equations (1)-(2) is closely related to, and inherits several important properties from, the classical discrete Laplacian operator on the line — leading to new arguments that have not been explored in the literature. In particular, we use this connection to show that an appropriate inversion of the system produce estimators that are concentrated enough so that σ and η may be recovered. The heart of the paper is the new definition of the estimators that exploits explicitly the shape of a certain data dependent operator, and the subsequent concentration analysis. In particular, this approach yields the first known *geometric* parameters of the data that can be used to bound convergence rates.

The rest of the paper is organized as follows: The related work is discussed in Section 2 and Section 3 contains the necessary definitions. In Section 4 we describe in general lines the methods and the main results of this paper. The technical estimates on certain operator spectra, that are critical to the analysis and may be of independent interest, are stated in Section 5. In Section 6 we present experimental results on synthetic and real data. Due to space constraints, while we outline the main arguments in the text, the full proofs are deferred to the Supplementary Material.

2 Literature

We refer to Chui and Chen (2017); Hamilton (1994); Anderson and Moore (1979); Shumway and Stoffer (2011) for a general background on LDSs, the Kalman Filter and maximum likelihood estimation.

Existing approaches to the variance estimation problem may be divided into three categories: (i) General methods for parameter identification in LDS, either via maximum likelihood estimation (MLE) (Hamilton, 1994), or via subspace identification (van Overschee and de Moor, 1996; Qin, 2006). In particular, finite sample bounds for system identification were given in (Campi and Weyer, 2005; Vidyasagar and Karandikar, 2006) and in the recent work Tsiamis and Pappas (2019). (ii) Methods designed specifically to learn the noise parameters of the system, developed primarily in the control theory community, in particular via the innovation auto-correlation function, such as the classical Mehra (1970); Belanger (1974), or for instance more recent Wang et al. (2017); Dunik et al. (2018). (iii) *Improper Learning* methods, such as Anava et al. (2013); Hazan et al. (2017); Kozdoba et al. (2019). In these approaches, one does not learn the LDS directly, but instead learns a model from a certain auxiliary class and shows that this auxiliary model produces forecasts that are as good as the forecasts of an LDS with “optimal” parameters.

Despite the apparent simplicity of the system (1)-(2), most of the above methods do not apply to this system. This is due to the fact that most of the methods are designed for time invariant, asymptotically stationary systems, where the observation operator (u_t in our notation) is constant and the Kalman gain (or, equivalently $C_t u_t$ in eq. (3)) converges with t . In particular this limitation exists in all the system identification results cited above, and is essential to the approaches taken there. However, if the observation vector sequence u_t changes with time – a necessary property for the dynamic regression problem – the system will no longer be asymptotically stationary. In particular, due to this reason, neither the subspace identification methods, nor any of the improper learning approaches above apply to system (1)-(2).

Among the methods that do apply to (1)-(2) are the general MLE estimation, and some of the auto-correlation methods (Belanger, 1974; Dunik et al., 2018). On one hand, both types of approaches may be applicable to systems apriori more general than (1)-(2). On the other hand, the situation with consistency guarantees – the guarantee that one recovers true parameters given enough observations – for these methods is somewhat complicated. Due to the non-convexity of the likelihood function, the MLE method is not guaranteed to find the true maximum, and as a result the whole method has no guarantees. The results in Belanger (1974); Dunik et al. (2018) do have *asymptotic* consistency guarantees. However, these rely on some explicit and implicit assumptions about the system, the sequence u_t in our case, which can not be easily verified. In particular, Belanger (1974); Dunik et al. (2018) assume *uniform observability* of the system, which we do not assume, and in addition rely on certain implicit assumption about invertibility and condition number of the matrices related to the sequence u_t . Moreover, even if one assumes that the assumptions hold, the results are purely asymptotic, and for any finite T , do not provide a bound of the expected estimation error as a function of T and $\{u_t\}_{t=1}^T$.

In addition, as mentioned earlier, MLE methods by definition must assume that the noises are Gaussian (or belong to some other predetermined parametric family) and autocorrelation based methods also strongly use the Gaussianity assumption. Our approach, on the other hand, requires only sub Gaussian noises with independent coordinates. We note that there are straightforward extensions of our methods to certain cases with dependencies. Indeed, the operator analysis part of this paper does not depend on the distribution of the noises. Therefore, to achieve such an extension, one would only need to correspondingly extend the main probabilistic tool, the Hanson-Wright inequality (Hanson and Wright, 1971; Rudelson et al., 2013, see also Section 4 and Supplementary Material Section E). One such extension, for vectors with the *convex concentration* property, was recently obtained in Adamczak (2015).

3 Notation

We refer to Bhatia (1997) and Vershynin (2018) as general references on the notation introduced below, for operators and sub Gaussian variables, respectively.

Let $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an operator with a singular value decomposition $A = U \cdot \text{Diag}(\lambda_1, \dots, \lambda_s) \cdot W$, where $s \leq \min\{m, n\}$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s > 0$. Note that singular values are strictly positive by definition (that is, vectors corresponding to the kernel of A do not participate in the decomposition $A = U \cdot \text{Diag}(\lambda_1, \dots, \lambda_s) \cdot W$). The Hilbert-Schmidt (Frobenius) norm is defined as $\|A\|_{HS} = \sqrt{\sum_{i=1}^s \lambda_i^2}$. The nuclear and the operator norms are given by $\|A\|_{nuc} = \sum_{i=1}^s \lambda_i$ and $\|A\|_{op} = \lambda_1$ respectively.

A centered ($\mathbb{E}X = 0$) scalar random variable X is sub-Gaussian with constant κ , denoted $X \sim SG(\kappa)$, if for all $t > 0$ it satisfies $\mathbb{P}(|X| > t) \leq 2 \exp -\frac{t^2}{\kappa^2}$. A random vector $X = (X_1, \dots, X_m)$ is κ sub-Gaussian, denoted $X \sim SG_m(\kappa)$, if for every $v \in \mathbb{R}^m$ with $|v| = 1$ the random variable $\langle v, X \rangle$ is κ sub-Gaussian. A random vector X is σ -isotropic if for every $v \in \mathbb{R}^m$ with $|v| = 1$, $\mathbb{E}\langle v, X \rangle = \sigma^2$.

Finally, a random vector $X = (X_1, \dots, X_m)$ is σ -isotropically κ sub-Gaussian with independent components, denoted $X \sim ISG_m(\sigma, \kappa)$ if X_i are independent, and for all $i \leq m$, $\mathbb{E}X_i = 0$, $\mathbb{E}X_i^2 = \sigma^2$ and $X_i \sim SG(\kappa)$. Clearly, if $X \sim ISG_m(\sigma, \kappa)$ then X is σ -isotropic. Recall also that $X \sim ISG_m(\sigma, \kappa)$ implies $X \sim SG_m(\kappa)$ (Vershynin, 2018). The noise variables we discuss in this paper are $ISG(\kappa, \sigma)$.

Throughout the paper, absolute constants are denoted by c, c', c'', \dots etc. Their values may change from line to line.

4 Overview of the approach

We begin by rewriting (1)-(2) in a vector form. To this end, we first encode sequences of T vectors in \mathbb{R}^n , $\{a_t\}_{t \leq T} \subset \mathbb{R}^n$, as a vector $a \in \mathbb{R}^{Tn}$, constructed by concatenation of a_t 's. Next, we define the summation operator $S' : \mathbb{R}^T \rightarrow \mathbb{R}^T$ which acts on any vector $(h_1, h_2, \dots, h_T) \in \mathbb{R}^T$ by

$$S'(h_1, h_2, \dots, h_T) = (h_1, h_1 + h_2, \dots, \sum_{i=1}^{T-1} h_i, \sum_{i=1}^T h_i). \quad (6)$$

Note that S' is an invertible operator. Next, we similarly define the summation operator $S : \mathbb{R}^{Tn} \rightarrow \mathbb{R}^{Tn}$, an n -dimensional extension of S' , which sums n -dimensional vectors. Formally, for $(h_i)_{i=1}^{Tn} \in \mathbb{R}^{Tn}$, and for $1 \leq j \leq n, 1 \leq t \leq T$, $(Sh)_{(t-1) \cdot n + j} = \sum_{i \leq t} h_{(i-1) \cdot n + j}$. Observe that if the sequence of process noise terms $h_1, \dots, h_T \in \mathbb{R}^n$ is viewed as a vector $h \in \mathbb{R}^{Tn}$, then by definition Sh is the \mathbb{R}^{Tn} encoding of the sequence X_t .

Next, given a sequence of observation vectors $u_1, \dots, u_T \in \mathbb{R}^n$, we define the observation operator $O_u : \mathbb{R}^{Tn} \rightarrow \mathbb{R}^T$ by $(O_u x)_t = u_t \cdot (x_{(t-1) \cdot n + 1}, \dots, x_{(t-1) \cdot n + n})$. In words, coordinate t of $O_u x$ is the inner product between u_t and t -th part of the vector $x \in \mathbb{R}^{Tn}$. Define also $Y = (Y_1, \dots, Y_T) \in \mathbb{R}^T$ to be the concatenation of Y_1, \dots, Y_T . With this notation, one may equivalently rewrite the system (1)-(2) as follows:

$$Y = O_u Sh + z, \quad (7)$$

where h and z are independent zero-mean random vectors in \mathbb{R}^{Tn} and \mathbb{R}^T respectively, with independent sub Gaussian coordinates. The variance of each coordinate of h is σ^2 and each coordinate of z has variance η^2 .

Up to now, we have reformulated our data model as a single vector equation. Note that in that equation, the observations Y and both operators O_u and S are known to us. Our problem may now be reformulated as follows: Given $Y \in \mathbb{R}^T$, assuming Y was generated by (7), provide estimates of σ, η .

As a motivation, we first consider taking the expectation of the norm squared of eq. (7). For any operator $A : \mathbb{R}^m \rightarrow \mathbb{R}^m$ and zero-mean vector h with independent coordinates and coordinate variance σ^2 , we have $\mathbb{E}|Ah|^2 = \|A\|_{HS}^2 \sigma^2$, where $\|A\|_{HS}$ is the Hilbert-Schmidt (or Frobenius) norm of A . Taking the norm and expectation of (7), and dividing by T^2 , we thus obtain

$$\frac{\mathbb{E}|Y|^2}{T^2} = \frac{\|O_u S\|_{HS}^2}{T^2} \sigma^2 + \frac{T}{T^2} \eta^2. \quad (8)$$

Next, note that $\|O_u S\|_{HS}^2$ is known, and an elementary computation shows that $\frac{\|O_u S\|_{HS}^2}{T^2}$ is of constant order (as a function of T ; see (25)), while the coefficient of η^2 is $\frac{1}{T}$. Thus, if the quantity $\frac{|Y|^2}{T^2}$ were close enough to its expectation with high probability, we could take this quantity as a (slightly biased) estimator of σ^2 . However, as it will become apparent later, the deviations of $\frac{|Y|^2}{T^2}$ around the expectation are also of constant order, and thus $\frac{|Y|^2}{T^2}$ can not be used as an estimator. The reason for these high deviations of $\frac{|Y|^2}{T^2}$ is that the spectrum of $O_u S$ is extremely peaked. The highest squared singular value of $O_u S$ is of order T^2 , the same order as sum of all of them, $\|O_u S\|_{HS}^2$. Contrast this with the case of identity operator, $Id : \mathbb{R}^{Tn} \rightarrow \mathbb{R}^{Tn}$: We have $\mathbb{E}|Id(h)|^2 = \mathbb{E}|h|^2 = Tn\sigma^2$, and one can also easily show that, for instance, $Var |Id(h)|^2 = Tn\sigma^2$, and thus the deviations are of order $\sqrt{Tn}\sigma$ – a smaller order than $\mathbb{E}|Id(h)|^2$. While for the identity operator the computation is elementary, for a general operator A the situation is significantly more involved, and the bounds on the deviations of $|Y|^2$ will be obtained from the Hanson-Wright inequality (Hanson and Wright, 1971, see also Rudelson et al. (2013)), combined with standard norm deviation bounds for isotropic sub Gaussian vectors.

With these observations in mind, we proceed to flatten the spectrum of $O_u S$ by taking the pseudo-inverse. Let $R : \mathbb{R}^T \rightarrow \mathbb{R}^{Tn}$ be the pseudo-inverse, or Moore-Penrose inverse of $O_u S$. Specifically, let

$$O_u S = U \circ \text{Diag}(\gamma_1, \dots, \gamma_T) \circ W, \quad (9)$$

be the singular value decomposition of $O_u S$, where $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_T$ are the singular values.

For the rest of the paper, we will assume that all of the observation vectors u_t are non-zero. This assumption is made solely for notational simplicity and may easily be avoided, as discussed later in this section. Under this assumption, since S is invertible and O_u has rank T , we have $\lambda_t > 0$ for all $t \leq T$. For $i \leq T$, denote $\chi_i = \gamma_{T+1-i}^{-1}$. Then χ_i are the singular values of R , arranged in a non-increasing order, and we have by definition

$$R = W^* \circ \text{Diag}(\chi_T, \chi_{T-1}, \dots, \chi_2, \chi_1) \circ U^*, \quad (10)$$

where W^*, U^* denote the transposed matrices of U, V , defined in (9).

Similarly to Eq. (8), we apply R to (7), and since $\|RO_u S\|_{HS} = T$, by taking the expectation of the squared norm we obtain

$$\frac{|RY|^2}{T} = \sigma^2 + \frac{\|R\|_{HS}^2}{T} \eta^2 + \frac{|RY|^2}{T} - \frac{\mathbb{E}|RY|^2}{T}. \quad (11)$$

In this equation, the deviation term $\frac{|RY|^2}{T} - \frac{\mathbb{E}|RY|^2}{T}$ is of order $O(\frac{1}{\sqrt{T}})$ with high probability (Theorem 1). Moreover, the coefficient of σ^2 is 1, and the coefficient of η^2 , which is $\frac{\|R\|_{HS}^2}{T}$, is of order at least $\Omega(\frac{1}{\log^2 T})$ (Theorem 3, see Section 5 for additional details) – much larger order than $\frac{1}{\sqrt{T}}$. Since $|RY|^2$ and $\|R\|_{HS}^2$ are known, it follows that we have obtained one equation satisfied by σ^2 and η^2 up to an error of $\frac{1}{\sqrt{T}}$, where both coefficients are of order larger than the error.

Algorithm 1 Spectrum Thresholding Variance Estimator (STVE)

- 1: **Input:** Observations Y_t , observation vectors u_t , with $t \leq T$, and $p = \alpha T$.
2: Compute the SVD of $O_u S$,

$$O_u S = U \circ \text{Diag}(\gamma_1, \dots, \gamma_T) \circ W,$$

where $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_T > 0$. Denote $\chi_i = \gamma_{T+1-i}^{-1}$ for $1 \leq i \leq T$.

- 3: Construct the operators

$$R = W^* \circ \text{Diag}(\chi_T, \dots, \chi_1) \circ U^*$$

and

$$R' = W^* \circ \text{Diag}(0, 0, \dots, 0, \chi_p, \dots, \chi_1) \circ U^*$$

- 4: Produce the estimates:

$$\hat{\eta}^2 = \frac{|R'Y|^2}{p} - \frac{|RY|^2}{T} \cdot \frac{\|R'\|_{HS}^2}{p} - \frac{\|R\|_{HS}^2}{T} \quad !$$

$$\hat{\sigma}^2 = \frac{|RY|^2}{T} - \frac{\|R\|_{HS}^2}{T} \hat{\eta}^2.$$

Next, we would like to obtain another linear relation between σ^2 , η^2 . To this end, choose some $p = \alpha T$, where $0 < \alpha < 1$ is of constant order. The possible choices of p are discussed later in this section. We define an operator $R' : \mathbb{R}^T \rightarrow \mathbb{R}^{Tn}$ to be a version of R truncated to the first p singular values. If (10) is the SVD decomposition of R , then

$$R' = W^* \circ \text{Diag}(0, 0, \dots, \chi_p, \chi_{p-1}, \dots, \chi_1) \circ U^*.$$

Similarly to the case for R , we have

$$\frac{|R'Y|^2}{p} = \sigma^2 + \frac{\|R'\|_{HS}^2}{p} \eta^2 + \frac{|R'Y|^2}{p} - \frac{\mathbb{E}|R'Y|^2}{p} \quad ! \quad (12)$$

The deviations in (12) are also described by Theorem 1. Note also that since $\|R'\|_{HS}^2$ is the sum of p largest squared singular values of R , by definition it follows that $\frac{\|R'\|_{HS}^2}{p} \geq \frac{\|R\|_{HS}^2}{T}$.

Now, given two equations in two unknowns, we can solve the system to obtain the estimates $\hat{\sigma}^2$ and $\hat{\eta}^2$. The full procedure is summarized in Algorithm 1, and the bounds implied by Theorem 1 on the estimators $\hat{\sigma}^2$ and $\hat{\eta}^2$ are given in Corollary 2. We first state these results, and then discuss in detail the various parameters appearing in the bounds.

Theorem 1. Consider a random vector $Y \in \mathbb{R}^T$ of the form $Y = O_u S h + z$ where $h \sim ISG_{Tn}(\sigma, \kappa)$ and $z \sim ISG_T(\eta, \kappa)$. Set $|u_{min}| = \min_t |u_t|$. Then for any $0 < \delta < 1$,

$$\mathbb{P} \left[\frac{|RY|^2}{T} - \sigma^2 + \frac{\|R\|_{HS}^2}{T} \eta^2 \geq c \frac{B}{\sqrt{T}} \right] \leq 4\delta, \quad (13)$$

$$\mathbb{P} \left[\frac{|R'Y|^2}{p} - \sigma^2 + \frac{\|R'\|_{HS}^2}{p} \eta^2 \geq c \frac{B}{\sqrt{p}} \right] \leq 4\delta \quad (14)$$

where B is given by

$$B = 1 + \kappa^2 \left(1 + |u_{min}|^{-2} \log \frac{1}{\delta} \right). \quad (15)$$

The bounds on the estimators of Algorithm 1 are given in the following Corollary. As discussed below, in addition to Theorem 1, the key to the derivation of this Corollary are the estimates of the spectrum of R , given in Theorem 3, Section 5.

(a) Spectrum ($\hat{\lambda}_i^2$) of R for electricity data (blue). (b) Similar figure for synthetic data $\bar{T} = 500$. $kR^0(p)k_{HS}^2 = p$ as a function of p (orange). The mean $kRk_{HS}^2 = T$ (green).

Figure 1: Spectra of R

Corollary 2. Let c_2^2 ; b_2 be the estimators of c_2^2 ; b_2 obtained from Algorithm 1 with $p = \frac{1}{4}T$. Set $j_{u_{\max}} = \max_t j_{u_t}$. Then for any $0 < \epsilon < 1$, with probability at least $1 - \epsilon$,

$$c_2^2 - c_2^2 \leq \frac{B}{p} \left(1 + \frac{pkRk_{HS}^2}{T kR^0k_{HS}^2} \right)^{1/2}; \quad (16)$$

$$b_2 - b_2 \leq \frac{B}{p} \left(1 + \frac{pkRk_{HS}^2}{T kR^0k_{HS}^2} \right)^{1/2} n^2 j_{u_{\max}} j^2 \log^2 T; \quad (17)$$

with B given by (15).

We first discuss the assumption $j_{u_{\min}} > 0$. This assumption is made solely for notational convenience, as detailed below. To begin, note that some form of lower bound on the norms of the observation vectors u_t must appear in the bounds. This is simply because if one had 0 for all T , then clearly no estimate of λ would have been possible. On the other hand, our use of the smallest value $j_{u_{\min}}$ may seem restrictive at first. We note however, that instead of considering the observation operator $O_u : R^{Tn} \rightarrow R^T$, one may consider the operator $O_{u_t} : R^{Tn} \rightarrow R^T$ for any subsequence $u_{t=1}^T$. The observation vector of O_{u_t} would be correspondingly restricted to the subsequence of indices. This allows us to treat missing values and to exclude any u_t with small norms. All the arguments in Theorems 1 and 3 hold for this modified O_{u_t} without change. The only price that will be paid is that T will be replaced by \bar{T} in the bounds. Moreover, we note that typically we have $j_{u_t} \geq 1$ by construction, see for instance Section 2.2. Additional discussion of missing values may be found in Supplementary Material Section 2.

Next, up to this point, we have obtained two equations (11)-(12), in two unknowns, c_2^2 ; b_2 . Note that in order to be able to obtain c_2^2 from these equations, at least one of the coefficients of either $\frac{kRk_{HS}^2}{T}$ or $\frac{kR^0k_{HS}^2}{p}$ must be of larger order than $\frac{1}{T}$, the order of deviations. Providing lower bounds on these quantities is one of the main technical contributions of this work. This analysis uses the connection between the operator O_u and the Laplacian on the line, and resolves the issue of translating spectrum estimates for the Laplacian into the spectral estimates for R . We note that there are no standard tools to study the spectrum of R , and our approach proceeds indirectly via the analysis of the nuclear norm $\|O_u\|_S$. These results are stated in Theorem 1. In particular, we show that $\frac{kRk_{HS}^2}{T}$ is $(\frac{1}{\log^2 T})$, which is the source of the log factor in (17).

Finally, in order to solve the equations (11)-(12), not only the equations must have large enough coefficients, but the equations must be different. This is reflected by the term $1 + \frac{pkRk_{HS}^2}{T kR^0k_{HS}^2}$ in

(16), (17). Equivalently, while $\frac{kR^0k_{HS}^2 = p}{kRk_{HS}^2 = T} \geq 1$ by definition, we would like to have

$$\frac{kR^0k_{HS}^2 = p}{kRk_{HS}^2 = T} \geq 1 + \text{const} \quad (18)$$

for the bounds (6), (17) to be stable. Note that since $\|R\|_{HS}^2$ and $\|kR\|_{HS}^2$ are computed in Algorithm 1, the condition (8) can simply be verified before the estimators \hat{c}_S, \hat{c}_R are returned.

It is worth emphasizing that for simple choices of ρ say $\rho = \frac{1}{4}T$, the condition (8) does hold in practice. Note that, for any ρ , we can have $\frac{\|kR\|_{HS}^2 = \rho}{\|kR\|_{HS}^2 = T} = 1$ only if the spectrum of R is constant. Thus (8) amounts to stating that the spectrum of R exhibits some decay. As we show in experiments below, the spectrum (squared) of u_t derived from daily temperature features, or for random Gaussian u_t , indeed decays. See Section 6, Figures 1a and 1b. In particular, in both cases (8) holds with $\text{const} > 1$. Additional bounds on the quantity $\frac{\|kR\|_{HS}^2}{\|R\|_{HS}^2}$ under various assumptions on the sequence u_t are given in Section 6 of the Supplementary Material.

5 Properties of $O_u S$ and R

As discussed in Section 4 (see the discussion following eq. (1)), one of the crucial points enabling Algorithm 1 and its analysis is the fact that the quantity $\frac{\|kR\|_{HS}^2}{T}$ is bounded below by an expression that is of much higher order than the noise magnitude $\frac{1}{T}$.

In this section we provide the formal statement of this and other associated results, and discuss the related arguments. First, we obtain the following bound on the spectral norm $\|O_u S\|$ (Lemma 4, Supplementary Material Section 5). Recall that the nuclear norm was defined in Section 3 and that for a sequence u_t we set $\|u\|_{\max} = \max_t |u_t|$ and $\|u\|_{\min} = \min_t |u_t|$. Then:

$$\|O_u S\|_{\text{nuc}} = \sum_{t=1}^T (O_u S)_{tt} \leq 4n \|u\|_{\max} j T \log T; \quad (19)$$

The proof of this bound exploits the connection between $O_u S$ and the Laplacian on the line. In particular, we use the fact that the eigenvalues of the Laplacian are known precisely, satisfying $\lambda_j = 2 \sin \frac{(T-1)j\pi}{2T}$. Next, we state the lower (and upper) bounds for

Theorem 3. Let $R : \mathbb{R}^T \rightarrow \mathbb{R}^{T \times n}$ be the pseudoinverse of $O_u S$. Then

$$\frac{1}{n \|u\|_{\max} j \log T} \leq \|kR\|_{\text{op}} \leq 2 \|u\|_{\min} j^{-1}; \quad (20)$$

$$\frac{1}{n^2 \|u\|_{\max} j^2 \log^2 T} \leq \|kR\|_{HS}^2 \leq 4 \|u\|_{\min} j^{-2} T; \quad (21)$$

$$\frac{1}{n^4 \|u\|_{\max} j^4 \log^4 T} \leq \|R\|_{HS}^2 \leq 16 \|u\|_{\min} j^{-4} T; \quad (22)$$

Due to the complicated structure of R as a pseudo-inverse of a composition of operators, there are no direct ways to control individual eigenvalues of R . Thus the main technical issue resolved in Theorem 3 is nevertheless obtaining lower bounds on $\|kR\|_{HS}^2$. Our approach is rather indirect, and we obtain these bounds from the nuclear norm bound (19) via a Markov type inequality on the eigenvalues.

6 Experiments

6.1 Synthetic Data

In this section the performance of STVE is evaluated on synthetic data. The data was generated by the LDS (1)-(2), using Gaussian noises with $\sigma^2 = 0.5$; $\sigma^2 = 2$. The input dimension was $n = 5$, and the input sequence u_t sampled from the Gaussian $\mathcal{N}(0, I_n)$.

We run the STVE algorithm for different values of T , where for each T we sampled the data 150 times. Figure 2a shows the average (over 150 runs) estimation error for both process and observation noise variances for various values of T . As expected from the bounds in Corollary 2, it may be observed in Figure 2a that the estimation errors decay roughly at the rate $\frac{1}{\text{const} \cdot T}$. A typical spectrum of R is shown in Figure 1b. For larger T , the spectra also exhibits similar decay.

(a) STVE variance estimation error vs T , and the $T^{-\frac{1}{2}}$ decay. (b) Load (y-axis) against Temperature (x-axis), both axis normalized. (c) Smoothed prediction errors. Stationary Regression trained on the first half (blue), regression learned on the second half (green). The full data (blue points), regression learned on the first half (orange), MLE (orange), STVE (orange), OG (red).

Figure 2: Evaluation

6.2 Temperatures and Electricity Consumption

In this section we examine the relation between daily temperatures and electricity consumption in the data from Hong et al. (2014) (see also Hong (2016)). The following forecasting methods are compared: a stationary regression, an online gradient, and a Kalman filter for a dynamic regression, with parameters learned via MLE or STVE. We find that the Kalman filter methods provide the best performance, with no significant difference between STVE and MLE derived systems.

The data consists of total daily electricity consumption (load) and the average daily temperature, v_t , for a certain region, for the period Jan 2004 to Jun 2008. Full details on the preprocessing of the data, as well as additional details on the experiments, are given in Supplementary Material Section 2. Here we note that the data contains missing load values, for 9 non-consecutive weeks (out of about 234 weeks total). All methods discussed here, including STVE, can naturally incorporate missing values, as discussed in Supplementary Material Section 2.

An elementary inspection of the data reveals that the load may be reasonably approximated as a quadratic function of the temperature, $y_t = x_{t,1} + x_{t,2} v_t + x_{t,3} v_t^2$, where $u_t = (1; v_t; v_t^2)$ is the observation vector (features), and $x = (x_{t,1}; x_{t,2}; x_{t,3})$ is the possibly time varying regression vector. This is shown in Figure 2b, where we fit a stationary (time invariant) regression of the above form, using either only the first or only the second half of the data. We note that these regressions differ – the regression vector changes with time. It is therefore of interest to track it via online regression.

We use the first half of the data (train set) to learn the parameters of the online regression (1)-(2) via MLE optimization and using STVE. We also use the train set to find the optimal learning rate for the OG forecaster described by the update equation (5). This learning rate is chosen as the rate that yields smallest least squares forecast error on the train set. In addition, we learn a time independent, stationary regression on the first half of the data.

We then employ the learned parameters to make predictions of the load given the temperature, by all four methods. The predictions for the system (2) are made with a Kalman filter (at time t we use the filtered state estimate, which depends only on $y_1; \dots; y_t$ and $u_1; \dots; u_t$, and make the prediction $\hat{y}_{t+1} = h_{t+1}^T \hat{x}_{t+1}$).

Daily squared prediction errors (that is $(y_t - \hat{y}_t)^2$) are shown in Figure 2c (smoothed with a moving average of 50 days). We see that the adaptive models (MLE, STVE, OG) outperform the stationary regression already on the train set (first half of the data), and that the difference in performance becomes dramatic on the second half (test). It is also interesting to note that the performance of the Kalman filter based methods (MLE, STVE) is practically identical, but both are somewhat better than the simpler OG approach.

We also note that by construction, we have $j_{\min} = 1$ in this experiment, due to the constant coordinate, and also $j_{\max} = 5$, due to the normalization. Since these operations are typical for any regression problem, we conclude that the direct invariance of j and j_{\max} on the bounds in Theorem 1 and Corollary 2 will not usually be significant.

7 Conclusion And Future Work

In this work we introduced the STVE algorithm for estimating the variance parameters of LDSs of type (1)-(2), and obtained the first sample complexity guarantees for such estimators. We have also shown how the shape of the spectrum can be exploited to obtain the estimators and the related bounds, thus providing the first explicit geometric parameter of the data that affects the bounds.

As discussed in Section 1 and demonstrated in Section 6 the system (1)-(2) is of independent interest in applications. However, we also believe that the analysis presented here is an important first step towards a finite time data-dependent quantitative understanding of general LDSs, and perhaps even non-linear dynamical systems.

Acknowledgments and Disclosure of Funding

This research was supported by the ISRAEL SCIENCE FOUNDATION (grant No. 2199/20).

References

- Adamczak, R. (2015). A note on the hanson-wright inequality for random vectors with dependencies. *Electronic Communications in Probability* 20.
- Anava, O., Hazan, E., Mannor, S., and Shamir, O. (2013). Online learning for time series prediction. In COLT 2013 - The 26th Annual Conference on Learning Theory, June 12-14, 2013, Princeton University, NJ, USA
- Anderson, B. and Moore, J. (1979). *Optimal Filtering*. Prentice Hall.
- Belanger, P. R. (1974). Estimation of noise covariance matrices for a linear time-varying stochastic process. *Automatica* 10(3).
- Bhatia, R. (1997). *Matrix Analysis Graduate Texts in Mathematics*. Springer New York.
- Campi, M. C. and Weyer, E. (2005). Guaranteed non-asymptotic confidence regions in system identification. *Automatica* 41(10):1751–1764.
- Chui, C. and Chen, G. (2017). *Kalman Filtering: with Real-Time Applications*. Springer International Publishing.
- Dunik, J., Kost, O., and Straka, O. (2018). Design of measurement difference autocovariance method for estimation of process and measurement noise covariance. *Automatica* 90.
- Gohberg, I. and Krein, M. (1969). *Introduction to the Theory of Linear Nonselfadjoint Operators*. Translations of mathematical monographs. American Mathematical Society.
- Hamilton, J. (1994). *Time Series Analysis*. Princeton University Press.
- Hanson, D. L. and Wright, E. T. (1971). A bound on tail probabilities for quadratic forms in independent random variables. *Ann. Math. Statist.* 42.
- Hazan, E. (2016). Introduction to online convex optimization. *Found. Trends Optim.*
- Hazan, E., Singh, K., and Zhang, C. (2017). Online learning of linear dynamical systems. In *Advances in Neural Information Processing Systems*, pages 6686–6696.
- Hong, T. (2016). Tao hong's blog. <http://blog.drhongtao.com/2016/07/gefcom2012-load-forecasting-data.html>. Accessed: 1/8/2019.
- Hong, T., Pinson, P., and Fan, S. (2014). Global energy forecasting competition. *International Journal of Forecasting* 30:357–363.
- Kozdoba, M., Marecek, J., Tchakian, T. T., and Mannor, S. (2019). On-line learning of linear dynamical systems: Exponential forgetting in Kalman Iteration. *AAAI*.

- Mehra, R. (1970). On the identification of variances and adaptive kalman filter. *IEEE Transactions on Automatic Control* 15(2).
- Mitchell, A. and Griffiths, D. (1980). *The finite difference method in partial differential equations* Wiley-Interscience publication. Wiley.
- Petris, G. (2010). An R package for dynamic linear models. *Journal of Statistical Software*
- Qin, S. J. (2006). An overview of subspace identification. *Computers & chemical engineering* 30(10-12):1502–1513.
- Rudelson, M., Vershynin, R., et al. (2013). Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability* 18.
- Shumway, R. and Stoffer, D. (2011). *Time Series Analysis and Its Applications* (3rd ed.)
- Tsiamis, A. and Pappas, G. J. (2019). Finite sample analysis of stochastic system identification. In *2019 IEEE 58th Conference on Decision and Control (CDC)* pages 3648–3654. IEEE.
- van Overschee, P. and de Moor, L. (1996). *Subspace identification for linear systems: theory, implementation, applications* Kluwer Academic Publishers.
- Vershynin, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science* Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Vidyasagar, M. and Karandikar, R. L. (2006). A learning theory approach to system identification and stochastic adaptive control. *Probabilistic and randomized methods for design under uncertainty*, pages 265–302.
- Wang, H., Deng, Z., Feng, B., Ma, H., and Xia, Y. (2017). An adaptive kalman filter estimating process noise covariance. *Neurocomputing* 223.
- West, M. and Harrison, J. (1997). *Bayesian Forecasting and Dynamic Models* (2nd ed.) Springer-Verlag.
- Zinkevich, M. (2003). Online convex programming and generalized in infinitesimal gradient ascent. *ICML*.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]**
 - (b) Did you describe the limitations of your work? **[Yes]**
 - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]** The main focus of this work is a theoretical analysis.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - (b) Did you include complete proofs of all theoretical results? **[Yes]** Thus Supplementary Material contains all the proofs. Proofs outlines are given in the main body of the paper.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** The data is publicly available, see the references in Section 6. The short code is not provided at the moment, but can be fully derived from Algorithm 1.

- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes](#) See Section 6 and I.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes](#)
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes](#)
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes](#)
 - (b) Did you mention the license of the assets? [N/A](#) The link to the data author's documentation is provided.
 - (c) Did you include any new assets either in the supplemental material or as a [URL](#)?
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A](#)
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A](#)

Figure 3: Lag and Over fitting to the most recent observation, for various variance values. See the discussion in Section 1 and A.

Finite Sample Analysis Of Dynamic Regression Parameter Learning - Supplementary Material

A Outline

This Supplementary Material is organized as follows: The proofs of the results of Section 1, including Theorem 3, as well as proofs of Theorems 4 and Corollary 2, are given in Section 6 to F. In

Section G we state and prove additional bounds on the quantity $\frac{\rho k_R k_{HS}^2}{T k_R \rho k_{HS}^2}$ that appears in Corollary 2. Section H contains the discussion of missing values in STVE. Additional details on the electricity consumption experiment are given in Section

Finally, we describe Figure 8. The data (observations, red) was generated from the system (2) (with one dimensional state $n = 1$), and we set $u_t = 1$. The states produced by the Kalman filter with ground truth values of θ are shown in green, while states obtained from other choices of parameters are shown in black and blue.

B Bounds on $O_u S$, Lemma 4

In this Section we prove the following bounds on the spectrum of $O_u S$. See Section 5 for a discussion of these bounds.

Lemma 4. The singular values of $O_u S$ satisfy the following:

$$\lambda_1(O_u S) \leq \sum_{t=1}^T |u_{\max}|^2 T \quad \text{and} \quad \lambda_T(O_u S) \geq \frac{1}{2} |u_{\min}|^2 \quad (23)$$

$$k_{O_u S} = \frac{\sum_{t=1}^T |u_{\max}|^2 T}{\lambda_T(O_u S)} \leq \frac{4n |u_{\max}|^2 T \log T}{\frac{1}{2} |u_{\min}|^2 T} \quad (24)$$

$$\frac{1}{4} |u_{\min}|^2 T^2 \leq k_{O_u S}^2 = \frac{\sum_{t=1}^T |u_t|^2}{\lambda_1(O_u S)^2} \leq \frac{1}{2} |u_{\max}|^2 T^2 \quad (25)$$

The proof of this lemma uses the following auxiliary result: Let $D_T : \mathbb{R}^T \rightarrow \mathbb{R}^T$ be the difference operator $(D_T x)_t = x_{t+1} - x_t$ for $t = 1, \dots, T-1$. In the field of Finite Difference methods, the operator $D_T D_T^T$ is known as the Laplacian on the line, or as the discrete derivative with Dirichlet boundary conditions, and is well studied. The eigenvalues of $D_T D_T^T$ may be derived by a direct computation, and correspond to the roots of the Chebyshev polynomial of second kind of order T . In particular, the following holds:

Lemma 5. The operator D_T has kernel of dimension 1 and singular values $\lambda_l(D_T) = 2 \sin \frac{(l-1)\pi}{2T}$ for $l = 1, \dots, T-1$.

We refer to Mitchell and Griffiths (1980) for the proof of Lemma 5. Next, in Lemma 6, we show that the inverse of the operator S^0 , defined in (6), is a one dimensional perturbation D_T , which implies bounds on singular values s_t^{-1} .

Lemma 6. The singular values s_t^{-1} satisfy

$$2 \sin \frac{(T-t)}{2(T+1)} \leq s_t(S_T^0)^{-1} \leq 2 \sin \frac{(T+1-t)}{2(T+1)} \quad (26)$$

for $1 \leq t \leq T-1$, and

$$\frac{1}{T} \leq s_T(S_T^0)^{-1} \leq 2 \sin \frac{1}{2(T+1)} \quad (27)$$

The proof of this is given in the next section. With the estimates of Lemma 6, the bounds in Lemma 4 follow. Proof of Lemma 4:

Proof. By Lemma 6,

$$kS_{k_{op}} \leq T \text{ and } |S_{k_{op}}^{-1} x| \leq \frac{1}{2} |x| \text{ for all } x \in \mathbb{R}^T \quad (28)$$

Note also that S is by definition a collection of n independent copies of S_T^0 , and therefore the spectrum of S is that of S_T^0 , but each singular value is taken with multiplicity n . In particular it follows that (28) holds also for S itself. Since clearly $kO_u k_{op} \leq j u_{\max} j$, the upper bound on $kO_u S k_{op}$ in (23) follows from (28).

For the lower bound, denote by V^0 the orthogonal complement to the kernel of O_u , $V^0 = (\text{Ker}(O_u))^\perp$. Denote by $P_{V^0} : \mathbb{R}^T \rightarrow V^0$ the orthogonal projection onto V^0 . We have in particular that $O_u S = O_u P_{V^0} S$. Next, the operator S maps the unit ball $B_{\mathbb{R}^T}$ of \mathbb{R}^T into an ellipsoid E_T , and by (28), we have $\frac{1}{2} B_{\mathbb{R}^T} \subset E_T$. It therefore follows that

$$\frac{1}{2} B_{V_1} \subset (P_{V^0} S)(B_{\mathbb{R}^T}); \quad (29)$$

where B_{V_1} is the unit ball of V_1 . It remains to observe that for every V_1 , we have

$$|O_u x| \geq j u_{\min} j |x| \quad (30)$$

Combining (29) and (30), we obtain the lower bound in (23).

To derive (24), recall that the nuclear norm is sub-multiplicative with respect to the operator norm:

$$kO_u S k_{nuc} \leq kO_u k_{op} kS k_{nuc} \leq j u_{\max} j kS k_{nuc} \quad (31)$$

This follows for instance from the characterization of the nuclear norm as trace dual of the operator norm (Bhatia 1997, Propositions IV.2.11, IV.2.12). Next, since the spectrum of S is the spectrum of S^0 taken with multiplicity n , we have

$$kS k_{nuc} = n kS^0 k_{nuc}; \quad (32)$$

and it remains to bound the nuclear norm of S^0 .

Using the inequality

$$\sin \frac{x}{2} \geq \frac{x}{4} \text{ for all } x \in [0; 1]; \quad (33)$$

and Lemma 6, we have

$$kS^0 k_{nuc} = T + 2 \sum_{t=T-1}^X \sin^{-1} \frac{T-t}{2(T+1)} \leq 2T \sum_{t=T}^X \frac{1}{t+1} \leq 4T \log T \quad (34)$$

Combining (31), (32) and (34), the inequality (24) follows.

It remains to estimate the Hilbert-Schmidt norm of S , which can be done by a direct computation. Recall that for any operator $A : \mathbb{R}^m \rightarrow \mathbb{R}^m$ the Hilbert-Schmidt norm satisfies

$$kA k_{HS}^2 = \text{tr}(A^* A) = \sum_{i=1}^m |A_{ii}|^2 + \sum_{i \neq j} |A_{ij}|^2 = \sum_{i=1}^m |A_{ii}|^2 + \sum_{i \neq j} |A_{ji}|^2; \quad (35)$$

for any orthonormal basis $\{e_i\}$ in \mathbb{R}^m . Let e_i be the standard basis vector in \mathbb{R}^n corresponding to coordinate i at time t . Let e_0, \dots, e_{T-1} denote the standard basis in \mathbb{R}^n . Then

$$O_u S e_i = \sum_{t=0}^{T-1} u_{t0_i} e_0; \quad (36)$$

where u_{t0_i} is the i -th coordinate of u_{t0} . It follows that

$$\|O_u S e_i\|^2 = \sum_{t=0}^{T-1} u_{t0_i}^2 \quad (37)$$

and hence

$$k O_u S k_{HS}^2 = \sum_{t=0}^{T-1} \sum_{i=1}^n \|O_u S e_i\|^2 = \sum_{t=0}^{T-1} \sum_{i=1}^n \|u_{t0}\|^2 = \sum_{t=0}^{T-1} \|u_{t0}\|^2 \quad (38)$$

The bounds (45) follow directly from (38). \square

C Proof of Lemma 6

Proof. Recall that the operator $S_T^{0,1}$ is given by $S_T^{0,1}x = (x_1; x_2 \dots x_T - 1)$ and the operator $D = D_T : \mathbb{R}^T \rightarrow \mathbb{R}^{T-1}$ is given by $Dx = (x_2 \dots x_T)$. Let $V = \text{span}\{e_2, \dots, e_T\}$ be the subspace spanned by all but the first coordinate. Let $P_V : \mathbb{R}^T \rightarrow \mathbb{R}^T$ be the projection onto V , i.e. a restriction to second to T -th coordinate. Observe that the action of $D P_V$ is equivalent to that of $S_T^{0,1}$. Therefore the singular values of $S_T^{0,1}$ are identical to those of $D P_V$. To obtain bounds on the singular values of $D P_V$, note that $(D P_V) D P_V = P_V D D P_V$ – that is, $P_V D D P_V$ is a compression of $D D$. Thus, by the Cauchy's Interlacing Theorem (Horn & Mirsky, 1997, Corollary III.1.5),

$$\lambda_t(D D) \leq \lambda_t(P_V D D P_V) \leq \lambda_{t+1}(D D) \quad (39)$$

for all $t = 1, \dots, T-1$. In conjunction with Lemma 5 this provides us with the estimates for all but the smallest singular value of $S_T^{0,1}$ (since $\lambda_T(D D) = 0$). We therefore estimate $\lambda_{T-1}(S_T^{0,1}) = \lambda_{T-1}(D P_V)$ directly, by bounding the norm of $S_T^{0,1}$. Indeed, for any T , by the Cauchy-Schwartz inequality,

$$\begin{aligned} \|S_T^{0,1}x\|^2 &= \sum_{t=1}^T \sum_{i=1}^n \left(\sum_{i=1}^n x_i A_{t,i} \right)^2 = \sum_{t=1}^T \sum_{i=1}^n \left(\sum_{i=1}^n x_i^2 A_{t,i}^2 \right) \\ &= \sum_{t=1}^T \sum_{i=1}^n x_i^2 A_{t,i}^2 = \sum_{t=1}^T \|A_{t,\cdot}\|^2 \end{aligned} \quad (40)$$

$$= \sum_{t=1}^T \sum_{i=1}^n x_i^2 A_{t,i}^2 = \sum_{t=1}^T \|A_{t,\cdot}\|^2 \quad (41)$$

Thus we have $\|S_T^{0,1}\|_{op} \leq \sqrt{\sum_{t=1}^T \|A_{t,\cdot}\|^2}$, which concludes the proof of the Lemma. \square

D Proof of Theorem 3

Proof. The bound $\text{orkR}k_{op}$ follows directly from the lower bound on the singular values of S in (23). Since R is of rank T , the upper bounds $\text{dkR}k_{HS}$ follow directly from the $\text{orkR}k_{op}$ bound.

The lower bounds $\text{orkR}k_{HS}^2$ and $\text{dkR}k_{HS}^2$ follow from the upper bounds on the nuclear norm in Lemma 4. Note that this argument would not have worked if we only had upper bounds on the Hilbert-Schmidt norm, rather than the nuclear norm in Lemma 4. Denote $\lambda_i = \lambda_i(O_u S)$. From (24) in Lemma 4, the number of λ_i that are larger than $\frac{1}{4n} \sum_{j=1}^n \log T$ satisfies

$$\#\{i : \lambda_i \geq \frac{1}{4n} \sum_{j=1}^n \log T\} \leq \frac{T}{2} \quad (42)$$

Since there are total T singular values λ_i overall, we can equivalently rewrite (42) as

$$\#\{i : \lambda_i \geq \frac{1}{4n} \sum_{j=1}^n \log T\} \leq \frac{T}{2} \quad (43)$$

This immediately implies the lower bounds $\text{orkR}k_{HS}^2$ and $\text{dkR}k_{HS}^2$. \square

E Proof of Theorem 1

The two main probabilistic tools that we use are the Hanson-Wright inequality (Hanson and Wright 1971), and a classical norm deviation inequality for sub Gaussian vectors, as follows:

Theorem 7 (Hanson-Wright Inequality) Let $X = (X_1, \dots, X_m) \in \mathbb{R}^m$ be a random vector such that the components X_i are independent and $X_i \sim \text{SG}(\sigma_i)$ for all $i \leq m$. Let A be an $m \times m$ matrix. Then, for every $t \geq 0$,

$$\mathbb{P}(|\langle AX, X \rangle - \mathbb{E} \langle AX, X \rangle| > t) \leq 2 \exp \left(-c \min \left\{ \frac{t^2}{4 \|A\|_{\text{HS}}^2}; \frac{t}{\|A\|_{\text{op}}} \right\} \right) \quad (44)$$

In particular, recall that we are interested in concentration of $\|R_j\|^2$. We may write:

$$\|R_j\|^2 = \|R_{0,u} S_j + R_z\|^2 = \|R_{0,u} S_j\|^2 + \|R_z\|^2 + 2 \langle R_{0,u} S_j, R_z \rangle \quad (45)$$

The deviations of the first two terms may be bounded via Theorem 7. For the third term, we use the following:

Lemma 8. For any $X \sim \text{SG}_m(\sigma)$ we have:

$$\mathbb{P}(\|X\| > 4 \sqrt{m} \sigma + t) \leq \exp \left(-\frac{ct^2}{2} \right) \quad (46)$$

Note that in Lemma 8 we do not require the coordinates X_i to be independent. This will be important in what follows. Lemma 8 is standard and can be proved via covering number estimates of the Euclidean ball, see for instance Vershynin (2018), Section 4.4.

In addition, the following observation is used throughout the text:

Lemma 9. Let $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an operator, and let $h = (h_1, \dots, h_m)$ have independent coordinates with $\mathbb{E} h_i^2 = \sigma^2$. Denote by σ_i , $k = \min(m, n)$, the singular values of A . Then

$$\mathbb{E} \|Ah\|^2 = \sum_{i=1}^k \|A\sigma_i\|^2 \sigma^2 \quad (47)$$

The elementary proof is omitted.

We now prove Theorem 1.

Proof. Let $0 < \epsilon < 1$ be given. We first bound the deviations from the expectation for the first term in (45), $\|R_{0,u} S_j\|^2$. We apply the Hanson-Wright inequality with $h = h$ and $A = (R_{0,u} S_j) R_{0,u} S_j$. By definition, A has a single eigenvalue with multiplicity T . Thus clearly $\|A\|_{\text{HS}}^2 = T$ and $\|A\|_{\text{op}} = 1$.

For an appropriate constant $c > 0$, set $t = c \sqrt{2} \sqrt{T} \log \frac{1}{\epsilon}$. Then,

$$\mathbb{P}(\|R_{0,u} S_j\|^2 > T + c^2 \sqrt{2} \sqrt{T} \log \frac{1}{\epsilon}) \leq \epsilon \quad (48)$$

The deviation of the second term in (45) is similarly bounded using $A = R R$. Recall that by Theorem 3 we have $\|R\|_{\text{op}} \leq 2 \sqrt{j_{\min}}$, and note that $\|R\|_{\text{HS}}^2 = T \|R\|_{\text{op}}^4 \leq c T j_{\min}^4$. Set $t = c \sqrt{2} \sqrt{j_{\min}} \sqrt{2} \sqrt{T} \log \frac{1}{\epsilon}$. With this choice it follows that both terms in the minimum (44) are larger than $c \log \frac{1}{\epsilon}$ and we have

$$\mathbb{P}(\|R_z\|^2 > c^2 \|R\|_{\text{HS}}^2 + c^2 \sqrt{2} \sqrt{j_{\min}} \sqrt{2} \sqrt{T} \log \frac{1}{\epsilon}) \leq \epsilon \quad (49)$$

Finally, we bound the third term in (45). Denote by D the event

$$D = \{ \langle R_{0,u} S_j, R_z \rangle > c \sqrt{j_{\min}} \sqrt{2} \sqrt{T} \log \frac{1}{\epsilon} \} \quad (50)$$

and by E the event

$$E = \left\{ \sum_{j=1}^p |z_j| > c \sqrt{\frac{1}{T} + \frac{1}{\log T}} \right\} \quad (51)$$

By Lemma 8 applied to z , and using the fact that $\|z\|_2 \leq \sqrt{2} \sqrt{\frac{1}{T} + \frac{1}{\log T}}$,

$$P(D) \leq P\left(\sum_{j=1}^p |z_j| > c \sqrt{\frac{1}{T} + \frac{1}{\log T}}\right) \quad (52)$$

Next, using independence of z and $\mathcal{S}_{T,n}$,

$$P(E | D^c) \quad (53)$$

where D^c is the complement of D . Therefore, combining (52) and (53),

$$P(E) = P(D)P(E|D) + P(D^c)P(E|D^c) \quad (54)$$

Combining (48), (49) and (54), we obtain via the union bound:

$$P\left(\sum_{j=1}^p \frac{|jR^Y j|^2}{T} \geq 2 + \frac{kR^2 k_{HS}^2}{T} + c \frac{1 + \frac{1}{\log T} + \frac{1}{\log T} \log \frac{1}{\log T}}{\sqrt{T}}\right) \leq 4 \quad (55)$$

Similarly, we obtain a bound for the equations involving

$$P\left(\sum_{j=1}^p \frac{|jR^{0Y} j|^2}{p} \geq 2 + \frac{kR^0 k_{HS}^2}{p} + c \frac{1 + \frac{1}{\log T} + \frac{1}{\log T} \log \frac{1}{\log T}}{\sqrt{p}}\right) \leq 4 \quad (56)$$

The only difference in the derivation of (56) compared to (55) is in the application of Lemma 8. In the later case, to replace \bar{T} with \bar{p} in (52), we apply Lemma 8 with $X = P_V z$ rather than with $X = z$, where P_V is the projection onto the range of \bar{p} – a p -dimensional space. Note that z does not necessarily have a structure of independent coordinates, but is sub Gaussian and isotropic. Therefore Lemma 8 still applies. □

F Proof of Corollary 2

We now turn to prove Corollary 2.

Proof. Denote

$$E(a) = \left\{ \sum_{j=1}^p \frac{|jR^Y j|^2}{T} \geq 2 + \frac{kR^2 k_{HS}^2}{T} + c \frac{1 + \frac{1}{\log T} + \frac{1}{\log T} \log \frac{1}{\log T}}{\sqrt{T}} \right\} \quad (57)$$

Using (56) and (55) we may write

$$\frac{|jR^Y j|^2}{T} + e_1 = 2 + \frac{kR^2 k_{HS}^2}{T} + \epsilon_1 \quad (58)$$

$$\frac{|jR^{0Y} j|^2}{p} + e_2 = 2 + \frac{kR^0 k_{HS}^2}{p} + \epsilon_2 \quad (59)$$

where e_1, e_2 are error terms such that $\mathbb{P}(e_1) \leq \frac{1}{8}$ and $\mathbb{P}(e_2) \leq \frac{1}{8}$ holds with probability at least $1 - \frac{1}{8}$.

It follows that

$$\sum_{j=1}^p \frac{|jR^{0Y} j|^2}{p} = \sum_{j=1}^p \frac{|jR^Y j|^2}{T} - \frac{kR^0 k_{HS}^2}{p} + \frac{kR^2 k_{HS}^2}{T} + (e_2 - e_1) \frac{kR^0 k_{HS}^2}{p} - \frac{kR^2 k_{HS}^2}{T} \quad (60)$$

and

$$\frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} \frac{|R^{\alpha} Y|^2}{\rho} + \frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} e_2 = \frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} \rho^2 + \frac{k R k_{HS}^2}{T} \rho^2. \quad (61)$$

Thus

$$\begin{aligned} \rho^2 &= \rho \left[\frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} \frac{|R^{\alpha} Y|^2}{\rho} + \frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} e_2 \right] \\ &= \rho \left[\frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} \frac{|R^{\alpha} Y|^2}{\rho} + \frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} e_2 \right] \rho. \end{aligned} \quad (62)$$

It remains to observe that

$$\frac{k R \alpha_{HS}^2}{\rho} \frac{k R k_{HS}^2}{T} = \rho \frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} \frac{\rho}{k R \alpha_{HS}^2} \quad (63)$$

$$\leq \rho \frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} n^2 j u_{\max} j^2 \log^2 T; \quad (64)$$

where the inequality follows from eq(43) in the proof of Theorem 3. \square

G Thresholding Gap Analysis

The main result of this section is the following Proposition.

Proposition 10. Given the sequence $\{u_t\}_{t=1}^T \subset \mathbb{R}^n$, define the scalar sequence $\rho = \frac{1}{T} \sum_{t=1}^T |u_t|^2$ and set

$$j_{\min} = \min_t |u_t| \text{ and } j_{\max} = \max_t |u_t|. \quad (65)$$

Then for $\rho = \frac{1}{4} T$,

$$\frac{k R \alpha_{HS}^2}{\rho} \frac{k R k_{HS}^2}{T} \leq \frac{j_{\min}}{n j_{\max}} \frac{k R k_{HS}^2}{T}. \quad (66)$$

The general idea behind the proof of Proposition 10 is to show that for $\rho = 1$, the ratio

$$\frac{\rho k R k_{HS}^2}{T k R \alpha_{HS}^2} \text{ can be controlled. This is done in Lemmas 11 and 12 below. In particular,}$$

Lemma 11 is a general statement about integrals of monotone real functions under certain order constraints, and Lemma 12 provides a relation of the spectrum of R to that of S^{α} . It is then shown that for arbitrary n , $O_u S$ contains a certain copy of an $n = 1$ -dimensional operator with parameters u_t , which implies the bounds.

For the case $\rho = 1$, stated in Lemma 12, the argument consists of showing that the spectrum of R is upper and lower bounded by appropriately decaying functions, and therefore can not be “too constant”. We first obtain general estimates for the integrals of such upper and lower bounded functions in the following Lemma:

Lemma 11. Let $f : [0; 1] \rightarrow \mathbb{R}$ be a monotone non-increasing function such that for $x \in [0; 1]$,

$$(1-x)^2 \leq f(x) \leq M(1-x)^2; \quad (67)$$

for some $M \geq 1$. Set $t_0 = \frac{1}{4}$. Then

$$r(f) := \frac{\int_{t_0}^1 f(x) dx}{\int_0^1 f(x) dx} \leq 1 + \frac{c}{M}. \quad (68)$$

Proof. Denote

$$I(a; b; f) = \int_a^b f(x) dx; \quad (69)$$

Write

$$r(f) = \frac{t_0^{-1} l(0; t_0; f)}{l(0; t_0; f) + l(t_0; 1; f)} = \frac{t_0^{-1}}{1 + \frac{l(t_0; 1; f)}{l(0; t_0; f)}} \quad (70)$$

and set $v := f(t_0)$. Then, among all f that satisfy (67) and $f(t_0) = v$, $r(f)$ is minimized iff f with maximal $l(t_0; 1; f)$ and minimal $l(0; t_0; f)$. Due to the form of the constraint (67) and monotonicity, this minimizer is given by

$$f_v(x) = \begin{cases} (1-x)^2 & x \in [0; t_v] \\ v & x \in [t_v; t_v^+] \\ M(1-x)^2 & x \in [t_v^+; 1]; \end{cases} \quad (71)$$

where $t_v := \max\{0; 1 - \sqrt{\frac{v}{M}}\}$ and $t_v^+ := 1 - \sqrt{\frac{v}{M}}$. Our problem is therefore now reduced from a minimization of $r(f)$ over the function space to a problem of minimizing $r(f_v) := r(v)$ over a single scalar variable v .

To this end, first note that we can assume w.l.o.g. that $(1 - t_0)^2 \leq v$. Indeed, for large M , the lower bound on $r(f)$ can only become smaller, since f would be minimized over a larger set. By construction, the value v must satisfy $(1 - t_0)^2 \leq v \leq M(1 - t_0)^2$, and for $M > (1 - t_0)^{-2}$, the value $v = 1$ satisfies these inequalities.

Next, by direct computation (taking the derivative) one verifies that for any $M > (1 - t_0)^{-2}$, $r(v)$ as a function of v is minimized at $v = 1 - \frac{1}{3M}$, which yields the statement of the Lemma. \square

We can now treat the $\theta = 1$ case.

Lemma 12. Consider the case $\theta = 1$, and $p = \frac{1}{4}T$. Then

$$\frac{kR^0 k_{HS}^2}{p} \leq \frac{kRk_{HS}^2}{T} \leq c \frac{j u_{\min} j}{j u_{\max} j} \leq \frac{kRk_{HS}^2}{T}. \quad (72)$$

Proof. For any two operators A, B we have

$$i(AB) \leq i(A) k_B k_{op} \text{ and } i(BA) \leq i(A) k_B k_{op}; \quad (73)$$

see Bhatia (1997); Gohberg and Krein (1969). Note that in the case $\theta = 1$ and $j u_{\min} j > 0$, $O_u S$ is invertible. Thus the singular values $\sigma_i(S^0)$ satisfy

$$j u_{\min} j \leq \sigma_i(S^0) \leq i(O_u S^0) \leq j u_{\max} j \leq \sigma_i(S^0); \quad (74)$$

where the first inequality follows by applying (73) to $(O_u S^0)^{-1}$, and the second by considering $O_u S^0$ itself. Equivalently, for $i = 1; \dots; T$

$$j u_{\max} j^{-1} \leq \sigma_{T+1-i}(S^0)^{-1} \leq \sigma_{T+1-i}(R) \leq j u_{\min} j^{-1} \leq \sigma_{T+1-i}(S^0)^{-1}; \quad (75)$$

and using Lemma 8,

$$c \frac{j u_{\max} j^{-1} i}{T+1} \leq \sigma_{T+1-i}(R) \leq j u_{\min} j^{-1} \frac{j u_{\min} j^{-1} i}{T+1}. \quad (76)$$

By changing the index and taking squares, we have

$$c j u_{\max} j^{-2} \leq \frac{i}{T+1} \leq \sigma_i^2(R) \leq c^0 j u_{\min} j^{-2} \leq \frac{i}{T+1}. \quad (77)$$

Now, using an elementary discretization argument for $\frac{1}{4}T$ by Lemma 11 it follows that

$$\frac{kR^0 k_{HS}^2}{p} \leq \frac{kRk_{HS}^2}{T} \leq 1 + c \frac{j u_{\max} j}{j u_{\min} j}, \quad (78)$$

which implies (72). \square

Finally, we prove Proposition 10.

Proof. As discussed earlier, the key to a statement such as (6) is to show that the spectrum of $O_u S$ is non-constant, and in particular has enough singular values. This is equivalent to providing appropriate lower bounds on the spectrum of $O_u S$. Here we derive such bounds by comparison with the $n = 1$ case as follows: Let $V \subset \mathbb{R}^n$ be a T -dimensional subspace spanned by vectors for which for every time t , all coordinates at time t are identical. Formally, for $x \in V$, for every t we require that $x_{(t-1)n+i} = x_{(t-1)n+j}$ for all $i, j \leq T$.

Observe that the restriction of $O_u S$ to V , the operator $O_u S|_V$, acts equivalently to the $n = 1$ case operator defined by $O_u S^0$. In particular, similarly to the argument in Lemma 12, it follows that

$$\lambda_{\min}^j(O_u S^0) \leq \lambda_{\min}^j(O_u S|_V). \quad (79)$$

Moreover, note that $(O_u S|_V)^T O_u S|_V = P_V S^T O_u O_u S P_V$ and $(O_u S)^T O_u S$ are non-negative operators and $(O_u S)^T O_u S - (O_u S|_V)^T O_u S|_V$ (that is, $(O_u S)^T O_u S - (O_u S|_V)^T O_u S|_V$ is non-negative). It follows that for all T ,

$$\lambda_{\min}^j(O_u S) \geq \lambda_{\min}^j(O_u S|_V). \quad (80)$$

The upper bounds on the spectrum of $(O_u S)$ may again be obtained via (73). Indeed,

$$\lambda_{\max}^j(O_u S) \leq \lambda_{\max}^j(O_u S|_V) + \lambda_{\max}^j(d_n^\perp e(S^0)); \quad (81)$$

where the first inequality follows from (73) while the second is due to the multiplicity of each singular value of S^0 in S . The rest of the argument proceeds as in Lemma 12. \square

H Missing Values in STVE

We now discuss the treatment of missing values in STVE. Recall that our starting point is the vector form of the system, (7), which we rewrite here:

$$Y = O_u S h + z; \quad (82)$$

where $h \in \mathbb{R}^n$ and $z \in \mathbb{R}^T$ are sub-Gaussian vectors and $y = (y_1, \dots, y_T) \in \mathbb{R}^T$ is the observation vector. Suppose that M out of T observation values are missing, at times t_1, \dots, t_M . Set $T^0 = T - M$ and define a projection operator $P_A : \mathbb{R}^T \rightarrow \mathbb{R}^T$ as the operator that omits the coordinates t_1, \dots, t_M . Formally, let $e(t), t \in T$, be the standard basis vector e_t , with 1 at coordinate t and zeros elsewhere. Then

$$P_A(e(t)) = \begin{cases} 0 & \text{if } t \in \{t_1, \dots, t_M\}; \\ e(t) & \text{otherwise.} \end{cases} \quad (83)$$

We can then rewrite (82) as

$$P_A Y = P_A O_u S h + P_A z; \quad (84)$$

Note that the vector $P_A Y$ contains only available, non-missing values of Y . Similarly to the case with no missing values, we define R as the Moore-Penrose inverse of $P_A O_u S$. Then we have

$$R P_A Y = R P_A O_u S h + R P_A z; \quad (85)$$

Note that since a Moore-Penrose inverse only acts on the image of $O_u S$, we have $R P_A = R$ and therefore

$$R Y = R O_u S h + R z \quad (86)$$

similarly to the case with no missing values. The only difference here will be that R now have T^0 rather than T non-zero singular values. R^0 will be defined similarly. The analysis of the spectrum of R concerns only non-zero singular values of $O_u S$ and holds with no change other than that should be replaced by T^0 . Consequently, the whole approach works identically with T^0 replaced by T everywhere, and in the bounds of Theorem 3 in particular.

