
On Multivariate Singular Spectrum Analysis and its Variants

Anish Agarwal*
MIT
anish90@mit.edu

Abdullah Alomar
MIT
aalomar@mit.edu

Devavrat Shah
MIT
devavrat@mit.edu

Abstract

We introduce and analyze a simpler, practically useful variant of multivariate singular spectrum analysis (mSSA), a known time series method to impute (or de-noise) and forecast a multivariate time series. Towards this, we introduce a spatio-temporal factor model to analyze mSSA. This model includes the usual components used to model dynamics in time series analysis such as trends (low order polynomials), seasonality (finite sum of harmonics) and linear time-invariant systems. We establish that given N time series and T observations per time series, the in-sample prediction error for both imputation and forecasting under mSSA scales as $1/\sqrt{\min(N, T)T}$. This is an improvement over: (i) the $1/\sqrt{T}$ error scaling of SSA, which is the restriction of mSSA to univariate time series; (ii) the $1/\min(N, T)$ error scaling for Temporal Regularized Matrix Factorized (TRMF), a matrix factorization based method for time series prediction. That is, mSSA exploits both the ‘temporal’ and ‘spatial’ structure in a multivariate time series. Our experimental results using various benchmark datasets confirm the characteristics of the spatio-temporal factor model and our theoretical findings—our variant of mSSA empirically performs as well or better compared to popular neural network based time series methods, LSTM and DeepAR.

*All authors are with Massachusetts Institute of Technology during the course of this work. Their affiliations include Department of EECS, LIDS, IDSS, Statistics and Data Science Center, and CSAIL. Their email addresses: {anish90, aalomar, devavrat}@mit.edu.

Contents

1	Introduction	4
1.1	Multivariate Singular Spectrum Analysis	4
1.2	Our Contributions	6
1.3	Literature Review	8
2	Model	9
2.1	Spatio-Temporal Factor Model	9
2.2	Comparison with Existing Time Series Models	10
2.3	A Diagnostic Test for the Spatio-Temporal Model	10
3	Main Results	11
4	Approximately Low-Rank Hankel Representation	12
4.1	Extending Model and Main Result	12
4.2	Hankel Calculus	13
4.3	Examples of (G, ϵ) -Hankel Time Series	13
5	Experiments	14
5.1	Imputation	14
5.2	Forecasting	16
6	Algorithmic Extensions of mSSA	16
6.1	Variance Estimation	16
6.2	Tensor SSA	19
7	Conclusion	20
A	Page vs. Hankel mSSA	23
B	Experiment Details	24
B.1	Datasets	25
B.2	Algorithms.	26
B.3	Parameters Selection	27
C	Time-varying Recommendation Systems	27
D	Proof of Proposition 4.1	28
E	Proofs For Section 4.2	29
E.1	Proof of Proposition 4.3	29
E.2	Proof of Proposition 4.4	30
E.3	Proof of Proposition 4.5	30

E.3.1	Helper Lemmas for Proposition 4.5	30
E.3.2	Completing Proof of Proposition 4.5	31
E.4	Proof of Proposition 4.6	31
F	Concentration Inequalities	33
G	Matrix Estimation via HSVT	33
G.1	Setup, Notations	33
G.2	Matrix Estimation using HSVT	34
G.3	A Useful Linear Operator	34
G.4	HSVT based Matrix Estimation: A Deterministic Bound	35
G.5	HSVT based Matrix Estimation: Deterministic To High-Probability	37
H	Proof of Theorem 4.1	41
I	Proof of Theorem 4.2	42
I.1	Proof of Proposition 4.2	45
J	Proof of Theorem 6.1	46
K	tSSA Proofs	49
K.1	Proof of Proposition 6.1	49
K.2	Proof of Proposition 6.2	49
K.3	Proof of Proposition C.1	50

1 Introduction

Multivariate time series data is of great interest across many application areas, including cyber-physical systems, finance, retail, healthcare to name a few. The goal across these domains can be summarized as accurate imputation and forecasting of a multivariate time series in the presence of noisy and/or missing data along with providing meaningful uncertainty estimates.

Setup. We consider a discrete time setting with time indexed as $t \in \mathbb{Z}$. For $N \in \mathbb{N}$, let $f_n : \mathbb{Z} \rightarrow \mathbb{R}$, $n \in [N] := \{1, \dots, N\}$ be the latent time series of interest. For $t \in [T]$ and $n \in [N]$, we observe $X_n(t)$ where for $\rho \in (0, 1]$,

$$X_n(t) = \begin{cases} f_n(t) + \eta_n(t) & \text{with probability } \rho \\ \star & \text{with probability } 1 - \rho. \end{cases} \quad (1)$$

Here \star represents a missing observation and $\eta_n(t)$ represents the per-step noise, which we assume to be an independent (across t, n) mean-zero random variable. Though $\eta_n(t)$ is independent, it is worth noting though the underlying time series, $f_n(\cdot)$, is of course strongly dependent across t, n . In fact, typically in time series analysis, one models structure in ‘time’ in the latent time series of interest, $f_n(\cdot)$, as composed of some mixture of a *trend*, a *seasonal*, and a *stationary* component. The typical model for *trend* is a low-order polynomial, for *seasonality* is a finite sum of harmonics, and for *stationarity* is a (variant of) autoregressive or linear time-invariant systems. Further, the dependence in ‘space’, i.e. across the multivariate components of a time series is typically captured through linear dependence. To capture both the temporal and spatial structure in a multivariate time series, we utilize a generic spatio-temporal factor model as described in detail in Section 2. As discussed in Section 2.2, it naturally captures the *trend*, *seasonality*, and *stationary* components in time as well as linear dependence in space.

We emphasize that the presence of per-step noise $\eta_n(t)$ and missing values (denoted by \star) represent an additional challenge of measurement error in our setup. The generic spatio-temporal factor model for $f_n(\cdot), n \in [N]$ described in Section 2 *without* additional noise $\eta_n(\cdot)$ or missingness already provides an expressive model for the various aspects of time series model analysis (trend, periodicity, stationarity).

Goal. The objective is two-folds, for $n \in [N]$: (i) imputation – estimating $f_n(t)$ for all $t \in [T]$; (ii) forecasting – learning a model to forecast $f_n(t)$ for $t > T$.

1.1 Multivariate Singular Spectrum Analysis

Multivariate singular spectrum analysis (mSSA) is a known method to impute and forecast a multivariate time series. The primary objective of this work is to introduce and theoretically analyze a simpler, practically useful variant of mSSA. Given the simplicity of the variant of mSSA we introduce, we start by describing it in detail below. In Section 1.3, we compare the original mSSA method with this variant and discuss key differences. See Figure 1 for a visual depiction of the key steps in mSSA.

Singular spectrum analysis (SSA). For ease of exposition and to build intuition, we start with $N = 1$, i.e. a univariate time series. There are two algorithmic parameters: $L \geq 1$ and $k \geq 1$. For simplicity and without loss of generality², assume that T is an integer multiple of L , i.e. $T/L \in \mathbb{N}$ and $k \leq \min(L, T/L)$. First, transform the time series $X_1(t)$, $t \in [T]$ into an $L \times T/L$ matrix where the entry of the matrix in row $i \in [L]$ and column $j \in [T/L]$ is $X_1(i + (j - 1) \times L)$. This matrix induced by the time series is called the Page matrix, and we denote it as $P(X_1, T, L)$.

²When $T/L \notin \mathbb{N}$, by applying both imputation and forecasting algorithm for two ranges, $1, \dots, \lfloor T/L \rfloor \times L$ and $(T \bmod L) + 1, \dots, T$, it will satisfy the assumed condition and will provide imputation and forecasting for the entire range.

Imputation. After replacing missing values (i.e. \star) in the matrix $P(X_1, T, L)$ by 0, we compute its singular value decomposition, which we denote as

$$P(X_1, T, L) = \sum_{\ell=1}^{\min(L, T/L)} s_{\ell} u_{\ell} v_{\ell}^T,$$

where $s_1 \geq s_2 \geq \dots \geq s_{\min(L, T/L)} \geq 0$ denote its ordered singular values, and $u_{\ell} \in \mathbb{R}^L, v_{\ell} \in \mathbb{R}^{T/L}$ denote its left and right singular vectors, respectively, for $\ell \in [\min(L, T/L)]$. Let $\hat{\rho}_1$ be the fraction of observed entries of X_1 , precisely defined as $(\max(1, \sum_{t=1}^T \mathbf{1}(X_1(t) \neq \star)))/T$. Let the normalized, truncated version of $P(X_1, T, L)$ be

$$\hat{P}(X_1, T, L; k) = \frac{1}{\hat{\rho}_1} \sum_{\ell=1}^k s_{\ell} u_{\ell} v_{\ell}^T, \quad (2)$$

i.e., we perform Hard Singular Value Thresholding (HSVT) on $P(X_1, T, L)$ to obtain $\hat{P}(X_1, T, L; k)$. We then define the *de-noised and imputed estimate* of the original time series, denoted by \hat{f}_1 , as: for $t \in [T]$, $\hat{f}_1(t)$ equals the entry of $\hat{P}(X_1, T, L; k)$ in row $(t - 1 \bmod L) + 1$ and column $\lceil t/L \rceil$.

Forecasting. To forecast, we learn a linear model $\hat{\beta}(X_1, T, L; k) \in \mathbb{R}^{L-1}$, which is the solution to

$$\text{minimize} \quad \sum_{m=1}^{T/L} (y_m - \beta^T x_m)^2 \quad \text{over} \quad \beta \in \mathbb{R}^{L-1},$$

where $y_m = (1/\hat{\rho}_1)X_1(L \times m)$ and $x_m = [\hat{f}_1(L \times (m-1) + 1) \dots \hat{f}_1(L \times (m-1) + L - 1)]^3$ for $m \in [T/L]$. Note, to define y_m , we impute missing values in X_1 by 0. Then the *forecasted estimate* at time $t = L \times m$ is given by $\hat{f}_1(L \times m) = \hat{\beta}(X_1, T, L; k)^T x_m$ for $m \in [T/L]$.

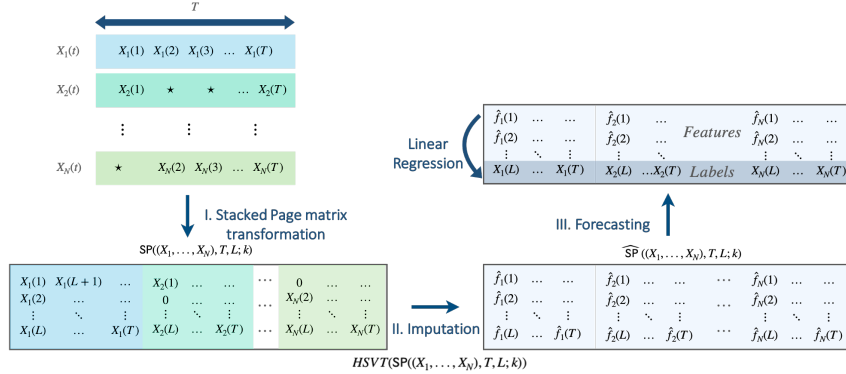


Figure 1: Key steps of our proposed variant of the mSSA algorithm.

Multivariate singular spectrum analysis (mSSA). Below we describe the variant of mSSA we propose, which is an extension of the SSA algorithm described above, to when we have a multivariate time series, i.e., $N > 1$. The key change is in the first step where we construct the Page matrix—instead of considering the Page matrix of a single time series, we now consider a ‘stacked’ Page matrix, which is obtained by a column-wise concatenation of the Page matrices induced by each time series separately. Specifically, like SSA, it has two algorithmic parameters, $L \geq 1$ and $k \geq 1$. For each time series, $n \in [N]$, create its $L \times T/L$ Page matrix $P(X_n, T, L)$, where the entry in row $i \in [L]$ and column $j \in [T/L]$ is $X_n(i + (j - 1) \times L)$. We then create a stacked Page matrix from these N time series by performing a column wise concatenation of the N matrices,

³To establish theoretical results for the forecasting algorithm, we produce estimates $(\hat{f}_1(L \times (m-1) + 1) \dots \hat{f}_1(L \times (m-1) + L - 1))$ for $m \in [T/L]$ by applying the imputation algorithm on $P(X_1, T, L)$ after setting its L th row equal to 0. Also, $\hat{\rho}_1$ in the definition of y_m is computed using only the first $L - 1$ rows of $P(X_1, T, L)$. This avoids dependences in the noise between y_m and x_m for $m \in [T/L]$.

$P(X_n, T, L)$, $n \in [N]$. We denote this matrix as $\text{SP}((X_1, \dots, X_N), T, L)$, and note that it has L rows and $N \times T/L$ columns.

Imputation. We replace missing values (i.e. \star s) in $\text{SP}((X_1, \dots, X_N), T, L)$ by 0. Similar to (2), we perform HSVT on $\text{SP}((X_1, \dots, X_N), T, L)$ and denote its normalized, truncated version as $\widehat{\text{SP}}((X_1, \dots, X_N), T, L; k)$ (instead of $\widehat{\rho}_1$, we now normalize by $\widehat{\rho} := (\max(1, \sum_{n=1}^N \sum_{t=1}^T \mathbf{1}(X_n(t) \neq \star)))/NT$). From $\widehat{\text{SP}}((X_1, \dots, X_N), T, L; k)$, like in SSA, we can read off $\hat{f}_n(t)$ for $n \in [N]$, $t \in [T]$, the *de-noised and imputed estimate* of the N time series over T time steps. In particular, let $\widehat{P}(X_n, T, L; k)$ refer to sub-matrix of $\widehat{\text{SP}}((X_1, \dots, X_N), T, L; k)$ induced by selecting only its $[(n-1) \times (T/L) + 1, \dots, n \times T/L]$ columns. Then for $t \in [T]$, $\hat{f}_n(t)$ equals the entry of $\widehat{P}(X_n, T, L; k)$ in row $(t-1 \bmod L) + 1$ and column $\lceil t/L \rceil$.

Forecasting. Similar to SSA, to forecast, we learn a linear model $\hat{\beta}((X_1, \dots, X_N), T, L; k) \in \mathbb{R}^{L-1}$, which is the solution to

$$\text{minimize} \quad \sum_{m=1}^{N \times T/L} (y_m - \beta^T x_m)^2 \quad \text{over} \quad \beta \in \mathbb{R}^{L-1}, \quad (3)$$

where y_m is the m th component of $(1/\widehat{\rho})[X_1(L), X_1(2 \times L), \dots, X_1(T), X_2(L), \dots, X_2(T), \dots, X_N(T)] \in \mathbb{R}^{N \times T/L}$, and $x_m \in \mathbb{R}^{L-1}$ corresponds to the vector formed by the entries of the first $L-1$ rows in the m th column of $\widehat{\text{SP}}((X_1, \dots, X_N), T, L; k)$ ⁴ for $m \in [N \times T/L]$. Note, to define y_m , we impute missing values in X_1, \dots, X_N by 0. Then the *forecasted estimate* at time $t = L \times m'$ for $m' \in [T/L]$ for time series $n \in [N]$ is $\hat{f}_n(L \times m') = \hat{\beta}((X_1, \dots, X_N), T, L; k)^T x_m$ where $m = m' + (n-1) \times T/L$. Recall, Figure 1 has a visual depiction of the key steps above.

Empirical performance of mSSA. This variant of the mSSA algorithm, fully described above, is arguably quite simple, with its major steps consisting of simply singular value thresholding and ordinary least squares. A key question is – how well does it perform empirically? In Table 1, we provide a summary comparison of mSSA’s performance for imputation and forecasting on benchmark datasets with respect to state-of-the-art time series algorithms. We find that by using the stacked Page matrix in mSSA, it greatly improves performance over SSA; thus, indicating that mSSA is effectively utilizing information *across* multiple time series. Further, it performs competitively or outperforms neural network based methods, such as LSTM and DeepAR⁵. However, aside from the recent theoretical analysis of SSA in [1, 3], theoretical explanation of the success of mSSA has been absent. This begs the question:

When and why does mSSA work?

1.2 Our Contributions

As our primary contribution, we provide an answer to the question posed above – under a spatio-temporal factor model that we introduce, the finite-sample analysis we carry out of mSSA’s estimation error for imputation and forecasting establishes its asymptotic consistency as well as its ability to effectively utilize both the spatial and temporal structure in a multivariate time series. Below, we detail the various aspects of our contribution with respect to the: (a) spatio-temporal factor model; (b) finite sample analysis of mSSA; (c) algorithmic extensions (and associated theoretical analysis) of mSSA to do time-varying variance estimation, and a tensor variant of SSA which we show can have a better imputation prediction error convergence rate compared to mSSA in certain regimes.

⁴Similar to the SSA forecasting algorithm, when creating a forecasting model in mSSA, we produce $\widehat{\text{SP}}((X_1, \dots, X_N), T, L; k)$ by first setting the L th row of $\text{SP}((X_1, \dots, X_N), T, L; k)$ equal to zero before performing the SVD and the subsequent truncation. Also, $\widehat{\rho}$ in the definition of y_m is computed only using the first $L-1$ rows of $\text{SP}((X_1, \dots, X_N), T, L; k)$.

⁵It is worth noting that for each benchmark dataset, the effective rank (see Section 2.3 for the definition of effective rank) of the stacked Page matrix effectively explains the relative performance of mSSA; the lower the effective rank, the better mSSA’s relative performance. See Section 2.3 and Table 3 for details.

Table 1: mSSA statistically outperforms SSA, other state-of-the-art algorithms, including LSTMs and DeepAR across many datasets. We use the average normalized root mean squared error (NRMSE) as our metric. Details of experiments run to produce results can be found in Section 5.

	Mean Imputation (NRMSE)					Mean Forecasting (NRMSE)				
	Electricity	Traffic	Synthetic	Financial	M5	Electricity	Traffic	Synthetic	Financial	M5
mSSA	0.398	0.508	0.416	0.238	0.883	0.485	0.536	0.281	0.251	1.021
SSA	0.514	0.713	0.675	0.467	0.958	0.632	0.696	0.665	0.303	1.068
LSTM	NA	NA	NA	NA	NA	0.558	0.478	0.559	1.205	1.034
DeepAR	NA	NA	NA	NA	NA	0.479	0.464	0.415	0.316	1.050
TRMF	0.641	0.460	0.564	0.430	0.916	0.495	0.508	0.422	0.291	1.032
Prophet	NA	NA	NA	NA	NA	0.569	0.614	1.010	1.286	1.010

Table 2: Comparison of finite-sample results with relevant algorithms in the literature.

Method	Functionality	Mean Estimation	
	Multivariate time series	Imputation	Forecasting
This Work	Yes	$1/\sqrt{\min(N, T)T}$	$1/\sqrt{\min(N, T)T}$
mSSA - Literature	Yes	—	—
SSA [11, 1]	No	$T^{-1/4}$	—
Neural Network [21, 7]	Yes	—	—
TRMF [34, 19]	Yes	$(\min(N, T))^{-1}$	—

Spatio-temporal factor model. Note that the collection of latent multivariate time series $f_n(t)$, for $n \in [N]$, $t \in [T]$ can be collectively viewed as a $N \times T$ matrix (in comparison, the stacked Page matrix used in mSSA is of dimension $L \times (N \times T/L)$, where L is a hyper-parameter). To capture the spatial structure, i.e. the relationship across rows, we model this matrix to be low-rank — there exists a low-dimensional latent factor (or feature) associated with each of N time series; analogously, there exists a low-dimensional latent factor associated with each of the T time steps. To capture the temporal structure, we further assume that each component of the latent temporal factor has a (approximately) *low-rank Hankel matrix*⁶ representation, i.e., the Hankel (and Page) matrix induced by each component of the latent temporal factor is (approximately) low-rank. For $N = 1$ (and hence rank = 1), this subsumes the model considered to explain the success of SSA in [1] as a special case. The model considered in [34] is similar in spirit to our spatio-temporal model, but they impose a specific linear structure for the temporal factors that is explicitly *known a priori* (as opposed to just requiring the existence of a latent representation as we do); the particular structure assumed affects the algorithm implementation in [34]. As stated earlier, the traditional modeling approach in time series analysis posits that a time series is a mixture of a trend (low-order polynomial), a seasonal (finite sum of harmonics), and a stationary (linear time-invariant function) component. We show that each of these components (low-order polynomial, harmonics, and linear time invariant function) indeed has a low-rank Hankel (or Page) matrix representation. Further, we establish that the set of time series that have a (approximately) low-rank Hankel representation is closed under component-wise addition and multiplication. Such a model *calculus* helps characterize the representational strength of the spatio-temporal factor model we introduce. It provides a natural generalization of the factor model considered in the panel data literature within econometrics by incorporating the temporal structure within the data; explicitly modeling the temporal structure is traditionally missing from this literature (see [29] for a textbook reference).

Finite sample analysis of mSSA. Under the spatio-temporal factor model, we establish that the imputation and (in-sample) forecasting prediction error scale as $1/\sqrt{\min(N, T)T}$, see Theorems 3.1 and 3.2. For $N = 1$, it implies that the SSA algorithm described above has imputation and forecasting error scaling as $1/\sqrt{T}$. That is, mSSA improves performance by a \sqrt{N} factor over SSA by utilizing information across the N time series. This also improves upon the prior work of [1] which established the weaker result that SSA has imputation error scaling as $1/T^{1/4}$ —also, [1]

⁶See Definition 2.1 for the Hankel matrix induced by a time series.

does not establish a result for the forecasting error of SSA. Further, existing matrix estimation or completion based methods applied to the $N \times T$ matrix of time series observations (i.e, without first performing the Page matrix transformation as done in mSSA) establish that the imputation prediction error scales as $1/\min(N, T)$. This is indeed the primary result of the works [34, 19] (precisely, see Theorem 2 of [19]⁷). That is, while the algorithm stated in [34, 19] utilizes the temporal structure in addition to the spatial structure, the theoretical guarantees do not reflect it. This explains why the guarantees provided by such methods are weaker (since $1/\min(N, T) \geq 1/\sqrt{\min(N, T)T}$) than that obtained by mSSA. See Table 2 for a summary of our theoretical results.

Extensions: variance and tensor SSA (tSSA).

First, we extend mSSA to estimate the latent time-varying variance, i.e. $\mathbb{E}[\eta_n^2(t)]$, $n \in [N]$, $t \in [T]$. We establish the efficacy of such an extension when the time-varying variance is also modeled through a spatio-temporal factor model. This, for example, enables meaningful uncertainty quantification for the estimation of the mean produced by mSSA. To the best of our knowledge, ours is the first result that provides provable finite-sample performance guarantees for estimating the time-varying variance of a time series. Second, we propose a tensor variant of SSA, termed tSSA, which exploits recent developments in tensor estimation. In tSSA, rather than doing a column-wise stacking of the Page Matrices induced by each of the N time series to form a larger matrix, we instead view each Page matrix as a slice of a $L \times T/L \times N$ order-three tensor. In other words, the entry of the tensor with indices $i \in [L]$, $j \in [T/L]$ and $n \in [N]$ equals the entry of $P(X_n, L, T)$ with indices i, j . In Proposition 6.2, with respect to imputation error, we characterize the relative performance of tSSA, mSSA, and “vanilla” matrix estimation (ME). We find that when $N = o(T^{1/3})$, mSSA outperforms tSSA; when $T^{1/3} = o(N)$, $N = o(T)$ tSSA outperforms mSSA; when $T = o(N)$, standard matrix estimation methods are equally as effective as mSSA and tSSA. See Figure 2 for a graphical depiction;

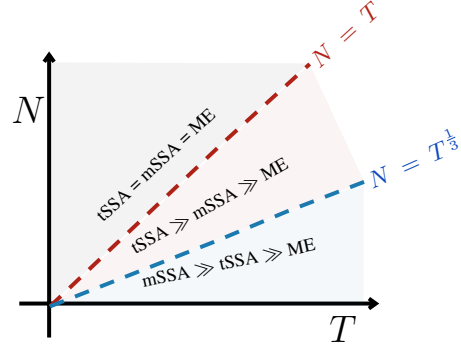


Figure 2: Relative effectiveness of tSSA, mSSA, ME for varying N, T .

1.3 Literature Review

Given the ubiquity of multivariate time series analysis, it will not be possible to do justice to the entire literature. We focus on a few techniques, most relevant to compare against, either theoretically or empirically.

SSA and mSSA. A good overview of the literature on SSA can be found in [11]. As alluded to earlier, the original SSA method differs from the variant discussed in [1] and in this work. The key steps of the original SSA method are: Step 1—create a Hankel matrix from the time series data; Step 2—do a Singular Value Decomposition (SVD) of it; Step 3—group the singular values based on user belief of the model that generated the process; Step 4—perform diagonal averaging to “Hankelize” the grouped rank-1 matrices outputted from the SVD to create a set of time series; and Step 5—learn a linear model for each “Hankelized” time series for the purpose of forecasting. The theoretical analysis of this original SSA method has been focused on proving that many univariate time series have a low-rank Hankel representation, and secondly on defining sufficient *asymptotic* conditions for when the singular values of the various time series components are separable, thereby justifying Step 3 of the method. Step 3 of the original SSA method requires user input and Steps 4 and 5 are not robust to noise and missing values due to the strong dependence across entries of the Hankel representation of the time series. To overcome these limitations, in [1] a simpler and practically useful version as described in Section 1.1 was introduced and a finite-sample analysis of it was done. As discussed earlier, this work improves upon the analysis of [1] by providing stronger bounds for imputation prediction error, and gives new bounds for forecasting prediction error, which were missing in [1]. The original mSSA method, like the original SSA method, involves the five

⁷There seems to be a typo in Corollary 2 of [34] in applying Theorem 2: square in Frobenius-norm error is missing.

steps described above, but first the Hankel matrices induced by each of the N time series are stacked either column-wise (horizontal mSSA) or row-wise (vertical mSSA) [15]. The theoretical analysis of the original mSSA is absent in the literature despite its empirical success and popularity (see [14, 13, 18]). In this work, as described in Section 1.1, we introduce and analyze an arguably simpler variant of mSSA that uses the Page (instead of the Hankel) matrix representation. In Appendix A, we compare our variant to the original variants of mSSA in both theoretical and practical aspects.

Matrix estimation based multivariate time series methods. There is a recent line of work in time series analysis (see [28, 34]), where multiple time series are viewed collectively as a matrix and some form of matrix factorization is done. Most such methods make strong prior model assumptions on the underlying time series and the algorithm changes based on the assumptions made on the time series dynamics that generated the data. Further, finite sample analysis, especially with respect to forecasting error, of such methods is usually lacking. We highlight one method, Temporal Regularized Matrix Factorization (TRMF) (see [34]), which we directly compare against due to its popularity, and as it achieves state-of-the-art imputation and forecasting empirical performance. The authors in [34] provide finite sample imputation analysis for an instance of the model considered in this work, but forecasting analysis is absent. As discussed earlier, it establishes that imputation error scales as $1/\min(N, T)$. This is a consequence of the low-rank structure of the time series matrix. But they fail to utilize, at least in the theoretical analysis, the temporal structure. Indeed, our analysis captures such temporal structure and hence our imputation error scales as $1/\sqrt{\min(N, T)T}$ which is a stronger guarantee. For example, for $N = \Theta(1)$, their error bound remains $\Theta(1)$ for any T , suggesting that TRMF [34] fails to utilize the temporal structure for better estimation, while mSSA does.

Other relevant literature. We take a brief note of popular time series methods in the recent literature. In particular, recently neural network (NN) based approaches have been popular and empirically effective. Some industry standard neural network methods include LSTMs, from the Keras library (a standard NN library, see [7]) and DeepAR (an industry leading NN library for time series analysis, see [21]).

2 Model

2.1 Spatio-Temporal Factor Model

Below, we introduce the spatio-temporal factor model we use to explain the success of mSSA. In short, the model requires that the underlying latent multivariate time series satisfies Properties 2.1 and 2.2, which capture the “spatial” and “temporal” structure within it, respectively. To that end, consider the matrix $\mathbf{M} \in \mathbb{R}^{N \times T}$, where its entry in row n and column t , M_{nt} is equal to $f_n(t)$, the value of the latent time series n at time t . We posit that this matrix \mathbf{M} is low-rank. Precisely,

Property 2.1. *Let $\text{rank}(\mathbf{M}) = R$. That is, for any $n \in [N], t \in [T]$, $M_{nt} = \sum_{r=1}^R U_{nr} W_{rt}$, where $|U_{nr}| \leq \Gamma_1, |W_{rt}| \leq \Gamma_2$ for constants $\Gamma_1, \Gamma_2 > 0$.*

Property 2.1 effectively captures the “spatial” structure amongst the N time series. The R dimensional latent factor $U_{n\cdot} \in \mathbb{R}^R$ characterizes the n -th time series; as a result, most (if not all) of the n time series can be expressed as a weighted linear combination of at most R other time series. This however does not capture the “temporal” structure within each time series. To do so, we impose additional structure on the latent factor $W_{\cdot t} \in \mathbb{R}^R$ associated with each time step $t \in [T]$. To that end, we introduce the notion of the Hankel matrix induced by a time series.

Definition 2.1 (Hankel Matrix). *Given a time series $f : \mathbb{Z} \rightarrow \mathbb{R}$, its Hankel matrix associated with observations over T time steps, $\{1, \dots, T\}$, is given by the matrix $H \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ with $H_{ij} = f(i + j - 1)$ for $i, j \in [\lfloor T/2 \rfloor]$.*

Now, for a given $r \in [R]$, consider the time series W_{rt} for $t \in [T]$. Let $H(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ denote its Hankel matrix restricted to $[T]$, i.e. $H(r)_{ij} = W_{r(i+j-1)}$ for $i, j \in [\lfloor T/2 \rfloor]$.

Property 2.2. *For each $r \in [R]$ and for any $T \geq 1$, the Hankel Matrix $H(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ associated with time series $W_{rt}, t \in [T]$ has rank at most G .*

Property 2.2 captures the temporal structure within the latent factors associated with time. As described in (1), for each $n \in [N]$ and $t \in [T]$, we observe $f_n(t) + \eta_n(t)$ with probability $\rho \in (0, 1]$ independently. We shall assume that noise $\eta_n(\cdot), n \in [N]$ satisfy the following property.

Property 2.3. *For $n \in [N], t \in [T]$, $\eta_n(t)$ are independent sub-gaussian random variables, with $\mathbb{E}[\eta_n(t)] = 0$ and $\|\eta_n(t)\|_{\psi_2} \leq \gamma^8$.*

2.2 Comparison with Existing Time Series Models

We compare the model introduced in Section 2.1 with typical models utilized for time series analysis. As discussed earlier, the ‘temporal’ structure in time series models in the literature is captured through a combination of trend, periodicity, and stationarity. Precisely, trend is typically modeled through low-degree polynomials, periodicity through a finite sum of harmonics, and stationarity through (variants of) autoregressive or linear-time invariant setup.

The spectral representation of generic stationary processes, which includes autoregressive processes, implies that *any* sample-path of a stationary process can be decomposed into a weighted sum (precisely integral) of harmonics, where the ‘weights’ in the sum are sample path dependent—see Property 4.1, Chapter 4 of [20]. That is, a finite (weighted) sum of harmonics provides a good model representation for stationary processes with the model becoming more expressive as the number of harmonics grows.

Therefore, all three components of trends, periodicity and stationarity in a time series can be captured through a low-order polynomial and finite (weighted) sum of harmonics. Indeed, as argued in Proposition 2.1 below, such a model representation satisfies Property 2.2.

In Section 4, we extend this model when Properties 2.1 and 2.2 are only approximately satisfied. In particular, with any generic stationary process in mind, we quantify the approximation error based on the smoothness of the underlying time series and the number of harmonics used in the summation to approximate it.

Proposition 2.1 (Proposition 5.2, [1]). *Consider a time series $f : \mathbb{Z} \rightarrow \mathbb{R}$ with its element at time t denoted as*

$$f(t) = \sum_{a=1}^A \exp(\alpha_a t) \cdot \cos(2\pi\omega_a t + \phi_a) \cdot P_{m_a}(t), \quad (4)$$

where $\alpha_a, \omega_a, \phi_a \in \mathbb{R}$ are parameters, P_{m_a} is a degree $m_a \in \mathbb{N}$ polynomial in t . Then $f(\cdot)$ satisfies Property 2.2. In particular, consider the Hankel matrix of f over $[T]$, denoted as $H(f) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ with $H(f)_{ij} = f(i+j-1)$ for $i, j \in [\lfloor T/2 \rfloor]$. For any T , the rank of $H(f)$ is at most $G = A(m_{\max} + 1)(m_{\max} + 2)$, where $m_{\max} = \max_{a \in A} m_a$.

2.3 A Diagnostic Test for the Spatio-Temporal Model

In Sections 3 and 4, under the model described in Section 2, we establish the efficacy of mSSA. Beyond this model though, our work does not provide any guarantees for mSSA. Therefore, to utilize the guarantees of this work, it would be useful to have a data-driven diagnostic test that can help identify scenarios when the model of Section 2 may or may not hold. We discuss one such test in this section. Recall that the primary representation utilized by mSSA, as described in Section 1.1, is the stacked Page matrix (with parameter L). Observe that the Page matrix of a univariate time series for any $L \leq \lfloor T/2 \rfloor$ is simply the sub-matrix of the associated Hankel matrix: precisely, the Page matrix can be obtained by restricting to the top L rows and columns $1, L+1, \dots$ of the Hankel matrix. Therefore, the rank of the Hankel matrix is a bound on the rank of the Page matrix. Under the spatio-temporal factor model satisfying Properties 2.1 and 2.2, we establish the following low-rank property of the Page matrix of any particular time series as well as that of the stacked Page matrix.

Proposition 2.2. *Let Properties 2.1 and 2.2 hold. Then for any $L \leq \lfloor T/2 \rfloor$ with any $T \geq 1$, the rank of the Page matrix induced by the univariate time series $f_n(\cdot)$ for $n \in [N]$ is at most $R \times G$. Further, the rank of the stacked Page matrix induced by all N time series $f_1(\cdot), \dots, f_N(\cdot)$ is also at most $R \times G$.*

⁸For definition of $\|\cdot\|_{\psi_\alpha}$ -norm, see [25], for example.

The proof is in Appendix D where a more general result is established in Proposition 4.1. Proposition 2.2 suggests a “data driven diagnosis test” to verify whether mSSA is likely to succeed as per the results of this work. Specifically, if the (effective) rank⁹ of the Page matrix associated with any of the univariate components $f_n(\cdot)$ and the (effective) rank of stacked Page matrix associated with the multivariate time series with N component are *very different*, then mSSA may not be effective compared to SSA, but if they are *very similar* then mSSA is likely to be more effective compared to SSA. Our finite-sample results in Sections 3 and 4 indicate that the optimal value for L is $\sqrt{\min(N, T)T}$. Thus as a further test, if the effective rank of the stacked Page matrix does not scale much slower than L for $L \sim \sqrt{\min(N, T)T}$, then SSA (and mSSA) are unlikely to be effective methods.

Table 3 compares the (effective) rank of the stacked Page matrices for different benchmark time series data sets. The value of T equals 3993, 26304, and 10560 for the Financial, Electricity, and Traffic datasets respectively (see Appendix B for details on the datasets). We set $L = \lfloor \sqrt{\min(N, T)T} \rfloor$ for all datasets. When $N = 1$, this corresponds to L equals 63, 162, and 102 for the Financial, Electricity, and Traffic datasets respectively. Table 1 shows the effective rank in each dataset as we vary N . As can be seen, for $N = 1$, the effective rank is much smaller than L (or T) suggesting that SSA is likely to be effective. For Electricity and Financial datasets, the rank does not change by much as we increase N . However, relatively the rank does increase substantially for the Traffic dataset. This might explain why mSSA is relatively less effective for the Traffic dataset in contrast to the Financial and Electricity datasets as noted in Table 1.

Table 3: Effective rank of stacked Page matrix across benchmarks as we vary N .

Dataset	N = 1	N = 10	N = 100	N = 350
Electricity	19	37	44	31
Financial	1	3	3	6
Traffic	14	32	69	116

3 Main Results

We now provide bounds on the imputation and forecasting prediction error for mSSA under the spatio-temporal model introduced in Section 2. We start by defining the metric by which we measure prediction error. For imputation, we define prediction error as

$$\text{ImpErr}(N, T) = \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}[(f_n(t) - \hat{f}_n(t))^2]. \quad (5)$$

Here, the imputed estimate $\hat{f}_n(\cdot)$, $n \in [N]$ are produced by the imputation algorithm of Section 1.1. For forecasting, we define prediction error as

$$\text{ForErr}(N, T, L) = \frac{L}{NT} \sum_{n=1}^N \sum_{m'=1}^{T/L} \mathbb{E}[(f_n(L \times m') - \bar{f}_n(L \times m'))^2]. \quad (6)$$

Again, the forecasted estimate $\bar{f}_n(\cdot)$, $n \in [N]$ are produced by the forecasting algorithm of Section 1.1. In (5) and (6), the expectation is with respect to the randomness in observations due to noise and missing-ness. To state the main results, we shall utilize the following model restriction.

Property 3.1 (Balanced spectra). *Denote the $L \times (NT/L)$ stacked Page matrix associated with all N time series $f_1(\cdot), \dots, f_N(\cdot)$ as $\text{SP}(f) := \text{SP}((f_1, \dots, f_N), T, L)$. Under the setup of Proposition 2.2, $\text{rank}(\text{SP}(f)) = k \geq 1$ and $k \leq R \times G$. Then, for $L = \sqrt{\min(N, T)T}$, $\text{SP}(f)$ is such that $\sigma_k(\text{SP}(f)) \geq c\sqrt{NT}/\sqrt{k}$ for some absolute constant $c > 0$, where σ_k is the k -th largest singular value of $\text{SP}(f)$.*

Note that if $\sigma_k = \Theta(\sigma_1)$, then one can verify that Property 3.1 holds. Indeed, assuming that the non-zero singular values are ‘well-balanced’ is standard in the matrix/tensor estimation literature. Now,

⁹Effective rank a matrix is defined as the minimum number of singular values capturing $> 90\%$ of its spectral energy.

we state the main results. In what follows, we let $C(c, \Gamma_1, \Gamma_2, \gamma)$ denote a constant that depends only (polynomially) on model parameters $c, \Gamma_1, \Gamma_2, \gamma$. We also remind the reader that R, Γ_1, Γ_2 are defined in Property 2.1, G in 2.2, γ in Property 2.3 and c in Property 3.1.

Theorem 3.1 (Imputation). *Let Properties 2.1, 2.2, 2.3 and 3.1 hold. For a large enough absolute constant $C > 0$, let $\rho \geq C \frac{\log NT}{\sqrt{NT}}$. Then with $L = \sqrt{\min(N, T)T}$,*

$$\text{ImpErr}(N, T) \leq C(c, \Gamma_1, \Gamma_2, \gamma) \left(\frac{R^3 G \log NT}{\rho^4 \sqrt{\min(N, T)T}} \right)$$

Existence of linear model for forecasting. Recall from (3) that in mSSA, we learn a linear model between the last row of $\text{SP}((X_1, \dots, X_N), T, L)$ and the $L - 1$ rows above it (after de-noising the sub-matrix induced these $L - 1$ rows via HSVT). Hence, we first establish that in the idealized scenario (no noise, no missing values), there does indeed exist a linear model between the last row and the $L - 1$ rows above. To that end, define the shorthand $\text{SP}(f) = \text{SP}((f_1, \dots, f_N), T, L)$; let $\text{SP}(f)_L$ denote the L th row of $\text{SP}(f)$ and $\text{SP}'(f) \in \mathbb{R}^{(L-1) \times (NT/L)}$ denote the sub-matrix of $\text{SP}(f)$ formed by selecting top $L - 1$ rows. In the proposition below, we show exists a linear relationship between $\text{SP}(f)_L$ and $\text{SP}'(f)$.

Proposition 3.1. *Let Properties 2.1 and 2.2 hold. Then there exists $\beta^* \in \mathbb{R}^{L-1}$ such that $\text{SP}(f)_L^T = \text{SP}'(f)^T \beta^*$. Further, $\|\beta^*\|_0 \leq RG$.*

Theorem 3.2 (Forecasting). *Let the conditions of Theorem 3.1 hold. Then, with $L = \sqrt{\min(N, T)T}$ and with β^* defined in Proposition 3.1, we have*

$$\text{ForErr}(N, T, L) \leq C(c, \gamma, \Gamma_1, \Gamma_2) \max(1, \|\beta^*\|_1^2) \left(\frac{R^3 G \log NT}{\rho^4 \sqrt{\min(N, T)T}} \right).$$

We remark that Theorems 3.1, 3.2 and Proposition 3.1 are special cases of Theorems 4.1, 4.2 and Proposition 4.2 below, respectively; their proofs are in Appendices H, I, and I, respectively.

4 Approximately Low-Rank Hankel Representation

4.1 Extending Model and Main Result

We first introduce the definition of the approximate rank of a matrix.

Definition 4.1 (ϵ -approximate rank). *Given $\epsilon > 0$, a matrix $M \in \mathbb{R}^{a \times b}$ is said to have ϵ -approximate rank at most $s \geq 1$ if there exists a rank s matrix $M_s \in \mathbb{R}^{a \times b}$ such that $\|M - M_s\|_\infty < \epsilon$.*

Definition 4.2 ((G, ϵ) -Hankel Time Series). *For a given $\epsilon \geq 0$ and $G \geq 1$, a time series $f : \mathbb{Z} \rightarrow \mathbb{R}$ is called a (G, ϵ) -Hankel time series if for any $T \geq 1$, its Hankel matrix has ϵ -approximate rank G .*

We extend the model of Section 2 by replacing Property 2.2 by the following.

Property 4.1. *For each $r \in [R]$ and for any $T \geq 1$, the Hankel Matrix $H(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ associated with time series W_{rt} , $t \in [T]$ has ϵ -approximate rank at most G for $\epsilon > 0$. That is, for each $r \in [R]$, $W_{r\cdot}$ is a (G, ϵ) -Hankel time series.*

We state an implication of the above stated properties.

Proposition 4.1. *Let Properties 2.1 and 4.1 hold. For any $L \leq \lfloor T/2 \rfloor$ with any $T \geq 1$, the stacked Page matrix induced by the N time series $f_1(\cdot), \dots, f_N(\cdot)$ has ϵ' -rank at most $R \times G$ for $\epsilon' = R\Gamma_1\epsilon$.*

Below, we provide generalizations of the results stated in Section 3; to do so, we utilize Property 4.2 which is analogous to Property 3.1 but for the approximate low-rank setting.

Property 4.2 (Approximately balanced spectra). *Under the setup of Proposition 4.1, we can represent the $L \times (NT/L)$ stacked Page matrix associated with all N time series $f_1(\cdot), \dots, f_N(\cdot)$ as $\text{SP}(f) = \tilde{M} + E$ with $\text{rank}(\tilde{M}) = k \geq 1$ and $k \leq R \times G$ and $\|E\|_\infty \leq R\Gamma_1\epsilon$. Then, for $L = \sqrt{\min(N, T)T}$, \tilde{M} is such that $\sigma_k(\tilde{M}) \geq c\sqrt{NT}/\sqrt{k}$ for some absolute constant $c > 0$, where σ_k is the k -th largest singular value of \tilde{M} .*

Theorem 4.1 (Imputation). *Let Properties 2.1, 4.1, 2.3 and 4.2 hold. For a large enough absolute constant $C > 0$, let $\rho \geq C \frac{\log NT}{\sqrt{NT}}$. Then, with $L = \sqrt{\min(N, T)T}$,*

$$\text{ImpErr}(N, T) \leq C(c, \Gamma_1, \Gamma_2, \gamma) \left(\frac{R^3 G \log NT}{\rho^4 \sqrt{\min(N, T)T}} + \frac{R^4 G(\epsilon + \epsilon^3)}{\rho^2} \right)$$

where $C(c, \Gamma_1, \Gamma_2, \gamma)$ is a positive constant dependent on model parameters including $\Gamma_1, \Gamma_2, \gamma$.

We remind the reader that R, Γ_1, Γ_2 are defined in Property 2.1, G in 2.2, γ in Property 2.3 and c in Property 4.2.

Existence of Linear Model. We now state Proposition 4.2, which is analogous to Proposition 3.1, but for the approximate low-rank setting.

Proposition 4.2. *Let Properties 2.1 and 4.1 hold. Then, there exists $\beta^* \in \mathbb{R}^{L-1}$, such that $\|\text{SP}(f)_L^T - \text{SP}'(f)^T \beta^*\|_\infty \leq R\Gamma_1(1 + \|\beta^*\|_1)\epsilon$, Further $\|\beta^*\|_0 \leq RG$.*

Theorem 4.2 (Forecasting). *Let the conditions of Theorem 4.1 hold. Then, with $L = \sqrt{\min(N, T)T}$ and with β^* defined in Proposition 4.2, we have*

$$\text{ForErr}(N, T, L) \leq C(c, \gamma, \Gamma_1, \Gamma_2) \max(1, \|\beta^*\|_1^2) \left(\frac{R^3 G \log NT}{\rho^4 \sqrt{\min(N, T)T}} + \frac{R^4 G(\epsilon + \epsilon^3)}{\rho^2} \right).$$

4.2 Hankel Calculus

We present a key property of the model class satisfying Property 4.1, i.e. time series that have an approximate low-rank Hankel matrix representation. To that end, we define ‘addition’ and ‘multiplication’ for time series. Given two time series $f_1, f_2 : \mathbb{Z} \rightarrow \mathbb{R}$, define their addition, denoted $f_1 + f_2 : \mathbb{Z} \rightarrow \mathbb{R}$ as $(f_1 + f_2)(t) = f_1(t) + f_2(t)$, for all $t \in \mathbb{Z}$. Similarly, their multiplication, denoted $f_1 \circ f_2 : \mathbb{Z} \rightarrow \mathbb{R}$ as $(f_1 \circ f_2)(t) = f_1(t) \times f_2(t)$, for all $t \in \mathbb{Z}$. Now, we state a key property for the model class satisfying Property 4.1 (proof in Appendix E).

Proposition 4.3. *For $i \in \{1, 2\}$, let f_i be a (G_i, ϵ_i) -Hankel time series for $G_i \geq 1, \epsilon_i \geq 0$. Then, $f_1 + f_2$ is a $(G_1 + G_2, \epsilon_1 + \epsilon_2)$ -Hankel time series and $f_1 \circ f_2$ is a $(G_1 G_2, 3 \max(\epsilon_1, \epsilon_2) \cdot \max(\|f_1\|_\infty, \|f_2\|_\infty))$ -Hankel time series.*

4.3 Examples of (G, ϵ) -Hankel Time Series

We establish that many important classes of time series dynamics studied in the literature are instances of (G, ϵ) -Hankel time series, i.e. they satisfy Property 4.1. In particular: any differentiable periodic function (Proposition 4.5); any time series with a Hölder continuous latent variable representation (Proposition 4.6). Proofs of Propositions 4.3, 4.4, 4.5 and 4.6 can be found in Appendix E.

Example 1. (G, ϵ) -LRF time series. We start by defining a linear recurrent formula (LRF), which is a standard model for linear time-invariant systems.

Definition 4.3 ((G, ϵ) -LRF). *For $G \in \mathbb{N}$ and $\epsilon \geq 0$, a time series f is said to be a (G, ϵ) -Linear Recurrent Formula (LRF) if for all $T \in \mathbb{Z}$ and $t \in [T]$, there exists $g : \mathbb{Z} \rightarrow \mathbb{R}$ such that*

$$f(t) = g(t) + h(t),$$

where for all $t \in \mathbb{Z}$, (i) $g(t) = \sum_{l=1}^G \alpha_l g(t-l)$ with constants $\alpha_1, \dots, \alpha_G$, and (ii) $|h(t)| \leq \epsilon$.

Now we establish a time series f that is a (G, ϵ) -LRF is also (G, ϵ) -Hankel.

Proposition 4.4. *If f is (G, ϵ) -LRF representable, then it is (G, ϵ) -Hankel representable.*

LRF’s cover a broad class of time series functions, including any finite sum of products of harmonics, polynomials and exponentials. In particular, it can be easily verified that a time series described by (4) is a $(G, 0)$ -LRF, where $G \leq A(m_{\max} + 1)(m_{\max} + 2)$ with $m_{\max} = \max_{a \in A} m_a$.

Example 2. “smooth” and periodic time series. We establish that any differentiable periodic function is (G, ϵ) -LRF and hence (G, ϵ) -Hankel for appropriate choices of G and ϵ .

Definition 4.4 ($C^k(R, \text{PER})$). For $k \geq 1$ and $R > 0$, we use $C^k(R, \text{PER})$ to denote the class of all time series $f : \mathbb{R} \rightarrow \mathbb{R}$ such that it is R periodic, i.e. $f(t + R) = f(t)$ for all $t \in \mathbb{R}$ and the k -th derivative of f , denoted $f^{(k)}$, exists and is continuous.

Proposition 4.5. Any $f \in C^k(R, \text{PER})$ is

$$\left(4G, C(k, R) \frac{\|f^{(k)}\|}{G^{k-0.5}}\right) - \text{Hankel representable},$$

for any $G \geq 1$. Here $C(k, R)$ is a term that depends only on k, R and $\|f^{(k)}\|^2 = \frac{1}{R} \int_0^R (f^{(k)}(t))^2 dt$.

Example 3. time series with latent variable model (LVM) structure. We now show that if a time series has a LVM representation, and the latent function is Hölder continuous, then it has a (G, ϵ) -Hankel representation for appropriate choice of $G \geq 1$ and $\epsilon \geq 0$. We first define the Hölder class of functions; this class of functions is widely adopted in the non-parametric regression literature [26]. Given a function $g : [0, 1)^K \rightarrow \mathbb{R}$, and a multi-index $\kappa \in \mathbb{N}^K$, let the partial derivate of g at $x \in [0, 1)^K$, if it exists, be denoted as $\nabla_\kappa g(x) = \frac{\partial^{|\kappa|} g(x)}{(\partial x)^\kappa}$ where $|\kappa| = \sum_{j=1}^K \kappa_j$ and $(\partial x)^\kappa = \partial^{\kappa_1} x_1 \cdots \partial^{\kappa_K} x_K$.

Definition 4.5 ((α, \mathcal{L})-Hölder Class). Given $\alpha, \mathcal{L} > 0$, the Hölder class $\mathcal{H}(\alpha, \mathcal{L})$ on $[0, 1)^K$ ¹⁰ is defined as the set of functions $g : [0, 1)^K \rightarrow \mathbb{R}$ whose partial derivatives satisfy for all $x, x' \in [0, 1)^K$, $\sum_{\kappa: |\kappa| = \lfloor \alpha \rfloor} \frac{1}{\kappa!} |\nabla_\kappa g(x) - \nabla_\kappa g(x')| \leq \mathcal{L} \|x - x'\|_\infty^{\alpha - \lfloor \alpha \rfloor}$. Here $\lfloor \alpha \rfloor$ refers to the greatest integer strictly smaller than α and $\kappa! = \prod_{j=1}^K \kappa_j!$.

Note that if $\alpha \in (0, 1]$, then the definition above is equivalent to the (α, \mathcal{L}) -Lipschitz condition, i.e., $|g(x) - g(x')| \leq \mathcal{L} \|x - x'\|_\infty^\alpha$, for $x, x' \in [0, 1)^K$. Given a time series $f : \mathbb{Z} \rightarrow \mathbb{R}$, for any $T \geq 1$, recall the Hankel matrix $H \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ is defined such that its entry in row $i \in [\lfloor T/2 \rfloor]$ and column $j \in [\lfloor T/2 \rfloor]$ is given by $H_{ij} = f(i + j - 1)$. We call a time series f to have (α, \mathcal{L}) -Hölder smooth LVM representation for $\alpha, \mathcal{L} > 0$ if for any given $T \geq 1$, the corresponding Hankel matrix H satisfies: for $i, j \in [\lfloor T/2 \rfloor]$, $H_{ij} = g(\theta_i, \omega_j)$, where $\theta_i, \omega_j \in [0, 1)^K$ are latent parameters and $g(\cdot, \omega) \in \mathcal{H}(\alpha, \mathcal{L})$ for any $\omega \in [0, 1)^K$. It can be verified that a $(G, 0)$ -Hankel time series is an instance of such a LVM representation with corresponding $g(x, y) = x^T y$. Thus in a sense, this model is a natural generalization of the $(G, 0)$ -Hankel matrix representation. The following proposition connects this LVM representation to the (G, ϵ) -Hankel representation for appropriately defined $G \geq 1, \epsilon > 0$.

Proposition 4.6. Given $\alpha, \mathcal{L} > 0$, let f have (α, \mathcal{L}) -Hölder smooth LVM representation. Then for all $\epsilon > 0$, f is

$$(C(\alpha, K) \left(\frac{1}{\epsilon}\right)^K, \mathcal{L} \epsilon^\alpha) - \text{Hankel representable}.$$

Here $C(\alpha, K)$ is a term that depends only on α and K .

5 Experiments

We describe experiments supporting our theoretical results for mSSA. In particular, we provide details of the experiments run to create the summary results described earlier in Table 1. In Appendix B, we describe the datasets utilized and the various algorithms we compare with as well as the procedure for selecting the hyper-parameters in each algorithm. In Section 5.1 and 5.2, we report the imputation and forecasting results. Note that in all experiments, we use the Normalized Root Mean Squared Error (NRMSE) as our accuracy metric. That is, we normalize all the underlying time series to have zero mean and unit variance before calculating the root mean squared error. We use this metric as it weighs the error on each time series equally.

5.1 Imputation

Setup. We test the robustness of the imputation performance by adding two sources of corruption to the data - varying the percentage of observed values and varying the amount of noise we perturb the

¹⁰The domain is easily extended to any compact subset of \mathbb{R}^K .

observations by. We test imputation performance by how accurately we recover missing values. We compare the performance of mSSA with TRMF, a method which achieves state-of-the-art imputation performance. Further, to analyze the added benefit of exploiting the spatial structure in a multivariate time series using mSSA, we compare with the SSA variant introduced in [1].

Results. Figures 3a, 3c, 3e, 4a, and 4c show the imputation error in the aforementioned datasets as we vary the fraction of missing values, while Figures 3b, 3d, 3f, 4b, and 4d show the imputation error as we vary σ , the standard deviation of the gaussian noise. We see that as we vary the fraction of missing values and noise levels, mSSA outperforms both TRMF and SSA in $\sim 75\%$ of experiments run. It is noteworthy the large empirical gain in mSSA over SSA, giving credence to the spatio-temporal model we introduce. The average NRMSE across all experiments for each dataset is reported in Table 1, where mSSA outperforms every other method across all datasets except for the Traffic dataset.

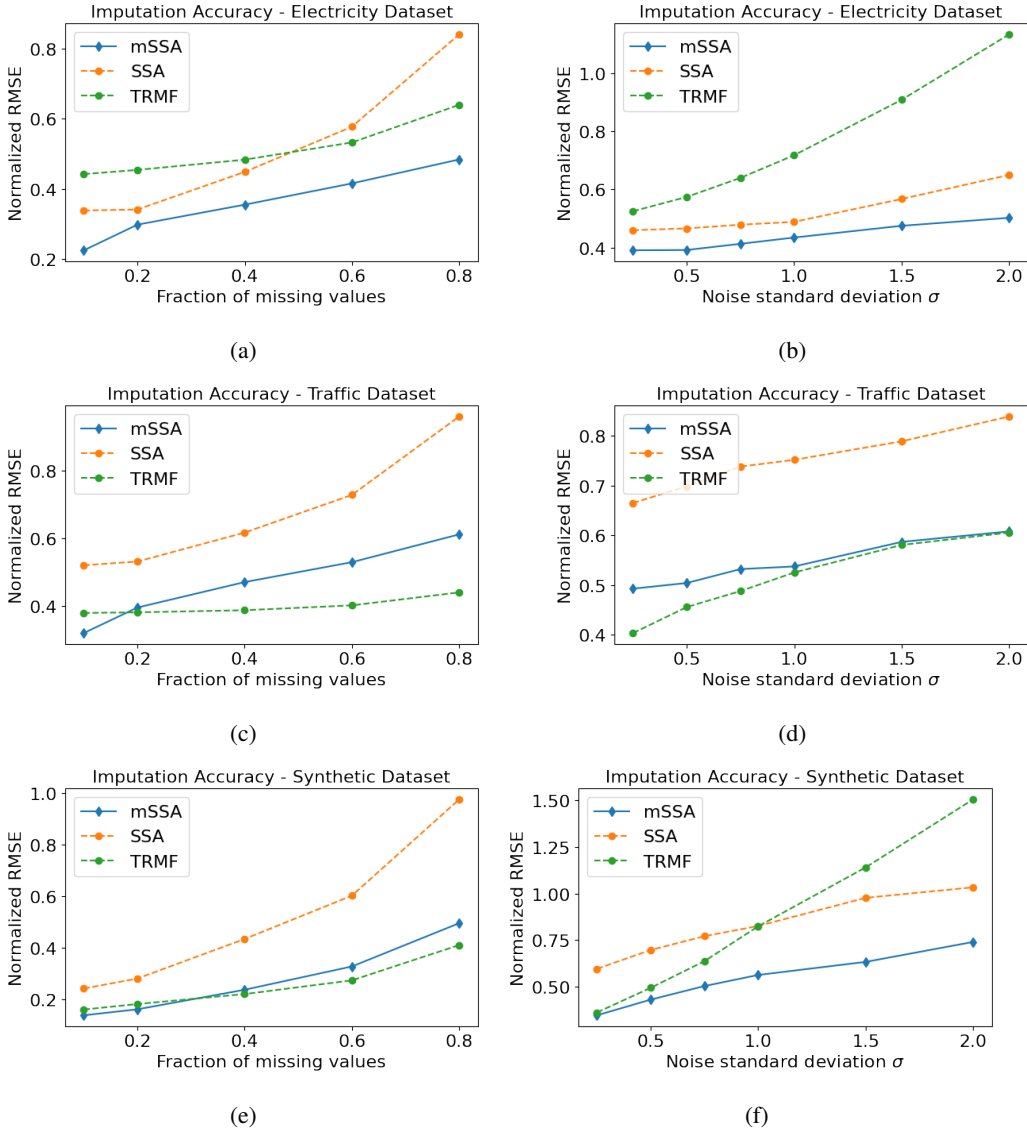


Figure 3: mSSA vs. TRMF vs. SSA - imputation performance on the Electricity, Traffic and Synthetic datasets. Figures 3a, 3c, and 3e, show imputation accuracy of mSSA, TRMF and SSA as we vary the fraction of missing values; Figures 3b, 3d, and 3f show imputation accuracy as we vary the noise level (and with 50% of values missing).

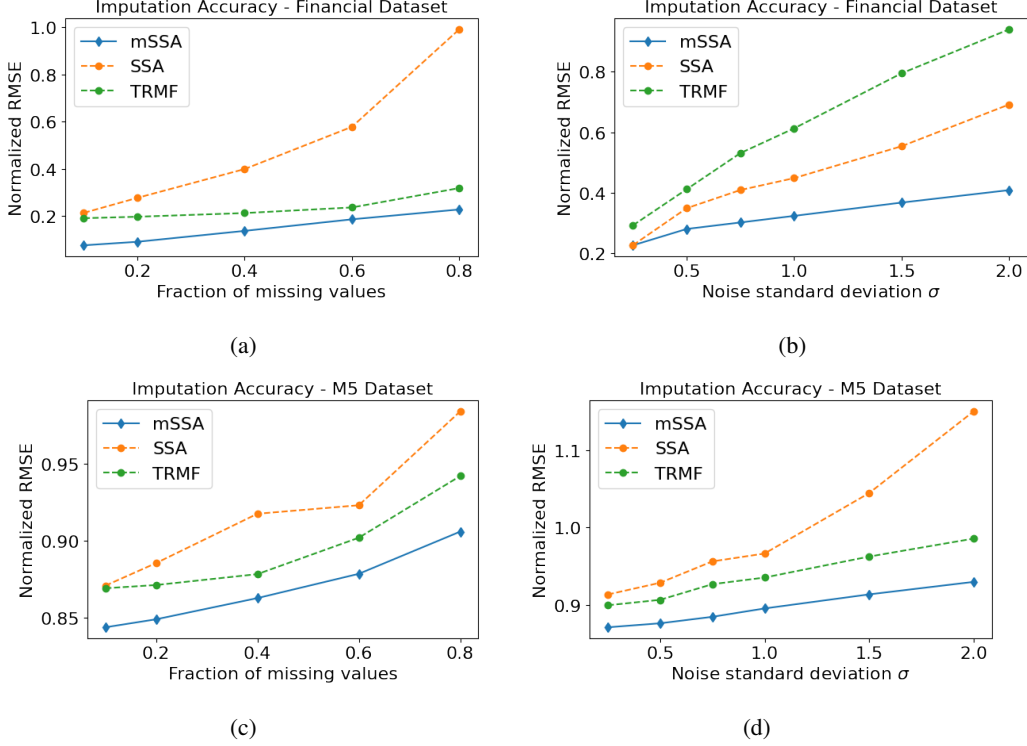


Figure 4: mSSA vs. TRMF vs. SSA - imputation performance on the Financial and M5 datasets. Figures 4a, and 4c show imputation accuracy of mSSA, TRMF and SSA as we vary the fraction of missing values; Figures 4b, and 4d show imputation accuracy as we vary the noise level (and with 50% of values missing).

5.2 Forecasting

Setup. We test the forecasting accuracy of the proposed mSSA against several state-of-the-art algorithms. For each dataset, we split the data into training, validation, and testing datasets as outlined in Appendix B.1. As was done in the imputation experiments, we vary how much each dataset is corrupted by varying the percentage of observed values and the noise levels.

Results. Figures 5a, 5c, 5e, 6a, and 6c show the forecasting accuracy of mSSA and other methods in the aforementioned datasets as we vary the fraction of missing values, while Figures 5b, 5d, 5f, 6b, and 6d show the forecasting accuracy as we vary the standard deviation of the added gaussian noise. We see that as we vary the fraction of missing values and noise level, mSSA is the best or comparable to the best performing method in $\sim 80\%$ of experiments. In terms of the average NRMSE across all experiments, we find that mSSA performs similar to or better than every other method across all datasets except for the traffic dataset as was reported in Table 1.

6 Algorithmic Extensions of mSSA

6.1 Variance Estimation

We extend the mSSA algorithm to estimate the time-varying variance of a time series by making the following simple observation. If we apply mSSA to the squared observations, $X_n^2(t)$, we will recover an estimate of $\mathbb{E}[X_n^2(t)]$ (for $\rho = 1$). However, observe that $\text{Var}[X_n(t)] = \mathbb{E}[X_n^2(t)] - \mathbb{E}[X_n(t)]^2$. Therefore, by applying mSSA twice, once on $X_n(t)$ and once on $X_n^2(t)$ for $n \in [N]$ and $t \in [T]$, and subsequently taking the component-wise difference of the two estimates will lead to an estimate of the variance. This suggests a simple algorithm which we describe next. We note this observation suggests *any* mean estimation algorithm (or imputation) in time series analysis can be converted to estimate the time varying variance – this ought to be of interest in its own right.

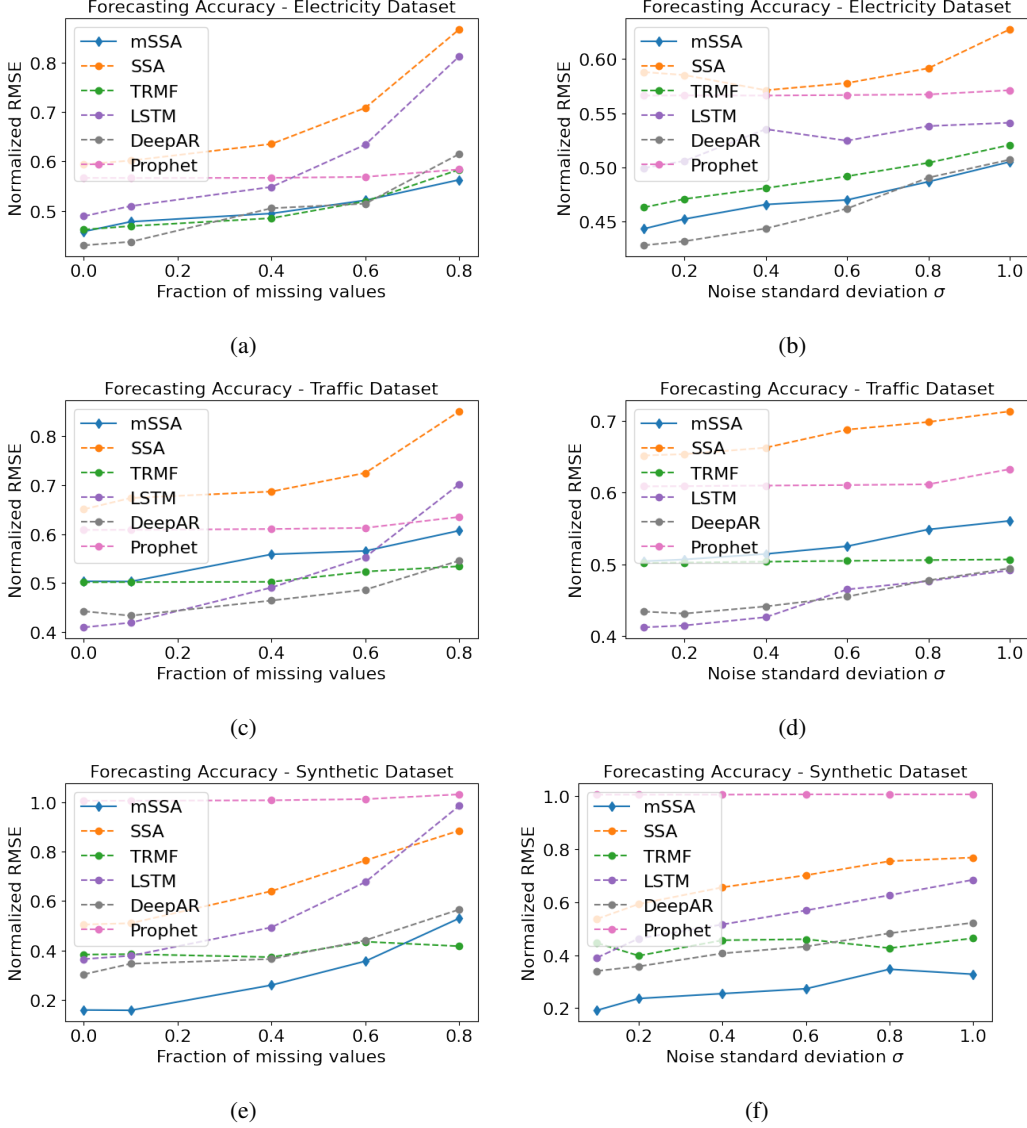


Figure 5: mSSA forecasting performance on standard multivariate time series benchmark is competitive with/outperforming industry standard methods as we vary the number of missing data and noise level. Figures 5a, 5c, and 5e show the forecasting accuracy of all methods on the Electricity, Traffic and Synthetic datasets with varying fraction of missing values; Figures 5b, 5d, and 5f show the forecasting accuracy on the same datasets with varying noise level.

Algorithm. As described in Section 1.1, let $L \geq 1$ and $k \geq 1$ be algorithm parameters. First, apply mSSA on observations $X_n(t)$, $n \in [N]$, $t \in [T]$ to produce imputed estimates $\hat{f}_n(t)$. Next, apply mSSA on observations $X_n^2(t)$, $n \in [N]$, $t \in [T]$ to produce imputed estimates $\hat{g}_n(t)$. Lastly, we denote $\hat{\sigma}_n^2(t) = \max(0, \hat{g}_n(t) - \hat{f}_n(t)^2)$, $n \in [N]$, $t \in [T]$ as our estimate of the time-varying variance.

Model. For $n \in [N]$, $t \in [T]$, let $\sigma_n^2(t) = \mathbb{E}[\eta_n^2(t)]$ be the time-varying variance of the time series observations, i.e. if $\rho = 1$ then $\sigma_n^2(t) = \text{Var}[X_n(t)] = \mathbb{E}[X_n^2(t)] - f_n^2(t)$. Let $\Sigma \in \mathbb{R}^{N \times T}$ be the matrix induced by the latent time-varying variances of the N time series of interest, i.e. the entry in row n at time t in Σ is $\Sigma_{nt} = \sigma_n^2(t)$. To capture the “spatial” and “temporal” structure across the N latent time-varying variances, we assume the latent variance matrix Σ satisfies Properties 6.1 and

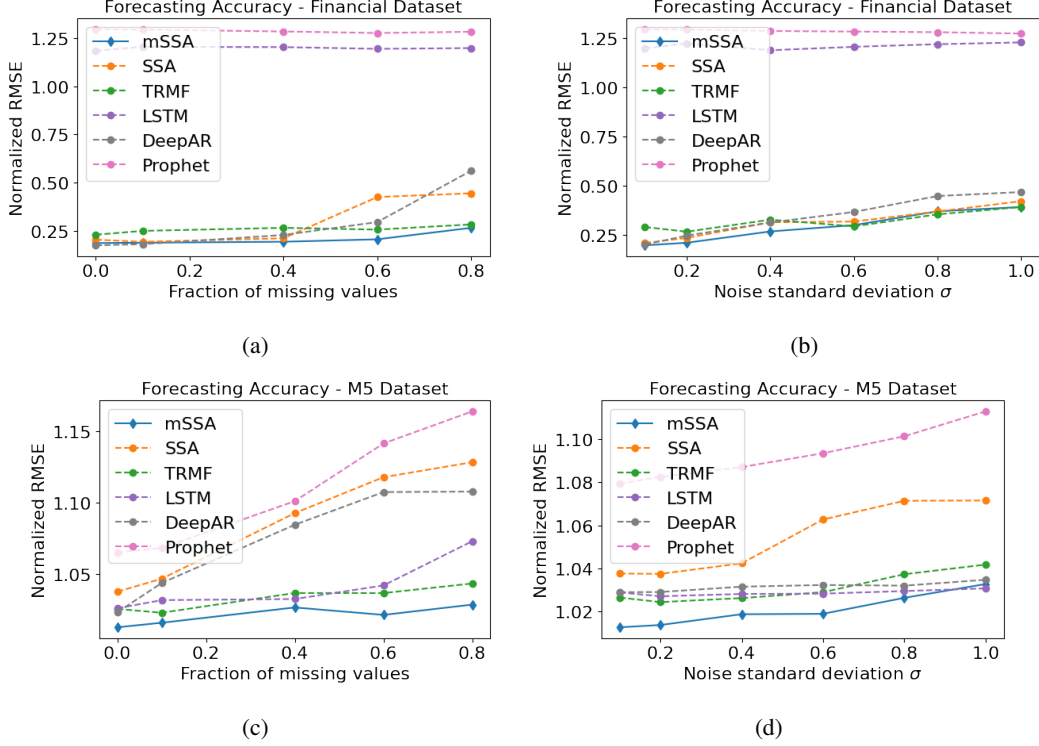


Figure 6: Figures 6a, and 6c show the forecasting accuracy of all methods on the financial and M5 datasets with varying fraction of missing values; Figures 6b, and 6d show the forecasting accuracy on the same datasets with varying noise levels.

6.2. These properties are analogous to those assumed about the latent mean matrix \mathbf{M} (defined in Section 2); in particular, Properties 2.1 and 2.2. We state them next.

Property 6.1. Let $R' = \text{rank}(\Sigma)$, i.e, for any $n \in [N], t \in [T]$, $\Sigma_{nt} = \sum_{r=1}^{R'} U'_{nr} W'_{rt}$, where the factorization is such that $|U'_{nr}| \leq \Gamma'_1$, $|W'_{rt}| \leq \Gamma'_2$ for $\Gamma'_1, \Gamma'_2 > 0$.

Like Property 2.1, the above property captures the “spatial” structure within N time series of variances. To capture the “temporal” structure, next we introduce an analogue of Property 2.2. To that end, for each $r \in [R']$, define the $\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor$ Hankel matrix of each time series W'_{rt} , $t \in [T]$ as $H'(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$, where $H'(r)_{ij} = W'_{r(i+j-1)}$ for $i, j \in [\lfloor T/2 \rfloor]$.

Property 6.2. For each $r \in [R']$, the Hankel Matrix $H'(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ associated with time series $W'_{rt}, t \in [T]$ has rank at most G' .

Result. To establish the estimation error for the variance estimation algorithm under the spatio-temporal model above, we need the following additional property (analogous to Property 3.1).

Property 6.3 (Balanced spectra). Denote the $L \times (NT/L)$ stacked Page matrix associated with all N time series $\sigma_1^2(\cdot), \dots, \sigma_N^2(\cdot)$ as $\text{SP}(\sigma^2) := \text{SP}(\sigma_1^2, \dots, \sigma_N^2, T, L)$. Due to Properties 6.1 and 6.2, and a simple variant of Proposition 2.2, we have $\text{rank}(\text{SP}(\sigma^2)) = k' \geq 1$ and $k' \leq R' \times G'$. Then, for $L = \sqrt{\min(N, T)T}$, $\text{SP}(\sigma^2)$ is such that $\sigma'_k(\mathbf{M}) \geq c\sqrt{NT}/\sqrt{k'}$ for some absolute constant $c > 0$, where σ'_k is the k -th singular value, order by magnitude, of $\text{SP}(\sigma^2)$.

Theorem 6.1 (Variance Estimation). Let Properties 2.1, 2.2, 2.3, 3.1, 6.1, 6.2, and 6.3 hold. Additionally let $|\hat{f}_n(t)| \leq \Gamma_3$ for all $n \in [N], t \in [T]$. Lastly, let $L = \sqrt{\min(N, T)T}$ and $\rho = 1$. Then the variance prediction error is bounded above as

$$\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}[(\sigma_n(t)^2 - \hat{\sigma}_n^2(t))^2] \leq \tilde{C} \left(\frac{(G^2 + G') \log^2 NT}{\sqrt{\min(N, T)T}} \right).$$

where \tilde{C} is a constant dependent (polynomially) on model parameters $\Gamma_1, \Gamma_2, \Gamma_3, \Gamma'_1, \Gamma'_2, \gamma, R, R'$.

Proof of Theorem 6.1 can be found in Appendix J.

6.2 Tensor SSA

Page tensor. We introduce an order-three tensor representation of a multivariate time series which we term the ‘Page tensor’. Given N time series, with observations over T time steps and hyperparameter $L \geq 1$, define $\mathbf{T} \in \mathbb{R}^{N \times T/L \times L}$ such that

$$\mathbf{T}_{n\ell s} = f_n((s-1) \times L + \ell), \quad n \in [N], \ell \in [L], s \in [T/L].$$

The corresponding observation tensor, $\mathbb{T} \in (\mathbb{R} \cup \{\star\})^{N \times T/L \times L}$, is

$$\mathbb{T}_{n\ell s} = X_n((s-1) \times L + \ell), \quad n \in [N], \ell \in [L], s \in [T/L]. \quad (7)$$

See Figure 7 for a visual depiction of \mathbb{T} . Under the model described in Section 2, we have the following properties.

Proposition 6.1. *Let Properties 2.1, 2.2, and 2.3 hold. Then, for any $1 \leq L \leq \sqrt{T}$, \mathbf{T} has canonical polyadic (CP)-rank¹¹ at most $R \times G$. Further, all entries of \mathbb{T} are independent random variables with each entry observed with probability $\rho \in (0, 1]$, and $\mathbb{E}[\mathbb{T}] = \rho \mathbf{T}$.*

tSSA: time series imputation using the Page tensor representation. The Page tensor representation and Proposition 6.1 collectively suggest that time series imputation can be reduced to low-rank tensor estimation, i.e., recovering a tensor of low CP-rank from its noisy, partial observations. Over the past decade, the field of low-rank tensor (and matrix) estimation has received great empirical and theoretical interest, leading to a large variety of algorithms including spectral, convex optimization, and nearest neighbor based approaches. We list a few works which have explicit finite-sample rates for noisy low-rank tensor completion [4, 31, 6, 33, 23]). As a result, we “blackbox” the tensor estimation algorithm used in tSSA as a pivotal subroutine. Doing so allows one the flexibility to use the tensor estimation algorithm of their choosing within tSSA. Consequently, as the tensor estimation literature continues to advance, the “meta-algorithm” of tSSA will continue to improve in parallel. To that end, we give a definition of a tensor estimation algorithm for a generic order- d tensor. Note that when $d = 2$, this reduces to standard matrix estimation (ME).

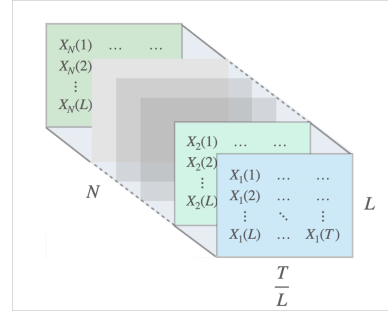


Figure 7: The observations Page tensor.

Definition 6.1 (Matrix/Tensor Estimation). *For $d \geq 2$, denote $\text{TE}_d : \{\star, \mathbb{R}\}^{n_1 \times n_2 \times \dots \times n_d} \rightarrow \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ as an order- d tensor estimation algorithm. It takes as input an order- d tensor \mathbb{G} with noisy, missing entries, where $\mathbb{E}[\mathbb{G}] = \rho \mathbf{G}$ and $\rho \in (0, 1]$ is the probability of each entry in \mathbb{G} being observed. TE_d then outputs an estimate of \mathbf{G} denoted as $\hat{\mathbf{G}} = \text{TE}_d(\mathbb{G})$.*

We assume the following ‘oracle’ error convergence rate for TE_d ; for ease of exposition, we restrict our attention to the setting where $\rho = 1$.

Property 6.4. *For $d \geq 2$, assume TE_d satisfies the following: the estimate $\hat{\mathbf{G}} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$, which is the output of $\text{TE}_d(\mathbb{G})$ with $\mathbb{E}[\mathbb{G}] = \mathbf{G}$, satisfies*

$$\frac{1}{n_1 \dots n_d} \|\hat{\mathbf{G}} - \mathbf{G}\|_F^2 = \tilde{\Theta} \left(1 / \min(n_1, \dots, n_d)^{\lceil d/2 \rceil} \right).$$

Here, $\tilde{\Theta}(\cdot)^{12}$ suppresses dependence on noise, i.e., $\mathbf{E} = \mathbf{G} - \mathbb{E}[\mathbb{G}]$, $\log(\cdot)$ factors, and CP-rank of \mathbf{G} .

¹¹The CP-rank of an order- d tensor $\mathbf{T} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ is the smallest value of $r \in \mathbb{N}$ such that $\mathbf{T}_{i_1, \dots, i_d} = \sum_{k=1}^r u_{i_1, k} \dots u_{i_d, k}$, where $u_{i_\ell, k}$ are latent factors for $\ell \in [d]$.

¹²We are using the standard Big- $\tilde{\Theta}$ notation where the dependence on poly logarithmic terms are suppressed for brevity of presentation.

Property 6.4 holds for a variety of matrix/tensor estimation algorithms. For $d = 2$, it holds for HSVT as we establish in the proof of Theorem 3.1 for mSSA of $\tilde{O}(1/\sqrt{\min(N, T), T})$. It is straightforward to show that this is the best rate achievable for TE_2 . For $d > 3$, it has recently been shown that Property 6.4 provably holds for a spectral gradient descent based algorithm [6] (see Corollary 1.5 of [6]), conditioned on certain standard “incoherence” conditions imposed on the latent factors of \mathbf{G} ; another spectral algorithm that achieved the same rate was furnished in [31], which the authors also establish is minimax optimal.

tSSA algorithm. We now define the “meta” tSSA algorithm; the two algorithmic hyper-parameters are $L \geq 1$ (defined in (7)) and TE_3 (the order-three tensor estimation algorithm one chooses). First, using $X_n(t)$ for $n \in [N], t \in [T]$, construct Page tensor \mathbb{T} as in (7). Second, obtain $\hat{\mathbb{T}}$ as the output of $\text{TE}_3(\mathbb{T})$ and read off $\hat{f}_n(t)$ by selecting appropriate entry in $\hat{\mathbb{T}}$.

Algorithmic comparison: tSSA vs. mSSA vs. ME. We now provide a unified view of tSSA, mSSA, and “vanilla” ME (which we describe below) to do time series imputation. All three methods have two key steps: (i) data transformation – converting the observations $X_n(t)$ into a particular data representation/structure; (ii) de-noising– applying some form of matrix/tensor estimation to de-noise the constructed data representation.

- tSSA – using $X_n(t)$, create the Page tensor $\mathbb{T} \in \mathbb{R}^{N \times L \times T/L}$ as in (7); apply $\text{TE}_3(\mathbb{T})$ to get $\hat{\mathbb{T}}$ (e.g. using the method in [6]); read off $\hat{f}_n(t)$ by selecting appropriate entry in $\hat{\mathbb{T}}$.
- mSSA – using $X_n(t)$, create the stacked Page matrix $\text{SP}((X_1, \dots, X_N), T, L) \in \mathbb{R}^{L \times (N \times T/L)}$ as detailed in Section 1.1; apply $\text{TE}_2(\text{SP}((X_1, \dots, X_N), T, L))$ to get $\widehat{\text{SP}}((X_1, \dots, X_N), T, L)$ (where we use HSVT for $\text{TE}_2(\cdot)$); read off $\hat{f}_n(t)$ by selecting appropriate entry in $\widehat{\text{SP}}((X_1, \dots, X_N), T, L)$.
- ME – using $X_n(t)$, create $\mathbf{X} \in \mathbb{R}^{N \times T}$, where \mathbf{X}_{nt} is equal to $X_n(t)$; apply $\text{TE}_2(\mathbf{X})$ (e.g. using HSVT as in mSSA) to get $\widehat{\mathbf{X}}$; read off $\hat{f}_n(t)$ by selecting appropriate entry in $\widehat{\mathbf{X}}$.

This perspective also suggests that one can use any “blackbox” matrix estimation routine to de-noise the constructed stacked Page matrix in mSSA; HSVT is one such choice that we analyze.

Theoretical comparison: tSSA vs. mSSA vs. ME. We now do a theoretical comparison of the relative effectiveness of tSSA, mSSA, and ME in imputing a multivariate time series $X_n(t)$ for $n \in [N], t \in [T]$, as we vary N and T . To that end, let $\text{ImpErr}(N, T; \text{tSSA})$, $\text{ImpErr}(N, T; \text{mSSA})$, and $\text{ImpErr}(N, T; \text{ME})$ denote the imputation error for tSSA, mSSA, and ME, respectively.

Proposition 6.2. *For tSSA and mSSA, pick hyper-parameter $L = \sqrt{T}$, $L = \sqrt{\min(N, T)T}$, respectively. Let Property 6.4 hold. Then if:*

- (i) $T = o(N)$: $\text{ImpErr}(N, T; \text{tSSA}), \text{ImpErr}(N, T; \text{mSSA}) = \tilde{O}(\text{ImpErr}(N, T; \text{ME}))$;
- (ii) $T^{1/3} = o(N)$, $N = o(T)$: $\text{ImpErr}(N, T; \text{tSSA}) = \tilde{o}(\text{ImpErr}(N, T; \text{mSSA})), \text{ImpErr}(N, T; \text{mSSA}) = \tilde{o}(\text{ImpErr}(N, T; \text{ME}))$;
- (iii) $N = o(T^{1/3})$: $\text{ImpErr}(N, T; \text{mSSA}) = \tilde{o}(\text{ImpErr}(N, T; \text{tSSA})), \text{ImpErr}(N, T; \text{tSSA}) = \tilde{o}(\text{ImpErr}(N, T; \text{ME}))$,

where $\tilde{o}(\cdot)$, $\tilde{O}(\cdot)$ suppresses dependence on noise parameters, CP-rank, poly-logarithmic factors.

We note given Property 6.4, $L = \sqrt{T}$ is optimal for tSSA and $L = \sqrt{\min(N, T)T}$ is optimal for mSSA. See Figure 2 in Section 1 for a graphical depiction of the different regimes in Proposition 6.2. Proofs of Proposition 6.1 and 6.2 below can be found in Appendix K.

Application to Time-varying Recommendation Systems In Appendix C, we discuss the extension of our spatio-temporal model and tSSA to time-varying recommendation systems.

7 Conclusion

We provide theoretical justification of a practical, simple variant of mSSA, a method heavily used in practice but with limited theoretical understanding. We show how to extend mSSA to estimate

time-varying variance and introduce a tensor variant, tSSA, which builds upon recent advancements in tensor estimation. We hope this work motivates future inquiry into the connections between the classical field of time series analysis and the modern, growing field of matrix/tensor estimation.

References

- [1] A. Agarwal, M. J. Amjad, D. Shah, and D. Shen. Model agnostic time series analysis via matrix estimation. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 2(3):40, 2018.
- [2] A. Agarwal, D. Shah, D. Shen, and D. Song. On robustness of principal component regression. In *Advances in Neural Information Processing Systems*, pages 9889–9900, 2019.
- [3] A. Agarwal, D. Shah, D. Shen, and D. Song. On robustness of principal component regression. *Accepted to appear in Journal of the American Statistical Association*, 2021.
- [4] B. Barak and A. Moitra. Noisy tensor completion via the sum-of-squares hierarchy. In V. Feldman, A. Rakhlin, and O. Shamir, editors, *Proceedings of the 29th Conference on Learning Theory, COLT 2016, New York, USA, June 23-26, 2016*, volume 49 of *JMLR Workshop and Conference Proceedings*, pages 417–445. JMLR.org, 2016.
- [5] S. Bernstein. *The Theory of Probabilities*. Gastehizdat Publishing House, 1946.
- [6] C. Cai, G. Li, H. V. Poor, and Y. Chen. Nonconvex low-rank tensor completion from noisy data. 32:1863–1874, 2019.
- [7] F. Chollet. keras. <https://github.com/fchollet/keras>, 2015.
- [8] C. Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- [9] Facebook. Prophet. <https://facebook.github.io/prophet/>, 2020. Online; accessed 25 February 2020.
- [10] M. Gavish and D. L. Donoho. The optimal hard threshold for singular values is $4/\sqrt{3}$. *IEEE Transactions on Information Theory*, 60(8):5040–5053, 2014.
- [11] N. Golyandina, V. Nekrutkin, and A. A. Zhigljavsky. *Analysis of time series structure: SSA and related techniques*. Chapman and Hall/CRC, 2001.
- [12] L. Grafakos. *Classical fourier analysis*, volume 2. Springer, 2008.
- [13] H. Hassani, S. Heravi, and A. Zhigljavsky. Forecasting uk industrial production with multivariate singular spectrum analysis. *Journal of Forecasting*, 32(5):395–408, 2013.
- [14] H. Hassani and R. Mahmoudvand. Multivariate singular spectrum analysis: A general view and new vector forecasting approach. *International Journal of Energy and Statistics*, 1(01):55–83, 2013.
- [15] H. Hassani and R. Mahmoudvand. *Singular spectrum analysis: Using R*. Springer, 2018.
- [16] R. J. Hyndman and G. Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2018.
- [17] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. The m5 accuracy competition: Results, findings and conclusions. *Int J Forecast*, 2020.
- [18] V. Oropeza and M. Sacchi. Simultaneous seismic data denoising and reconstruction via multi-channel singular spectrum analysis. *Geophysics*, 76(3):V25–V32, 2011.
- [19] N. Rao, H.-F. Yu, P. K. Ravikumar, and I. S. Dhillon. Collaborative filtering with graph information: Consistency and scalable methods. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems* 28, pages 2107–2115. Curran Associates, Inc., 2015.
- [20] D. S. S. Robert H. Shumway. *Time Series Analysis and It’s Applications*. Blue Printing, 3rd edition, 2015.
- [21] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 2019.
- [22] R. Sen, H.-F. Yu, and I. S. Dhillon. Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting. In *Advances in Neural Information Processing Systems*, pages 4838–4847, 2019.

- [23] D. Shah and C. L. Yu. Iterative collaborative filtering for sparse noisy tensor estimation. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 41–45. IEEE, 2019.
- [24] A. Trindade. UCI machine learning repository - individual household electric power consumption data set. 2014.
- [25] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [26] L. Wasserman. *All of nonparametric statistics*. Springer, 2006.
- [27] P.-Å. Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111, 1972.
- [28] K. W. Wilson, B. Raj, and P. Smaragdis. Regularized non-negative matrix factorization with temporal dependencies for speech denoising. In *Ninth Annual Conference of the International Speech Communication Association*, 2008.
- [29] J. M. Wooldridge. *Econometric Analysis of Cross Section and Panel Data*. The MIT Press, 2010.
- [30] WRDS. The trade and quote (taq) database. 2021.
- [31] D. Xia, M. Yuan, and C.-H. Zhang. Statistically optimal and computationally efficient low rank tensor completion from noisy entries, 2018.
- [32] J. Xu. Rates of convergence of spectral methods for graphon estimation. *arXiv preprint arXiv:1709.03183*, 2017.
- [33] C. L. Yu. Tensor estimation with nearly linear samples. *arXiv preprint arXiv:2007.00736*, 2020.
- [34] H.-F. Yu, N. Rao, and I. S. Dhillon. Temporal regularized matrix factorization for high-dimensional time series prediction. In *Advances in neural information processing systems*, pages 847–855, 2016.

A Page vs. Hankel mSSA

This section discusses the benefits and drawbacks of using the Page matrix representation, as we propose in our variant, instead of the Hankel representation used in the original mSSA. Recall the key steps of the original SSA method in Section 1.3. The extension to mSSA is done by stacking the Hankel matrices induced by each of the N time series either column-wise (horizontal mSSA) or row-wise (vertical mSSA) [15]. In this section, we will use mSSA to denote our mSSA variant, and hSSA/vSSA to denote the original horizontal/vertical mSSA. In what follows, we will compare our mSSA variant with hSSA/vSSA in terms of their: (i) theoretical analysis; (ii) computational complexity; and (iii) empirical performance.

Theoretical analysis. We re-emphasize that to the best of our knowledge, the theoretical analysis of the mSSA algorithm, both hSSA and vSSA, have been absent from the literature, despite their popularity. We do a comprehensive theoretical analysis of the variant of mSSA we propose. By utilizing the Page matrix, it allows us to invoke results from random matrix theory to prove our imputation and forecasting results. However, extending our analysis to the Hankel matrix representation is challenging as the Hankel matrix has repeated entries of the same time series observation. This leads to correlation in the noise in the observation of the entries of the Hankel matrix, which prevents us from invoking the results from random matrix theory in a straightforward way. The Page matrix representation does not have repeated entries of the same observation, and thus allows us to circumvent this issue in our theoretical analysis.

Computational complexity. Our mSSA variant is computationally far more efficient than both hSSA and vSSA. This is because the Page matrix representation of a multivariate time series with N time series and T time steps is a matrix of dimension $\sqrt{NT} \times \sqrt{NT}$ (with $L = \sqrt{NT}$), i.e., it has a total of $\mathcal{O}(NT)$ entries. In contrast, the Hankel matrix representation is of dimension $T/4 \times 3NT/4$ for hSSA and $NT/4 \times 3T/4$ for vSSA¹³, i.e., both variants of the Hankel matrix have $\mathcal{O}(NT^2)$ entries. This makes computing the SVD (the most computationally intensive step of mSSA) prohibitive for hSSA and mSSA even for the standard time series benchmarks we consider in Section 5.

To empirically demonstrate the computational efficiency of our variant of mSSA, we compare its training time to that of hSSA and vSSA. Specifically, we measure the training time for mSSA, hSSA, and vSSA as we increase the number of time steps $T \in [400, 10000]$. We perform this experiment on two datasets: (i) the synthetic dataset; (ii) a subset of the electricity dataset, where we choose only 50 of the available 370 time series. Both datasets are described in details in Appendix B. Figure 8 shows that in both datasets, the training time of both hSSA and vSSA can be as 600-1000x as high as the training time of our mSSA variant as we increase T .

Empirical performance. Here, we compare the forecasting performance of mSSA to that of hSSA and vSSA. We report performance in terms of the NRMSE of the three methods as we increase the number of time steps $T \in [400, 10000]$ in the aforementioned synthetic and electricity dataset. The goal in the synthetic dataset is to predict the next 50 time steps using one step ahead forecasts, while the goal in the electricity dataset is to predict the next three days using day-ahead forecasts. For hSSA and vSSA, we choose $L = T/4$ as recommended in [15]; and for mSSA, we choose $L = \lfloor \sqrt{NT} \rfloor$. For all three methods, we choose the number of retained singular values based on the thresholding procedure outlined in [10].

Figure 9 shows the performance of the three methods in both datasets. We find that initially, with few data points ($T < 600$ in the synthetic data and $T < 4000$ in the electricity data), both hSSA and vSSA outperform mSSA. As we increase T , mSSA performance significantly improves and eventually outperforms vSSA. In the electricity dataset, mSSA performs similar to hSSA for $T = 10000$. These experiments suggest that if only a few observations were available, hSSA and vSSA might provide better performance. However, if the number of observations were relatively large, then the performance of mSSA would be superior to vSSA and relatively similar to hSSA.

¹³We set the parameter L to $T/4$ as recommended in [15].

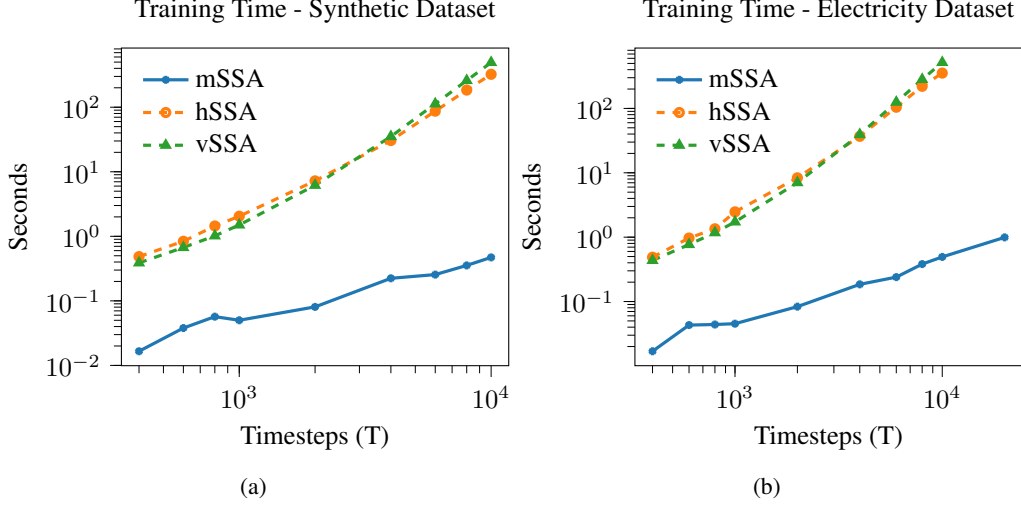


Figure 8: The training time of the original mSSA variants (hSSA in the orange dotted line and vSSA in the green dotted line) are orders of magnitude higher than that of the mSSA variant we propose (blue solid line).

Importantly, the electricity dataset experiment illustrates a critical advantage of our mSSA variant. Specifically, when T is large such that running hSSA or vSSA is computationally infeasible, then one can achieve better accuracy using mSSA. For example, while we could not run the hSSA and vSSA on the electricity dataset with $T = 20000$ due to memory constraints, we were able to run mSSA and achieve a lower NRMSE. This suggests that our mSSA variant is the more practical mSSA algorithm when it comes to efficiently utilizing large multivariate time series.

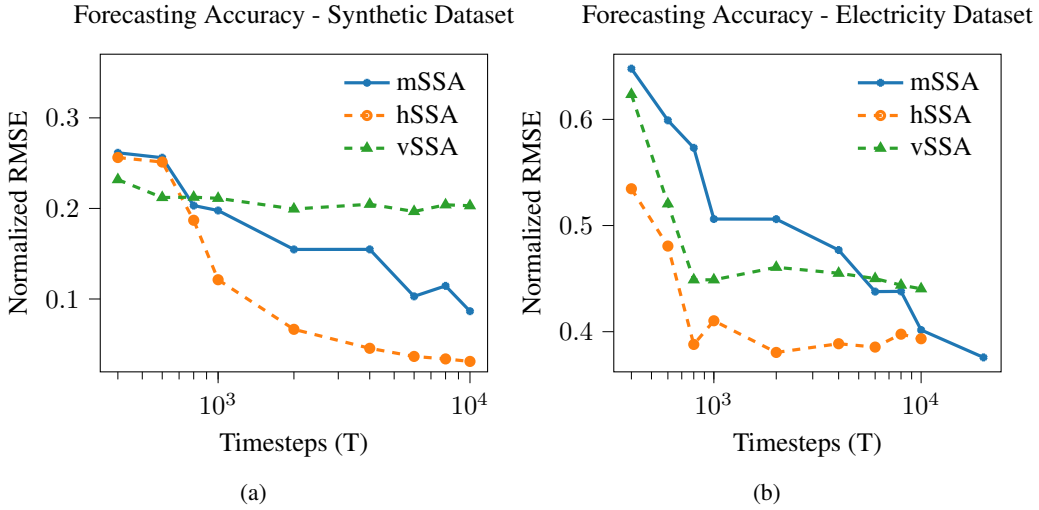


Figure 9: The forecasting error of the original mSSA variants (hSSA in the orange dotted line and vSSA in the green dotted line) and the proposed mSSA variant (blue solid line) as we increase T .

B Experiment Details

In Appendix B.1, we describe the datasets utilized. In Appendix B.2, we describe the various algorithms we compare with as well as the choice of hyper-parameters used for each of them.

Table 4: Dataset and training/validation/test split details.

Dataset	No.time series	Observations per time series	Forecast horizon (h)	Training period	No. validation windows W_{val}	Validation period	No. test windows	Test period
Electricity	370	26136	24	1 to 25824	2	25825 to 25968	7	25969 to 26136
Traffic	963	10560	24	1 to 10248	2	10249 to 10392	7	10393 to 10560
Synthetic	50	15000	10	1 to 13700	10	13701 to 14000	100	14001 to 15000
Financial	839	3993	1	1 to 3693	40	3694 to 3813	180	3814 to 3993
M5	15678	1941	28	1 to 1829	1	1830 to 1913	1	1914 to 1941

B.1 Datasets

We use four real-world datasets and one synthetic dataset. The description and preprocessing we do for each of these datasets are as follows.

Electricity Dataset. This is a public dataset obtained from the UCI repository which shows the 15-minutes electricity load of 370 households [24]. As was done in [34],[22],[21], we aggregate the data into hourly intervals and use the first 25824 time-points for training, the next 288 points for validation, and the last 168 points for testing in the forecasting experiments. Specifically, in our testing period, we do 24-hour ahead forecasts for the next seven days (i.e. 24-step ahead forecast). See Table 4 for more details.

Traffic Dataset. This public dataset obtained from the UCI repository shows the occupancy rate of traffic lanes in San Francisco [24]. The data is sampled every 15 minutes but to be consistent with previous work in [34], [22], we aggregate the data into hourly data and use the first 10248 time-points for training, the next 288 points for validation, and the last 168 points for testing in the forecasting experiments. Specifically, in our testing period, we do 24-hour ahead forecasts for the next seven days (i.e. 24-step ahead forecast). See Table 4 for more details.

Financial Dataset. This dataset is obtained from the Wharton Research Data Services (WRDS) and contains the average daily stocks prices of 839 companies from October 2004 till November 2019 [30]. The dataset was preprocessed to remove stocks with any null values, or those with an average price below 30\$ across the aforementioned period. This was simply done to constrain the number of time series for ease of experimentation and we end up with 839 time series (i.e. stock prices of listed companies) each with 3993 readings of daily stock prices. In our forecasting experiments, we train on the first 3693 time points, validate on the next 120 time points, while for testing we consider the task of predicting 180 time-points ahead one point at a time. That is, the goal here is to do one-day ahead forecasts for the next 180 days (i.e. 1-step ahead forecast). We choose to do so as this is a standard goal in finance. See Table 4 for more details.

M5 Dataset. This public dataset obtained from Kaggle’s M5 Forecasting competition include daily sales data of 30490 items across different Walmart stores for 1941 days [17]. The dataset was preprocessed to only include items that has more than zero sales in at least 500 days. For forecasting, as is the goal in the Kaggle competition, we consider the task of predicting the sales for the next 28 days (i.e. 28-step ahead forecast). We use the first 1829 points for training, the next 84 points for cross validation, and the last 28 points for testing.

Synthetic Dataset. We generate the observation tensor $X \in \mathbb{R}^{n \times m \times T}$ by first randomly generating the two matrices $U \in \mathbb{R}^{r \times n} = [u_1, \dots, u_n]$ and $V \in \mathbb{R}^{r \times m} = [v_1, \dots, v_m]$; we do so by randomly sampling each coordinate of U, V independently from a standard normal. Then, we generate r mixtures of harmonics where each mixture $g_k(t), k \in [r]$, is generated as: $g_k(t) = \sum_{h=1}^4 \alpha_h \cos(\omega_h t/T)$ where the parameters α_h, ω_h are selected uniformly at randomly from the ranges $[-1, 10]$ and $[1, 1000]$, respectively. Then each value in the observation tensor is constructed as follows: $X_{i,j}(t) = \sum_{k=1}^r u_{ik} v_{jk} g_k(t)$, where r is the tensor rank, $i \in [n], j \in [m]$. In our experiment, we select $n = 5, m = 10, T = 15000$, and $r = 4$. This gives us $N = n \times m = 50$ time series each with 15000 observations per time series. In the forecasting experiments, we use the first 13700 points for training, the next 300 points for validation, while for testing, we do 10-step ahead forecasts for the final 1000 points. See Table 4 for more details.

B.2 Algorithms.

In this section, we describe the algorithms used throughout the experiments in more detail and the hyper-parameters/implementation used for each method.

mSSA & SSA. Note that since the SSA’s variant described in [1] is a special case of our proposed mSSA algorithm, we use our mSSA’s implementation to perform the SSA experiments; key difference in SSA is that we do not “stack” the various Page matrices induced by each time series. For all experiments we choose the parameters through the cross validation process detailed in Appendix B.3, where we perform a grid search for the following parameters:

1. *The number of retained singular values, k .* This parameter is chosen using one of the following data-driven methods: (i) we choose k based on the thresholding procedure outlined in [10], where the threshold is determined by the median of the singular values and the shape of the matrix; (ii) we choose k as the minimum number of singular values capturing $> 90\%$ of its spectral energy; (iii) we choose a constant low rank, specifically $k = 3$.
2. *The shape of the Page matrix.* For mSSA, we vary the shape of the Page matrix by choosing $L \in \{500, 1000, 2000, 3000\}$ for the electricity and Traffic datasets, $L \in \{500, 700, 800\}$ for the synthetic dataset, $L \in \{250, 500, 1000, 1500\}$ for the financial dataset, and $L \in \{10, 50, 100, 500\}$ for the M5 dataset. For SSA, we choose $L \in \{50, 100, 150\}$ in the electricity and Traffic datasets, $L \in \{30, 50, 100\}$ in the synthetic dataset, $L \in \{20, 30, 50\}$ in the financial dataset, and $L \in \{5, 10, 20, 40\}$ in the M5 dataset.
3. *Missing values initialization.* Initializing the missing values is done according to one of two methods: (i) set the missing values to zero; (ii) perform forward filling where each missing value is replaced by the nearest preceding observation, followed by backward filling to accommodate the situation when the first observation is missing.

DeepAR. We use the “DeepAREstimator” algorithm provided by the GluonTS package. We choose the parameters through a grid search for the following parameters:

1. *Context length.* This parameter determines the number of steps to unroll the RNN for before computing predictions. We choose this from the set $\{h \text{ (default)}, 2h, 3h\}$, where h is the prediction horizon.
2. *Number of Layers.* This parameter determines the number of RNN layers. We choose this from the set $\{2 \text{ (default)}, 3\}$.

TRMF. We use the implementation provided by the authors in the Github repository associated with the paper ([34]). We choose the parameters through a grid search, as suggested by the authors in their codebase, for the following parameters:

1. *Matrix rank k .* This parameter represents the chosen rank for the $T \times N$ time series matrix, we choose k from the set $\{5, 10, 20, 40, 60\}$.
2. *Regularization parameters $\lambda_f, \lambda_x, \lambda_w$.* We choose these parameters from $\{0.05, 0.5, 5, 50\}$ as suggested in the authors repository.

For the lag indices, we include the last day and the same weekday in the last week for the traffic and electricity data, the last 30 points for the financial and synthetic dataset, and the last 10 points for the M5 dataset.

LSTM. Across all datasets, we use an LSTM network with $H \in \{2, 3, 4\}$ hidden layers each, with 45 neurons per layer, as is done in [22]. We use the Keras implementation of LSTM. As with other methods’ parameters, H is chosen via cross validation.

Prophet. We used Prophet’s Python library with the parameters selected using a grid search of the following parameters as suggested in [9]:

1. *Changepoint prior scale.* This parameter determines how much the trend changes at the detected trend changepoints. We choose this parameter from $\{0.001, 0.05, 0.2\}$.
2. *Seasonality prior scale.* This parameter controls the magnitude of the seasonality. We choose this parameter from $\{0.01, 10\}$.
3. *Seasonality Mode.* Which is chosen to be either ‘additive’ or ‘multiplicative’.

B.3 Parameters Selection

In all experiments, we choose the hyperparameters for our method and for the baselines by using cross-validation. Below, we detail the procedure for both imputation and forecasting experiments.

Imputation Experiments. To select the parameters in our imputation experiments, we additionally mask 10% of the observed data uniformly at random. Then, we evaluate the performance of each parameter choice in recovering these additionally masked observations. This process is repeated 3 times, and the choice of parameters that achieves the best performance (in NRMSE) across these runs is selected. In our results, we report the accuracy of the selected parameters in recovering the original missing values.

Forecasting Experiments. For parameters selection in the forecasting experiments, we use cross-validation on a rolling basis as typically used in time-series forecasting models [16]. In this procedure, there are multiple validation sets. For each validation set, we train the model only on previous observations. That is, no future observations can be used in training the model, which will occur when a typical cross-validation procedure is followed for time series data. In our experiments, we start with a subset of the data used for training, then we forecast the first validation set using h -step ahead forecasts for W_{val} windows, where the horizon h and the number of validation windows W_{val} are detailed in Table 4. We do this for three validation sets, each of length $h \times W_{val}$, and select the choice of parameters that achieves the best performance (in NRMSE) for evaluation on the test set. When evaluating on the test set, both the training and validation periods are used for training.

C Time-varying Recommendation Systems

In Section 6.2, we considered the setting where the $N \times T$ matrix \mathbf{M} induced by the latent time series $f_1(\cdot), \dots, f_N(\cdot)$ is low-rank; in particular, Property 2.1 captures this spatial structure across these N time series. However, in many settings there is *additional* spatial structure across the N time series.

Recommendation systems – time-varying matrices/tensors. For example, in recommendation systems, for each $t \in T$, there is a $N_1 \times N_2$ matrix, $\mathbf{M}^{(t)} \in \mathbb{R}^{N_1 \times N_2}$ of interest. The n_1 -th row and n_2 -th column of $\mathbf{M}^{(t)}$ denotes the latent rating user n_1 has for product n_2 , i.e., $M_{n_1, n_2}^{(t)}$ denotes the value of the latent time series $f_{n_1, n_2}(\cdot)$ at time step t . To capture the latent structure across users and products, one typically assumes that each $\mathbf{M}^{(t)}$ is low-rank. More generally, at each time step t , $\mathbf{M}^{(t)} \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_d}$ could be an order- d tensor. That is, $M_{n_1, \dots, n_d}^{(t)}$ denotes the value of the latent time series $f_{n_1, \dots, n_d}(\cdot)$ at time step t for $n_1, \dots, n_d \in [N_1] \times \dots \times [N_d]$. For example, if $d = 3$, $\mathbf{M}^{(t)}$ might represent the t -th measurement for a collection of (x, y, z) -spatial coordinates. Let $\mathbf{N} \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_d \times T}$ denote the $d + 1$ order tensor induced by viewing each order- d tensor $\mathbf{M}^{(t)}$ as the t -th ‘slice’ of \mathbf{N} , for $t \in [T]$. Again, to capture the spatial and temporal structure of these latent time series, we posit the following spatio-temporal model for \mathbf{N} , which is a higher-order analog of the model assumed in Property 2.1.

Property C.1. Let \mathbf{N} have CP-rank at most R . That is, for any $n_1, \dots, n_d \in [N_1] \times \dots \times [N_d]$

$$N_{n_1, \dots, n_d, t} = \sum_{r=1}^R U_{n_1, r} \dots U_{n_d, r} W_{rt},$$

where the factorization is such that $|U_{n_1, r}|, \dots, |U_{n_d, r}| \leq \Gamma_1$, $|W_{rt}| \leq \Gamma_2$ for constants $\Gamma_1, \Gamma_2 > 0$.

As before, to explicitly model the temporal structure, we continue to assume Property 2.2 holds for the latent time factors W_r , for $r \in [R]$.

Order- $d + 2$ Page tensor representation. We now consider the following order- $d + 2$ Page tensor representation of \mathbf{N} . In particular, given the hyper-parameter $L \geq 1$, define $\mathbf{HT} \in \mathbb{R}^{N_1 \times \dots \times N_d \times T/L \times L}$ such that for $n_1, \dots, n_d \in [N_1] \times \dots \times [N_d]$, $\ell \in [L]$, $s \in [T/L]$,

$$\mathbf{HT}_{n_1, \dots, n_d, \ell, s} = f_{n_1, \dots, n_d}((s-1) \times L + \ell).$$

The corresponding observation tensor, $\mathbb{HT} \in (\mathbb{R} \cup \{\star\})^{N_1 \times \dots \times N_d \times T/L \times L}$, is

$$\mathbb{HT}_{n_1, \dots, n_d, \ell, s} = X_{n_1, \dots, n_d}((s-1) \times L + \ell). \quad (8)$$

Recall from (1) that $X_{n_1, \dots, n_d}(t)$ is the noisy, missing observation we get of $f_{n_1, \dots, n_d}(t)$. \mathbf{HT} and \mathbb{HT} then have the following property:

Proposition C.1. *Let Properties C.1, 2.2, and 2.3 hold. Then, for any $1 \leq L \leq \sqrt{T}$, \mathbf{HT} has CP-rank at most $R \times G$. Further, all entries of \mathbb{HT} are independent random variables with each entry observed with probability $\rho \in (0, 1]$, and $\mathbb{E}[\mathbb{HT}] = \rho \mathbf{HT}$.*

Analogous to Proposition 6.1, Proposition C.1 also establishes that order- $d+2$ Page tensor representation of the various latent time series $f_{n_1, \dots, n_d}(\cdot)$ has CP-rank that continues to be bounded by $R \times G$. Proof of Proposition C.1 can be found in Appendix K.

Higher-order tensor singular spectrum analysis (htSSA). Proposition C.1 motivates the following algorithm, which exploits the further spatial structure amongst the N time series. We now define the “meta” htSSA algorithm. The two algorithmic hyper-parameters are $L \geq 1$ (defined in (7)) and TE_{d+2} (the order- $d+2$ tensor estimation algorithm one chooses). First, using the observations $X_{n_1, \dots, n_d}(t)$ for $n_1, \dots, n_d \in [N_1] \times \dots \times [N_d], t \in [T]$ we construct the higher-order Page tensor \mathbb{HT} as in (8). Second, we obtain $\widehat{\mathbf{HT}}$ as the output of $\text{TE}_{d+2}(\mathbb{HT})$, and read off $\hat{f}_{n_1, \dots, n_d}(t)$ by selecting the appropriate entry in $\widehat{\mathbf{HT}}$.

Relative effectiveness of mSSA, htSSA, and tensor estimation (TE). Again, for ease of exposition, we consider the case where $\rho = 1$. We now briefly discuss the relative effectiveness of htSSA, mSSA, and “vanilla” tensor estimation (TE) in imputing $X_{n_1, \dots, n_d}(\cdot)$ to estimate $f_{n_1, \dots, n_d}(\cdot)$. mSSA and htSSA have been previously described. In TE, one directly de-noises the original order- $d+1$ tensor induced by the noisy observations, which we denote $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_d \times T}$, where $\mathbf{X}_{n_1, \dots, n_d, t} = X_{n_1, \dots, n_d}(t)$. In particular, one produces an estimate of $\widehat{\mathbf{N}} = \text{TE}_{d+1}(\mathbf{X})$, and then produces the estimates $\hat{f}_{n_1, \dots, n_d}(t)$ by reading off the appropriate entry of $\widehat{\mathbf{N}}$. Let $\text{ImpErr}(N, T; \text{htSSA})$, $\text{ImpErr}(N, T; \text{mSSA})$, and $\text{ImpErr}(N, T; \text{TE})$ denote the imputation error for htSSA, mSSA, and TE, respectively. Now if we assume Property 6.4 holds, we have

$$\begin{aligned} \text{ImpErr}(N, T; \text{htSSA}) &= \tilde{\Theta} \left(\frac{1}{\min(N_1, \dots, N_d, \sqrt{T})^{\lceil \frac{d+2}{2} \rceil}} \right), \\ \text{ImpErr}(N, T; \text{mSSA}) &= \tilde{\Theta} \left(\frac{1}{\sqrt{\min(N, T)T}} \right), \\ \text{ImpErr}(N, T; \text{TE}) &= \tilde{\Theta} \left(\frac{1}{\min(N_1, \dots, N_d, T)^{\lceil \frac{d+1}{2} \rceil}} \right). \end{aligned}$$

Then just as was done in the proof of Proposition 6.2, for any given d , one can reason about the relative effectiveness of htSSA, mSSA, and TE for different asymptotic regimes of the relative ratio of N and T .

D Proof of Proposition 4.1

Below, we present the proof of Proposition 4.1. First we define the stacked Hankel matrix of N time series over T time steps. Precisely, given N latent time series f_1, \dots, f_N , consider the stacked Hankel matrix induced by each of them over T time steps, $[T]$, defined as follows. It is $\text{SH} \in \mathbb{R}^{\lfloor T/2 \rfloor \times N \lfloor T/2 \rfloor}$ where its entry in row $i \in [\lfloor T/2 \rfloor]$ and column $j \in [N \lfloor T/2 \rfloor]$, is given by

$$\text{SH}_{ij} = f_{n(i,j)}(i + (j \bmod \lfloor T/2 \rfloor) - 1), \text{ where } n(i,j) = \left\lceil \frac{j}{\lfloor T/2 \rfloor} \right\rceil.$$

We now establish Proposition D.1, which immediately implies Proposition 4.1 – the stacked Page matrix can be viewed as a sub-matrix of SH , by selecting the appropriate columns.

Proposition D.1. *Let Properties 2.1 and 4.1 hold for N latent time series of interest, f_1, \dots, f_N . Then for any $T \geq 1$, the stacked Hankel Matrix of these N time series has ϵ' -approximate rank $R \times G$ with $\epsilon' = R\Gamma_1\epsilon$.*

Proof. We have N latent time series f_1, \dots, f_N satisfying Properties 2.1 and 4.1. Consider their stacked Hankel matrix over $[T]$, $\text{SH} \in \mathbb{R}^{\lfloor T/2 \rfloor \times N \lfloor T/2 \rfloor}$. By definition for $i \in \llbracket \lfloor T/2 \rfloor \rrbracket$ and $j = (n-1) \times \lfloor T/2 \rfloor + j'$ for $j' \in \llbracket \lfloor T/2 \rfloor \rrbracket$, we have

$$\text{SH}_{ij'} = f_n(i + j' - 1).$$

That is,

$$\begin{aligned} \text{SH}_{ij} &= f_n(i + j' - 1) \\ &= \sum_{r=1}^R U_{nr} W_{r(i+j'-1)}. \end{aligned} \quad (9)$$

Let $H(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ be the Hankel matrix associated with W_r over $[T]$. Due to Property 4.1, there exists a low-rank matrix $M(r) \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ such that (a) $\text{rank}(M(r)) \leq G$, (b) $\|H(r) - M(r)\|_\infty \leq \epsilon$. That is, for any $i, j' \in \llbracket \lfloor T/2 \rfloor \rrbracket$, we have that $M(r)_{ij'} = \sum_{g=1}^G a_{ig}^r b_{j'g}^r$ for some $a_{ig}^r, b_{j'g}^r \in \mathbb{R}^G$. Therefore, for any $i, j' \in \llbracket \lfloor T/2 \rfloor \rrbracket$, we have that

$$\begin{aligned} W_{r(i+j'-1)} &= H(r)_{ij'} = M(r)_{ij'} + (H(r)_{ij'} - M(r)_{ij'}) \\ &= \sum_{g=1}^G a_{ig}^r b_{j'g}^r + (H(r)_{ij'} - M(r)_{ij'}). \end{aligned} \quad (10)$$

From (9) and (10), we conclude that

$$\begin{aligned} \text{SH}_{ij} &= \sum_{r=1}^R \sum_{g=1}^G U_{nr} a_{ig}^r b_{j'g}^r + \sum_{r=1}^R U_{nr} (H(r)_{ij'} - M(r)_{ij'}) \\ &= \sum_{(r,g) \in [R] \times [G]} a_{ig}^r \times (U_{nr} b_{j'g}^r) + \sum_{r=1}^R U_{nr} (H(r)_{ij'} - M(r)_{ij'}). \end{aligned}$$

Define matrix $\mathbf{M} \in \mathbb{R}^{\lfloor T/2 \rfloor \times N \lfloor T/2 \rfloor}$ with its entry for row $i \in \llbracket \lfloor T/2 \rfloor \rrbracket$ and column $j = (n-1) \times \lfloor T/2 \rfloor + j'$ for $j' \in \llbracket \lfloor T/2 \rfloor \rrbracket$ given by

$$\begin{aligned} \mathbf{M}_{ij} &= \sum_{(r,g) \in [R] \times [G]} a_{ig}^r \times (U_{nr} b_{j'g}^r) \\ &= \sum_{(r,g) \in [R] \times [G]} \alpha_{i(r,g)} \beta_{j(r,g)}, \end{aligned}$$

where $\alpha_{i(r,g)} = a_{ig}^r$ and $\beta_{j(r,g)} = U_{nr} b_{j'g}^r$. Further,

$$\begin{aligned} |\text{SH}_{ij} - \mathbf{M}_{ij}| &\leq \sum_{r=1}^R |U_{nr}| |H(r)_{ij'} - M(r)_{ij'}| \\ &\leq \sum_{r=1}^R \Gamma_1 \|H(r) - M(r)\|_\infty \leq R\Gamma_1\epsilon. \end{aligned}$$

That is, the stacked Hankel matrix SH of N time series of $[T]$ has ϵ' -approximate rank $G \times R$ with $\epsilon' = R\Gamma_1\epsilon$. This completes the proof. \square

E Proofs For Section 4.2

E.1 Proof of Proposition 4.3

Proof. f_1, f_2 have a (G_1, ϵ_1) and (G_2, ϵ_2) -Hankel representation, respectively. For any $T \geq 1$, let $\mathbf{H}_1, \mathbf{H}_2 \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ be the Hankel matrices of f_1, f_2 , respectively, over the time interval

$[T]$. By definition, there exists matrices $M_1, M_2 \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ such that $\text{rank}(M_1) \leq G_1$, $\|M_1 - H_1\|_\infty \leq \epsilon_1$ and $\text{rank}(M_2) \leq G_2$, $\|M_2 - H_2\|_\infty \leq \epsilon_2$.

Component-wise addition. Note the Hankel matrix of $f_1 + f_2$ over $[T]$ is $H_1 + H_2$. Then, matrix $M = M_1 + M_2$ has rank at most $G_1 + G_2$ since for any two matrices A and B , it is the case that $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$. Further, $\|H_1 + H_2 - (M_1 + M_2)\|_\infty \leq \epsilon_1 + \epsilon_2$. Therefore it follows that $f_1 + f_2$ has $(G_1 + G_2, \epsilon_1 + \epsilon_2)$ -Hankel representation.

Component-wise multiplication. For $f_1 \circ f_2$, its Hankel over $[T]$ is given by $H_1 \circ H_2$ where we abuse notation of \circ in the context of matrices as the Hadamard product of matrices. Let $M = M_1 \circ M_2$. Then $\text{rank}(M) \leq G_1 \times G_2$ since for any two matrices A and B , $\text{rank}(A \circ B) \leq \text{rank}(A)\text{rank}(B)$. Now

$$\begin{aligned} \|H_1 \circ H_2 - M_1 \circ M_2\|_\infty &\leq \|H_1 \circ H_2 - H_1 \circ M_2\|_\infty + \|H_1 \circ M_2 - M_1 \circ M_2\|_\infty \\ &\leq \|H_1\|_\infty \|H_2 - M_2\|_\infty + \|M_2\|_\infty \|H_1 - M_1\|_\infty \\ &\leq \|f_1\|_\infty \epsilon_2 + (\|M_2 - H_2\|_\infty + \|H_2\|_\infty) \epsilon_1 \\ &\leq \|f_1\|_\infty \epsilon_2 + (\|f_2\|_\infty + \epsilon_2) \epsilon_1 \\ &= \|f_1\|_\infty \epsilon_2 + \|f_2\|_\infty \epsilon_1 + \epsilon_1 \epsilon_2 \leq 3 \max(\epsilon_1, \epsilon_2) \max(\|f_1\|_\infty, \|f_2\|_\infty). \end{aligned}$$

This completes the proof of Proposition 4.3. \square

E.2 Proof of Proposition 4.4

Proof. Proof is immediate from Definitions 4.2 and 4.3. \square

E.3 Proof of Proposition 4.5

E.3.1 Helper Lemmas for Proposition 4.5

We begin by stating some classic results from Fourier Analysis. To do so, we introduce some notation. Throughout, we have $R > 0$.

$C[0, R]$ and $L^2[0, R]$ functions. $C[0, R]$ is the set of real-valued, continuous functions defined on $[0, R]$. $L^2[0, R]$ is the set of square integrable functions defined on $[0, R]$, i.e. $\int_0^R f^2(t) dt < \infty$

Inner Product of functions in $L^2[0, R]$. $L^2[0, R]$ is a space endowed with inner product defined as $\langle f, g \rangle := \frac{1}{R} \int_0^R f(t)g(t)dt$, and associated norm as $\|f\| := \sqrt{\frac{1}{R} \int_0^R f^2(t)dt}$.

Fourier Representation of functions in $L^2[0, R]$. For $f \in L^2[0, R]$, define its $G \geq 1$ -order Fourier representation, $\mathcal{F}(f, G) \in L^2[0, R]$ as

$$\mathcal{F}(f, G)(t) = a_0 + \sum_{g=1}^G (a_g \cos(2\pi gt/R) + b_g \sin(2\pi gt/R)), \quad t \in [0, R], \quad (11)$$

where a_0, a_g, b_g with $g \in [G]$ are called the Fourier coefficients of f , defined as

$$\begin{aligned} a_0 &:= \langle f, 1 \rangle = \frac{1}{R} \int_0^R f(t) dt, \\ a_g &:= \langle f, \cos(2\pi gt/R) \rangle = \frac{1}{R} \int_0^R f(t) \cos(2\pi gt/R) dt, \\ b_g &:= \langle f, \sin(2\pi gt/R) \rangle = \frac{1}{R} \int_0^R f(t) \sin(2\pi gt/R) dt. \end{aligned}$$

We now state a classic result from Fourier analysis.

Theorem E.1 ([12]). Given $k \geq 1, R > 0$, let $f \in C^k(R, \mathbb{R})$. Then, for any $t \in [0, R]$ (or more generally $t \in \mathbb{R}$),

$$\lim_{G \rightarrow \infty} \mathcal{F}(f, G)(t) \rightarrow f(t).$$

We next argue that if $f \in C^k(R, \text{PER})$, then its Fourier coefficients decay rapidly.

Lemma E.1. *Given $k \geq 1, R > 0$, let $f \in C^k(R, \text{PER})$. Then, for $j \in [k]$, the G -order Fourier coefficient of $f^{(j)}$, the j -th derivative of f , recursively satisfy the following relationship: for $g \in [G]$,*

$$a_g^{(j)} = -\left(\frac{2\pi g}{R}\right)b_g^{(j-1)}, \quad b_g^{(j)} = \left(\frac{2\pi g}{R}\right)a_g^{(j-1)}. \quad (12)$$

Proof. We establish (12) for $a_g^{(1)}, g \in [G]$. Notice that an identical argument applies to establish (12) for any $a_g^{(j)}, b_g^{(j)}$ for $j \in [k]$ and $g \in [G]$.

$$\begin{aligned} a_g^{(1)} &= \langle f^{(1)}, \cos(2\pi gt/R) \rangle = \frac{1}{R} \int_0^R f^{(1)}(t) \cos(2\pi gt/R) dt \\ &\stackrel{(a)}{=} \frac{1}{R} \left(\left[f(t) \cos(2\pi gt/R) \right]_0^R - \frac{2\pi g}{R} \left[\frac{1}{R} \int_0^R f(t) \sin(2\pi gt/R) dt \right] \right) \\ &= -\left(\frac{2\pi g}{R}\right)b_g^{(0)}. \end{aligned}$$

(a) follows by integration by parts. \square

E.3.2 Completing Proof of Proposition 4.5

Proof. For $G \in \mathbb{N}$, let $\mathcal{F}(f, G)$ be defined as in (11). Then for $t \in \mathbb{R}$

$$\begin{aligned} |f(t) - \mathcal{F}(f, G)(t)| &\stackrel{(a)}{=} \left| \sum_{g=G+1}^{\infty} (a_g \cos(2\pi gt/R) + b_g \sin(2\pi gt/R)) \right| \\ &\leq \sum_{g=G+1}^{\infty} |a_g| + |b_g| \\ &\stackrel{(b)}{\leq} \sum_{g=G+1}^{\infty} \left(\frac{R}{2\pi g}\right)^k (|a_g^{(k)}| + |b_g^{(k)}|) \\ &\stackrel{(c)}{\leq} \sqrt{2} \left(\frac{R}{2\pi}\right)^k \sqrt{\sum_{g=G+1}^{\infty} \left(\frac{1}{g}\right)^{2k}} \sqrt{\sum_{g=G+1}^{\infty} (|a_g^{(k)}|^2 + |b_g^{(k)}|^2)} \\ &\stackrel{(d)}{\leq} \sqrt{2} \left(\frac{R}{2\pi}\right)^k \frac{1}{G^{k-0.5}} \sqrt{\sum_{g=G+1}^{\infty} (|a_g^{(k)}|^2 + |b_g^{(k)}|^2)} \\ &\stackrel{(e)}{\leq} \sqrt{2} \left(\frac{R}{2\pi}\right)^k \frac{\|f^{(k)}\|}{G^{k-0.5}} \\ &= C(k, R) \frac{\|f^{(k)}\|}{G^{k-0.5}}, \end{aligned}$$

where $C(k, R)$ is a constant that depends only on k and R ; (a) follows from Theorem E.1; (b) follows from Lemma E.1; (c) follows from Cauchy-Schwarz inequality and fact that $(\alpha + \beta)^2 \leq 2(\alpha^2 + \beta^2)$ for any $\alpha, \beta \in \mathbb{R}$; (d) $\sum_{g=G+1}^{\infty} g^{-2k} \leq \int_G^{\infty} x^{-2k} dx$ which can be bounded as $G^{-2k+1}/(2k-1)$ which is at most G^{-2k+1} since $k \geq 1$; (e) follows from Bessel's inequality, i.e. $\|f^{(k)}\|^2 \geq \sum_{g=0}^{\infty} (|a_g^{(k)}|^2 + |b_g^{(k)}|^2)$.

Thus, for any $t \in \mathbb{R}$, we have a uniform error bound for f being approximated by $\mathcal{F}(f, G)$ which is a sum of $2G$ harmonics. Noting $2G$ harmonics can be represented by an order- $4G$ LRF (by Proposition 2.1), we complete the proof. \square

E.4 Proof of Proposition 4.6

This analysis is adapted from [32].

Proof. Step 1: Partitioning the space $[0, 1)^K$. Consider an equal partition of $[0, 1)^K$. Precisely, for any $k \in \mathbb{N}$, we partition the set $[0, 1)$ into $1/k$ half-open intervals of length $1/k$, i.e., $[0, 1) = \cup_{i=1}^k [(i-1)/k, i/k)$. It follows that $[0, 1)^K$ can be partitioned into k^K cubes of forms $\otimes_{j=1}^K [(i_j - 1)/k, i_j/k)$ with $i_j \in [k]$. Let \mathcal{E}_k be such a partition with I_1, I_2, \dots, I_{k^K} denoting all such cubes and $z_1, z_2, \dots, z_{k^K} \in \mathbb{R}^K$ denoting the centers of those cubes.

Step 2: Taylor Expansion of $g(\cdot, \omega)$. Consider a fixed ω . To reduce notational overload, we suppress dependence of g on ω , and abuse notation by using $g(\cdot) = g(\cdot, \omega)$ in what follows.

For every I_i with $1 \leq i \leq k^K$, define $P_{I_i, \ell}(x)$ as the degree- ℓ Taylor's series expansion of $g(x)$ at point z_i :

$$P_{I_i, \ell}(x) = \sum_{\kappa: |\kappa| \leq \ell} \frac{1}{\kappa!} (x - z_i)^\kappa \nabla_\kappa g(z_i), \quad (13)$$

where $\kappa = (\kappa_1, \dots, \kappa_d)$ is a multi-index with $\kappa! = \prod_{i=1}^K \kappa_i!$, and $\nabla_\kappa g(z_i)$ is the partial derivative defined in Section 4.3. Note similar to g , $P_{I_i, \ell}(x)$ really refers to $P_{I_i, \ell}(x, \omega)$.

Now we define a degree- ℓ piecewise polynomial

$$P_{\mathcal{E}_k, \ell}(x) = \sum_{i=1}^{k^K} P_{I_i, \ell}(x) \mathbb{1}(x \in I_i).$$

For the remainder of the proof, let $\ell = \lfloor \alpha \rfloor$ (recall $\lfloor \alpha \rfloor$ refers to the largest integer strictly smaller than α). Since $f \in \mathcal{H}(\alpha, L)$, it follows that

$$\begin{aligned} \sup_{x \in [0, 1)^K} |g(x) - P_{\mathcal{E}_k, \ell}(x)| &= \max_{1 \leq i \leq k^K} \sup_{x \in I_i} |g(x) - P_{I_i, \ell}(x)| \\ &\stackrel{(a)}{=} \max_{1 \leq i \leq k^K} \sup_{x \in I_i} \left| \sum_{\kappa: |\kappa| \leq \ell-1} \frac{\nabla_\kappa g(z_i)}{\kappa!} (x - z_i)^\kappa + \sum_{\kappa: |\kappa| = \ell} \frac{\nabla_\kappa g(\tilde{z}_i)}{\kappa!} (x - z_i)^\ell - P_{I_i, \ell}(x) \right| \\ &\stackrel{(b)}{=} \max_{1 \leq i \leq k^K} \sup_{x \in I_i} \left| \sum_{\kappa: |\kappa| = \ell} \frac{\nabla_\kappa g(\tilde{z}_i)}{\kappa!} (x - z_i)^\ell - \sum_{\kappa: |\kappa| = \ell} \frac{\nabla_\kappa g(z_i)}{\kappa!} (x - z_i)^\ell \right| \\ &= \max_{1 \leq i \leq k^K} \sup_{x \in I_i} \left| \sum_{\kappa: |\kappa| = \ell} \frac{\nabla_\kappa g(\tilde{z}_i) - \nabla_\kappa g(z_i)}{\kappa!} (x - z_i)^\ell \right| \\ &\stackrel{(c)}{\leq} \max_{1 \leq i \leq k^K} \sup_{x \in I_i} \|x - z_i\|_\infty^\ell \sup_{x \in I_i} \sum_{\kappa: |\kappa| = \ell} \frac{1}{\kappa!} |\nabla_\kappa g(\tilde{z}_i) - \nabla_\kappa g(z_i)| \\ &\stackrel{(d)}{\leq} \mathcal{L} k^{-\alpha}. \end{aligned} \quad (14)$$

where (a) follows from multivariate version of Taylor's theorem (and using the Lagrange form for the remainder) and $\tilde{z}_i \in [0, 1)^K$ is a vector that can be represented as $z_i + cx$ for $c \in (0, 1)$; (b) follows from (13); (c) follows from Holder's inequality; (d) follows from Definition 4.5.

Step 3: Construct Low-Rank Approximation of Time Series Hankel Using $P_{\mathcal{E}_k, \ell}$. Recall the Hankel matrix, $\mathbf{H} \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$ induced by the original time series over $[T]$, where $\mathbf{H}_{ts} = g(\theta_t, \omega_s)$, $t, s \in [\lfloor T/2 \rfloor]$ with $g(\cdot, \omega) \in \mathcal{H}(\alpha, \mathcal{L})$ for any ω . We now construct a low-rank approximation of it using $P_{\mathcal{E}_k, \ell} = P_{\mathcal{E}_k, \ell}(\cdot, \omega)$. Define $\widetilde{\mathbf{H}} \in \mathbb{R}^{\lfloor T/2 \rfloor \times \lfloor T/2 \rfloor}$, where $\widetilde{\mathbf{H}}_{ts} = P_{\mathcal{E}_k, \ell}(\theta_t, \omega_s)$, $t, s \in [\lfloor T/2 \rfloor]$.

By (14), we have that for all $t, s \in [\lfloor T/2 \rfloor]$,

$$|\mathbf{H}_{ts} - \widetilde{\mathbf{H}}_{ts}| \leq \mathcal{L} k^{-\alpha}.$$

It remains to bound the rank of $\widetilde{\mathbf{H}}$. Note that since $P_{\mathcal{E}_k, \ell}(\cdot, \omega)$ is a piecewise polynomial of degree $\ell = \lfloor \alpha \rfloor$ for any given ω , it has the following decomposition: for $t, s \in \llbracket T/2 \rrbracket$,

$$\widetilde{\mathbf{H}}_{ts} = P_{\mathcal{E}_k, \ell}(\theta_t, \omega_s) = \sum_{i=1}^{k^K} \langle \Phi(\theta_t), \beta_{I_i, s} \rangle \mathbb{1}(\theta_t \in I_i)$$

where for any $\theta \in \mathbb{R}^K$,

$$\Phi(\theta) = \left(1, \theta_1, \dots, \theta_K, \dots, \theta_1^\ell, \dots, \theta_K^\ell \right)^T,$$

the vector of all monomials of degree less than or equal to ℓ , and $\beta_{I_i, s}$ is a vector collecting the corresponding coefficients. The number of such monomials is easily show to be equal to $C(\alpha, K) := \sum_{i=1}^{\lfloor \alpha \rfloor} \binom{i+K-1}{i}$. That is, $\widetilde{\mathbf{H}}_{ts} = u_t^T v_s$ where u_t, v_s are of dimension at most $k^K C(\alpha, K)$ for each $t, s \in \llbracket T/2 \rrbracket$. That is, $\widetilde{\mathbf{H}}$ has rank at most $k^K C(\alpha, K)$. Setting $k = \left\lceil \frac{1}{\epsilon} \right\rceil$ completes the proof. \square

F Concentration Inequalities

We recall known concentration inequalities that will be useful throughout.

Theorem F.1 (Bernstein's Inequality [5]). *Suppose that X_1, \dots, X_n are independent random variables with zero mean, and M is a constant such that $|X_i| \leq M$ with probability one for each i . Let $S := \sum_{i=1}^n X_i$ and $v := \text{Var}(S)$. Then for any $t \geq 0$,*

$$\mathbb{P}(|S| \geq t) \leq 2 \exp\left(-\frac{3t^2}{6v + 2Mt}\right).$$

Theorem F.2 (Norm of matrices with sub-gaussian entries [25]). *Let \mathbf{A} be an $m \times n$ random matrix whose entries A_{ij} are independent, mean zero, sub-gaussian random variables. Then, for any $t > 0$, we have*

$$\|\mathbf{A}\| \leq CK(\sqrt{m} + \sqrt{n} + t)$$

with probability at least $1 - 2 \exp(-t^2)$. Here, $K = \max_{i,j} \|A_{ij}\|_{\psi_2}$.

Lemma F.1 (Maximum of sequence of random variables [25]). *Let X_1, X_2, \dots, X_n be a sequence of random variables, which are not necessarily independent, and satisfy $\mathbb{E}[X_i^{2p}]^{\frac{1}{2p}} \leq Kp^{\frac{\beta}{2}}$ for some $K, \beta > 0$ and all i . Then, for every $n \geq 2$,*

$$\mathbb{E} \max_{i \leq n} |X_i| \leq CK \log^{\frac{\beta}{2}}(n).$$

We note that Lemma F.1 implies that if X_1, \dots, X_n are ψ_α random variables with $\|X_i\|_{\psi_\alpha} \leq K_\alpha$ for all $i \in [n]$, then

$$\mathbb{E} \max_{i \leq n} |X_i| \leq CK_\alpha \log^{\frac{1}{\alpha}}(n).$$

G Matrix Estimation via HSVT

This section describes and analyzes a well-known matrix estimation method, Hard Singular Value Thresholding (HSVT). While the analysis utilizes known arguments from the literature, we need to adapt it for the setting where the underlying ‘signal’ is only approximately low-rank.

G.1 Setup, Notations

Setup. Given a deterministic matrix $\mathbf{M} \in \mathbb{R}^{q \times p}$ with $p, q \in \mathbb{N}$ and $q \leq p$, a random matrix $\mathbf{Y} \in \mathbb{R}^{q \times p}$ is such that all of its entries, Y_{ij} , $i \in [q]$, $j \in [p]$ are mutually independent and for any given $i \in [q]$, $j \in [p]$,

$$Y_{ij} = \begin{cases} M_{ij} + \epsilon_{ij} & \text{w.p. } \rho, \text{ (i.e. observed)} \\ 0 & \text{w.p. } 1 - \rho, \text{ (i.e. not observed)} \end{cases}$$

for some $\rho \in (0, 1]$ with ε_{ij} are independent random variables with $\mathbb{E}[\varepsilon_{ij}] = 0$ and $\|\varepsilon_{ij}\|_{\psi_2} \leq \sigma$. Given this, we have $\mathbb{E}[\mathbf{Y}] = \rho \mathbf{M}$. Define

$$\hat{\rho} = \max \left(1/(q p), \left(\sum_{i=1}^q \sum_{j=1}^p \mathbf{1}(Y_{ij} \text{ is obs.}) \right) / (q p) \right).$$

Goal of Matrix Estimation. The goal of matrix estimation is to produce an estimate $\widehat{\mathbf{M}}$ from observation \mathbf{Y} so that $\widehat{\mathbf{M}}$ is close to \mathbf{M} . In particular, we will be interested in bounding the error between $\widehat{\mathbf{M}}$ and \mathbf{M} using the following metric: $\|\widehat{\mathbf{M}} - \mathbf{M}\|_{2,\infty}$.

G.2 Matrix Estimation using HSVT

Hard Singular Value Thresholding (HSVT) Map. We define the HSVT map. For any $q, p \in \mathbb{N}$, consider a matrix $\mathbf{B} \in \mathbb{R}^{q \times p}$ such that $\mathbf{B} = \sum_{i=1}^{q \wedge p} \sigma_i(\mathbf{B}) x_i y_i^T$. Here for $i \in [q \wedge p]$, $\sigma_i(\mathbf{B})$ is the i th largest singular value of \mathbf{B} and x_i, y_i are the corresponding left and right singular vectors respectively. Then, for given any $\lambda > 0$, we define the map $\text{HSVT}_\lambda : \mathbb{R}^{q \times p} \rightarrow \mathbb{R}^{q \times p}$, which simply shaves off the singular values of the input matrix that are below the threshold λ . Precisely,

$$\text{HSVT}_\lambda(\mathbf{B}) = \sum_{i=1}^{q \wedge p} \sigma_i(\mathbf{B}) \mathbf{1}(\sigma_i(\mathbf{B}) \geq \lambda) x_i y_i^T.$$

Matrix Estimating using HSVT map. We define a matrix estimation method using the HSVT map that is utilized by mSSA for imputation. Precisely, we estimate \mathbf{M} from \mathbf{Y} as follows: given parameter $k \geq 1$,

$$\widehat{\mathbf{M}} = \frac{1}{\hat{\rho}} \text{HSVT}_{\lambda_k}(\mathbf{Y}). \quad (15)$$

where $\lambda_k = \sigma_k(\mathbf{Y})$, i.e. the k th largest singular value of \mathbf{Y} .

G.3 A Useful Linear Operator

We define a linear map associated to HSVT. For a specific choice of $\lambda \geq 0$, define $\varphi_\lambda^{\mathbf{B}} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ as follows: for any vector $w \in \mathbb{R}^p$ (i.e. $w \in \mathbb{R}^{p \times 1}$),

$$\varphi_\lambda^{\mathbf{B}}(w) = \sum_{i=1}^{q \wedge p} \mathbf{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i y_i^T w. \quad (16)$$

Note that $\varphi_\lambda^{\mathbf{B}}$ is a linear operator and it depends on the tuple (\mathbf{B}, λ) ; more precisely, the singular values and the right singular vectors of \mathbf{B} , as well as the threshold λ . If $\lambda = 0$, then we will adopt the shorthand notation: $\varphi^{\mathbf{B}} = \varphi_0^{\mathbf{B}}$. The following is a simple, but curious relationship between $\varphi_\lambda^{\mathbf{B}}$ and HSVT_λ that will be useful subsequently.

Lemma G.1 (Lemma 35 of [2, 3]). *Let $\mathbf{B} \in \mathbb{R}^{q \times p}$ and $\lambda \geq 0$ be given. Then for any $j \in [q]$,*

$$\varphi_\lambda^{\mathbf{B}}(\mathbf{B}_{j\cdot}^T) = \text{HSVT}_\lambda(\mathbf{B})_{j\cdot}^T,$$

where $\mathbf{B}_{j\cdot} \in \mathbb{R}^{1 \times p}$ represents the j th row of \mathbf{B} , and $\text{HSVT}_\lambda(\mathbf{B})_{j\cdot} \in \mathbb{R}^{1 \times p}$ represents the j th row of the matrix obtained after applying HSVT over \mathbf{B} with threshold λ .

Proof. By (16), the orthonormality of the right singular vectors and noting $\mathbf{B}_{j\cdot}^T = \mathbf{B}^T \mathbf{e}_j$ with $\mathbf{e}_j \in \mathbb{R}^p$ with j th entry 1 and everything else 0, we have

$$\begin{aligned}
\varphi_\lambda^B(\mathbf{B}_{j\cdot}^T) &= \sum_{i=1}^{q \wedge p} \mathbb{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i y_i^T \mathbf{B}_{j\cdot}^T = \sum_{i=1}^{q \wedge p} \mathbb{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i y_i^T \mathbf{B}^T \mathbf{e}_j \\
&= \sum_{i=1}^{q \wedge p} \mathbb{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i y_i^T \left(\sum_{i'=1}^{q \wedge p} \sigma_{i'}(\mathbf{B}) x_{i'} y_{i'}^T \right)^T \mathbf{e}_j = \sum_{i,i'=1}^{q \wedge p} \sigma_{i'}(\mathbf{B}) \mathbb{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i y_i^T y_{i'} x_{i'}^T \mathbf{e}_j \\
&= \sum_{i,i'=1}^{q \wedge p} \sigma_{i'}(\mathbf{B}) \mathbb{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i \delta_{ii'} x_{i'}^T \mathbf{e}_j = \sum_{i=1}^{q \wedge p} \sigma_i(\mathbf{B}) \mathbb{1}(\sigma_i(\mathbf{B}) \geq \lambda) y_i x_i^T \mathbf{e}_j \\
&= \text{HSVT}_\lambda(\mathbf{B})^T \mathbf{e}_j = \text{HSVT}_\lambda(\mathbf{B})_{j\cdot}^T.
\end{aligned}$$

□

G.4 HSVT based Matrix Estimation: A Deterministic Bound

We state the following result about property of the estimator.

Lemma G.2. For $k \geq 1$, let $\mathbf{M} = \mathbf{M}_k + \mathbf{E}_k$ with $\text{rank}(\mathbf{M}_k) = k$. Let $\varepsilon = \max(\hat{\rho}/\rho, \rho/\hat{\rho}) \geq 1$. Then, the HSVT estimate $\widehat{\mathbf{M}}$ with parameter k is such that for all $j \in [q]$,

$$\begin{aligned}
\|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^2 &\leq \frac{2\|\mathbf{Y} - \rho\mathbf{M}\|_2^2 + 2\rho^2\|\mathbf{E}_k\|_2^2}{(\sigma_k(\rho\mathbf{M}_k))^2} \left(2\|[\mathbf{M}_k]_{j\cdot}^T\|_2^2 + \frac{4\varepsilon^2(\|\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T\|_2)^2}{\rho^2} \right) \\
&\quad + \frac{4\varepsilon^2}{\rho^2} \left\| \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T) \right\|_2^2 + 2(\varepsilon - 1)^2 \|\mathbf{M}_{j\cdot}^T\|_2^2 + 2\|[\mathbf{E}_k]_{j\cdot}^T\|_2^2. \quad (17)
\end{aligned}$$

Proof. We prove our lemma in four steps.

Step 1. Decomposing $\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T$ in two terms. Fix a row index $j \in [q]$. Let λ_k be the k th largest singular value of \mathbf{Y} , as used by HSVT algorithm with parameter $k \geq 1$.

$$\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T = \left(\widehat{\mathbf{M}}_{j\cdot}^T - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) \right) + \left(\varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) - \mathbf{M}_{j\cdot}^T \right).$$

By definition per (16), $\varphi_{\lambda_k}^{\mathbf{Y}} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is the projection operator onto $\text{span}\{u_1, \dots, u_k\}$, the span of top k right singular vectors of \mathbf{Y} , denoted as u_1, \dots, u_k . Therefore,

$$\varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) - \mathbf{M}_{j\cdot}^T \in \text{span}\{u_1, \dots, u_k\}^\perp.$$

By design, $\text{rank}(\widehat{\mathbf{M}}) = k$. Therefore, by Lemma G.1

$$\widehat{\mathbf{M}}_{j\cdot} - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) = \frac{1}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T) - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) \in \text{span}\{u_1, \dots, u_k\}.$$

Therefore, $\langle \widehat{\mathbf{M}}_{j\cdot}^T - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T), \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) - \mathbf{M}_{j\cdot}^T \rangle = 0$, and hence

$$\left\| \widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T \right\|_2^2 = \left\| \widehat{\mathbf{M}}_{j\cdot}^T - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) \right\|_2^2 + \left\| \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) - \mathbf{M}_{j\cdot}^T \right\|_2^2 \quad (18)$$

by the Pythagorean theorem.

Step 2. Bounding Term 1, $\left\| \widehat{\mathbf{M}}_{j\cdot}^T - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) \right\|_2$. We begin by bounding the first term on the right hand side of (18). By Lemma G.1,

$$\begin{aligned}
\widehat{\mathbf{M}}_{j\cdot} - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) &= \frac{1}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T) - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) = \varphi_{\lambda_k}^{\mathbf{Y}} \left(\frac{1}{\widehat{\rho}} \mathbf{Y}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T \right) \\
&= \frac{1}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T) + \frac{\rho - \widehat{\rho}}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T).
\end{aligned}$$

Using the Parallelogram Law (or, equivalently, combining Cauchy-Schwartz and AM-GM inequalities), we obtain

$$\begin{aligned}
\|\widehat{\mathbf{M}}_{j\cdot}^T - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T)\|_2^2 &= \left\| \frac{1}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) + \frac{\rho - \widehat{\rho}}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) \right\|_2^2 \\
&\leq 2 \left\| \frac{1}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) \right\|_2^2 + 2 \left\| \frac{\rho - \widehat{\rho}}{\widehat{\rho}} \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T) \right\|_2^2 \\
&\leq \frac{2}{\widehat{\rho}^2} \|\varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T)\|_2^2 + 2 \left(\frac{\rho - \widehat{\rho}}{\widehat{\rho}} \right)^2 \|\mathbf{M}_{j\cdot}^T\|_2^2 \\
&\leq \frac{2\varepsilon^2}{\rho^2} \|\varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T)\|_2^2 + 2(\varepsilon - 1)^2 \|\mathbf{M}_{j\cdot}^T\|_2^2. \tag{19}
\end{aligned}$$

From definition of ε , $\frac{1}{\widehat{\rho}} \leq \frac{\varepsilon}{\rho}$ and $\left(\frac{\rho - \widehat{\rho}}{\widehat{\rho}} \right)^2 \leq (\varepsilon - 1)^2$. The first term of (19) can be decomposed as,

$$\begin{aligned}
&\|\varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T)\|_2^2 \\
&\leq 2 \left\| \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) - \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) \right\|_2^2 + 2 \left\| \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) \right\|_2^2. \tag{20}
\end{aligned}$$

In above, we have used notation $\varphi^{\mathbf{M}_k} = \varphi_0^{\mathbf{M}_k}$. Given that \mathbf{M}_k is rank k matrix, $\varphi^{\mathbf{M}_k} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is the projection operator mapping any element in \mathbb{R}^p to the projection onto the subspace spanned by $\{\mu_1, \dots, \mu_k\}$, where $\mu_1, \dots, \mu_k \in \mathbb{R}^p$ are the k non-trivial right singular vectors of \mathbf{M}_k . Similarly, by definition $\varphi_{\lambda_k}^{\mathbf{Y}}$ is a map $\mathbb{R}^p \rightarrow \mathbb{R}^p$ mapping any element in \mathbb{R}^p to its projection onto the subspace spanned by $\{u_1, \dots, u_k\}$, the top k right singular vectors of \mathbf{Y} —this can be seen by noting $\lambda_k = \sigma_k(\mathbf{Y})$ is the k -th top singular value of \mathbf{Y} . Recall $\sigma_j(\mathbf{Y})$, $j \in [q \wedge p]$ is the j th largest singular value of \mathbf{Y} .

Next, we bound the first term on the right hand side of (20). To that end, by Wedin sin Θ Theorem (see [8, 27]) and recalling $\text{rank}(\mathbf{M}_k) = k$,

$$\begin{aligned}
\|\varphi_{\lambda_k}^{\mathbf{Y}} - \varphi^{\mathbf{M}_k}\|_2 &\leq \frac{\|\mathbf{Y} - \rho \mathbf{M}_k\|_2}{\sigma_k(\rho \mathbf{M}_k)} \\
&\leq \frac{\|\mathbf{Y} - \rho \mathbf{M}\|_2}{\sigma_k(\rho \mathbf{M}_k)} + \frac{\rho \|\mathbf{M} - \mathbf{M}_k\|_2}{\sigma_k(\rho \mathbf{M}_k)} \\
&\leq \frac{\|\mathbf{Y} - \rho \mathbf{M}\|_2}{\sigma_k(\rho \mathbf{M}_k)} + \frac{\rho \|\mathbf{E}_k\|_2}{\sigma_k(\rho \mathbf{M}_k)}. \tag{21}
\end{aligned}$$

Then it follows that

$$\begin{aligned}
\left\| \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) - \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) \right\|_2 &\leq \|\varphi_{\lambda_k}^{\mathbf{Y}} - \varphi^{\mathbf{M}_k}\|_2 \|\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T\|_2 \\
&\leq \frac{(\|\mathbf{Y} - \rho \mathbf{M}\|_2 + \rho \|\mathbf{E}_k\|_2) (\|\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T\|_2)}{\sigma_k(\rho \mathbf{M}_k)}. \tag{22}
\end{aligned}$$

Using (20) and (22) in (19),

$$\begin{aligned}
\|\widehat{\mathbf{M}}_{j\cdot}^T - \varphi_{\lambda_k}^{\mathbf{Y}}(\mathbf{M}_{j\cdot}^T)\|_2^2 &\leq \frac{4\varepsilon^2}{\rho^2} \frac{(\|\mathbf{Y} - \rho \mathbf{M}\|_2 + \rho \|\mathbf{E}_k\|_2)^2 (\|\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T\|_2)^2}{(\sigma_k(\rho \mathbf{M}_k))^2} \\
&\quad + \frac{4\varepsilon^2}{\rho^2} \left\| \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) \right\|_2^2 + 2(\varepsilon - 1)^2 \|\mathbf{M}_{j\cdot}^T\|_2^2. \tag{23}
\end{aligned}$$

Step 3. Bounding Term 2, $\left\| \varphi_{\lambda_k}^Y(\mathbf{M}_{j\cdot}^T) - \mathbf{M}_{j\cdot}^T \right\|_2^2$. Recall $\mathbf{M} = \mathbf{M}_k + \mathbf{E}_k$ and using (21),

$$\begin{aligned}
\left\| \varphi_{\lambda_k}^Y(\mathbf{M}_{j\cdot}^T) - \mathbf{M}_{j\cdot}^T \right\|_2^2 &= \left\| \varphi_{\lambda_k}^Y([\mathbf{M}_k]_{j\cdot}^T + [\mathbf{E}_k]_{j\cdot}^T) - [\mathbf{M}_k]_{j\cdot}^T - [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2 \\
&\leq 2 \left\| \varphi_{\lambda_k}^Y([\mathbf{M}_k]_{j\cdot}^T) - [\mathbf{M}_k]_{j\cdot}^T \right\|_2^2 + 2 \left\| \varphi_{\lambda_k}^Y([\mathbf{E}_k]_{j\cdot}^T) - [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2 \\
&= 2 \left\| \varphi_{\lambda_k}^Y([\mathbf{M}_k]_{j\cdot}^T) - \varphi_{\lambda_k}^{\mathbf{M}_k}([\mathbf{M}_k]_{j\cdot}^T) \right\|_2^2 + 2 \left\| \varphi_{\lambda_k}^Y([\mathbf{E}_k]_{j\cdot}^T) - [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2 \\
&\leq 2 \left\| \varphi_{\lambda_k}^Y - \varphi_{\lambda_k}^{\mathbf{M}_k} \right\|_2^2 \left\| [\mathbf{M}_k]_{j\cdot}^T \right\|_2^2 + 2 \left\| [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2 \\
&\leq 2 \frac{(\|\mathbf{Y} - \rho\mathbf{M}\|_2 + \rho\|\mathbf{E}_k\|_2)^2}{(\sigma_k(\rho\mathbf{M}_k))^2} \left\| [\mathbf{M}_k]_{j\cdot}^T \right\|_2^2 + 2 \left\| [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2. \tag{24}
\end{aligned}$$

Step 4. Putting everything together. Inserting (23) and (24) back to (18), we have that for each $j \in [q]$,

$$\begin{aligned}
\left\| \widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T \right\|_2^2 &\leq 2 \frac{(\|\mathbf{Y} - \rho\mathbf{M}\|_2 + \rho\|\mathbf{E}_k\|_2)^2}{(\sigma_k(\rho\mathbf{M}_k))^2} \left\| [\mathbf{M}_k]_{j\cdot}^T \right\|_2^2 + 2 \left\| [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2 \\
&\quad + \frac{4\varepsilon^2 (\|\mathbf{Y} - \rho\mathbf{M}\|_2 + \rho\|\mathbf{E}_k\|_2)^2 (\|\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T\|_2)^2}{\rho^2 (\sigma_k(\rho\mathbf{M}_k))^2} \\
&\quad + \frac{4\varepsilon^2}{\rho^2} \left\| \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T) \right\|_2^2 + 2(\varepsilon - 1)^2 \|\mathbf{M}_{j\cdot}^T\|_2^2 \\
&\leq \frac{2\|\mathbf{Y} - \rho\mathbf{M}\|_2^2 + 2\rho^2\|\mathbf{E}_k\|_2^2}{(\sigma_k(\rho\mathbf{M}_k))^2} \left(2 \left\| [\mathbf{M}_k]_{j\cdot}^T \right\|_2^2 + \frac{4\varepsilon^2 (\|\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T\|_2)^2}{\rho^2} \right) \\
&\quad + \frac{4\varepsilon^2}{\rho^2} \left\| \varphi^{\mathbf{M}_k}(\mathbf{Y}_{j\cdot}^T - \rho\mathbf{M}_{j\cdot}^T) \right\|_2^2 + 2(\varepsilon - 1)^2 \|\mathbf{M}_{j\cdot}^T\|_2^2 + 2 \left\| [\mathbf{E}_k]_{j\cdot}^T \right\|_2^2,
\end{aligned}$$

where we used $(a + b)^2 \leq 2a^2 + 2b^2$. This completes the proof. \square

G.5 HSVT based Matrix Estimation: Deterministic To High-Probability

Next, we convert the bound obtained in Lemma G.2 to a bound in expectation (as well as one in high-probability) for our metric of interest: $\|\widehat{\mathbf{M}} - \mathbf{M}\|_{2,\infty}$. In particular, we establish

Theorem G.1. *For $k \geq 1$, let $\mathbf{M} = \mathbf{M}_k + \mathbf{E}_k$ with $\text{rank}(\mathbf{M}_k) = k$. Let $\epsilon = \|\mathbf{E}_k\|_\infty$ and $\Gamma = \|\mathbf{M}_k\|_\infty$. Let $\rho \geq C \log(qp)/q$ for C large enough and $q \leq p$. Then, the HSVT estimate $\widehat{\mathbf{M}}$ with parameter k is such that*

$$\mathbb{E} \left[\max_{j \in [q]} \frac{1}{p} \left\| \widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T \right\|_2^2 \right] \leq \frac{p(C\sigma^2 + \rho^2\epsilon q)}{\rho^2\sigma_k(\mathbf{M}_k)^2} \left(\Gamma^2 + \frac{\sigma^2}{\rho^2} \right) + \frac{C\sigma^2 k \log p}{p\rho^2} + \frac{C(\Gamma + \epsilon)^2}{p} + 2\epsilon^2 + \frac{C}{(pq)^2}.$$

Proof. We start by identifying certain high probability events. Subsequently, using these events and Lemma G.2, we shall conclude the proof.

High Probability Events. For some positive absolute constant $C > 0$, define

$$E_1 := \left\{ |\hat{\rho} - \rho| \leq \rho/20 \right\},$$

$$E_2 := \left\{ \|\mathbf{Y} - \rho \mathbf{M}\|_2 \leq C\sigma\sqrt{p} \right\}, \quad (25)$$

$$E_3 := \left\{ \|\mathbf{Y} - \rho \mathbf{M}\|_{\infty,2}, \|\mathbf{Y} - \rho \mathbf{M}\|_{2,\infty} \leq C\sigma\sqrt{p} \right\}, \quad (26)$$

$$E_4 := \left\{ \max_{j \in [q]} \|\varphi_{\sigma_k(\mathbf{B})}^{\mathbf{B}}(\mathbf{Y}_j^T - \rho \mathbf{M}_j^T)\|_2^2 \leq C\sigma^2 k \log(p) \right\}, \quad (27)$$

$$E_5 := \left\{ \left(1 - \sqrt{\frac{20 \log(qp)}{\rho qp}}\right) \rho \leq \hat{\rho} \leq \frac{1}{1 - \sqrt{\frac{20 \log(qp)}{\rho qp}}} \rho \right\}.$$

In (27) above, $\mathbf{B} \in \mathbb{R}^{q \times p}$ is a deterministic matrix. Let the singular value decomposition of \mathbf{B} be given as $\mathbf{B} = \sum_{i=1}^q \sigma_i(\mathbf{B}) x_i x_i^T$, where $\sigma_i(\mathbf{B})$ are the singular values of \mathbf{B} in decreasing order and x_i, y_i are the left and right singular vectors respectively. Recall the definition of $\varphi_{\lambda}^{\mathbf{B}}$ in (16). In particular, we choose $\lambda = \sigma_k(\mathbf{B})$, the k th singular value of \mathbf{B} in (27). As a result, in effect, we are bounding norm of projection of random vector $\mathbf{Y}_j - \rho \mathbf{M}_j$ for any given deterministic subspace of \mathbb{R}^p of dimension k .

Lemma G.3. For some positive constant $c_1 > 0$ and $C > 0$ large enough in definitions of E_1, \dots, E_5 ,

$$\begin{aligned} \mathbb{P}(E_1) &\geq 1 - 2e^{-c_1 pq\rho} - (1 - \rho)^{pq}, \\ \mathbb{P}(E_2) &\geq 1 - 2e^{-p}, \\ \mathbb{P}(E_3) &\geq 1 - 2e^{-p}, \\ \mathbb{P}(E_4) &\geq 1 - \frac{2}{(qp)^{10}}, \\ \mathbb{P}(E_5) &\geq 1 - \frac{2}{(qp)^{10}}. \end{aligned} \quad (28)$$

Proof. We bound the probability of events E_1, \dots, E_5 in that order.

Bounding E_1 . Let

$$\hat{\rho}_0 = \left(\sum_{i=1}^q \sum_{j=1}^p \mathbf{1}(Y_{ij} \text{ is obs.}) \right) / (qp).$$

That is, $\hat{\rho} = \max(\hat{\rho}_0, 1/(pq))$ and $\mathbb{E}[\hat{\rho}_0] = \rho$. We define the event $E_6 := \{\hat{\rho}_0 = \hat{\rho}\}$. Thus, we have that

$$\begin{aligned} \mathbb{P}(E_1^c) &= \mathbb{P}(E_1^c \cap E_6) + \mathbb{P}(E_1^c \cap E_6^c) \\ &= \mathbb{P}(|\hat{\rho}_0 - \rho| \geq \rho/20) + \mathbb{P}(E_1^c \cap E_6^c) \\ &\leq \mathbb{P}(|\hat{\rho}_0 - \rho| \geq \rho/20) + \mathbb{P}(E_6^c) \\ &= \mathbb{P}(|\hat{\rho}_0 - \rho| \geq \rho/20) + (1 - \rho)^{qp}, \end{aligned}$$

where the final equality follows by the independence of observations assumption and the fact that $\hat{\rho}_0 \neq \hat{\rho}$ only if we do not have any observations. By Bernstein's Inequality, we have that

$$\mathbb{P}(|\hat{\rho}_0 - \rho| \geq \rho/20) \leq 2e^{-c_1 \rho qp}.$$

Bounding E_2 . To start with, $\mathbb{E}[\mathbf{Y}] = \rho \mathbf{M}$. For any $i \in [q], j \in [p]$, the Y_{ij} are independent, 0 with probability $1 - \rho$ and with probability ρ equal to $M_{ij} + \epsilon_{ij}$ with $\|\epsilon_{ij}\|_{\psi_2} \leq \sigma$. Therefore, it follows that $\|Y_{ij} - \rho M_{ij}\|_{\psi_2} \leq C'\sigma$ for a constant $C' > 0$. Since $q \leq p$, using Theorem F.2 it follows that for an appropriately large constant $C > 0$,

$$\mathbb{P}(E_2) \geq 1 - 2e^{-p}.$$

Bounding E_3 . Recall that we assume $q \leq p$. Observe that for any matrix $A \in \mathbb{R}^{q \times p}$, $\|A\|_{\infty,2}, \|A\|_{2,\infty} \leq \|A\|_2$. Thus using the argument to bound E_2 , we have (28).

Bounding E_4 . Consider for $j \in [q]$,

$$\|\varphi_{\sigma_k(B)}^B(\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T)\|_2^2 = \sum_{i=1}^k \|y_i y_i^T (\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T)\|_2^2 \leq \sum_{i=1}^k \left(y_i^T (\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T) \right)_2^2 = \sum_{i=1}^k Z_i^2,$$

where $Z_i = y_i^T (\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T)$. By definition of the ψ_2 norm of a random variable and since y_i is unit norm vector that is deterministic (and hence independent the of random vector $\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T$), it follows that

$$\|Z_i\|_{\psi_2} = \|y_i^T (\mathbf{Y}_{j\cdot} - \rho \mathbf{M}_{j\cdot})\|_{\psi_2} \leq \|(\mathbf{Y}_{j\cdot} - \rho \mathbf{M}_{j\cdot})\|_{\psi_2}.$$

Since the coordinates of $\mathbf{Y}_{j\cdot}^T - \rho \mathbf{M}_{j\cdot}^T$ are mean-zero and independent, with ψ_2 norm bounded by $\sqrt{C}\sigma$ for some absolute constant $C > 0$, using arguments from [2, 3], it follows that

$$\mathbb{P}\left(\sum_{i=1}^k Z_i^2 > t\right) \leq 2k \exp\left(-\frac{t}{kC\sigma^2}\right).$$

Therefore, for choice of $t = C\sigma^2 k \log p$ with large enough constant $C > 0$, $q \leq p$, and taking a union bound over all $j \in [p]$, we have that

$$\mathbb{P}(E_4^c) \leq \frac{2}{(qp)^{10}}.$$

Bounding E_5 . Recall the definition of $\hat{\rho}$. By the binomial Chernoff bound, for $\varepsilon > 1$,

$$\begin{aligned} \mathbb{P}\left(\hat{\rho} > \varepsilon \rho\right) &\leq \exp\left(-\frac{(\varepsilon - 1)^2}{\varepsilon + 1} qp\rho\right), \quad \text{and} \\ \mathbb{P}\left(\hat{\rho} < \frac{1}{\varepsilon} \rho\right) &\leq \exp\left(-\frac{(\varepsilon - 1)^2}{2\varepsilon^2} qp\rho\right). \end{aligned}$$

By the union bound,

$$\mathbb{P}\left(\frac{1}{\varepsilon} \rho \leq \hat{\rho} \leq \varepsilon \rho\right) \geq 1 - \mathbb{P}\left(\hat{\rho} > \varepsilon \rho\right) - \mathbb{P}\left(\hat{\rho} < \frac{1}{\varepsilon} \rho\right).$$

Noticing $\varepsilon + 1 < 2\varepsilon < 2\varepsilon^2$ for all $\varepsilon > 1$, and substituting $\varepsilon = \left(1 - \sqrt{\frac{20 \log(qp)}{qp\rho}}\right)^{-1}$ completes the proof. \square

The following are immediate corollaries of the above stated bounds.

Corollary G.1. Let $E := E_1 \cap E_2$. Then, for $\rho \geq C \log(qp)/q$,

$$\mathbb{P}(E^c) \leq C_1 e^{-c_2 p},$$

where C_1 and c_2 are positive constants.

Corollary G.2. Let $E := E_2 \cap E_3 \cap E_4 \cap E_5$. Then,

$$\mathbb{P}(E^c) \leq \frac{C_1}{(qp)^{10}},$$

where C_1 is an absolute positive constant.

Probabilistic Bound for HSVT based Matrix Estimation. Recall $\epsilon = \|\mathbf{E}_k\|_\infty$. Then $\|\mathbf{E}_k\|_F^2 \leq \epsilon qp$. And $\|\mathbf{E}_k\|_2^2 \leq \|\mathbf{E}_k\|_F^2 \leq \epsilon qp$. Let $\rho \geq C \log(qp)/q$ for C large enough and recall $q \leq p$. Further, recall $\Gamma = \|\mathbf{M}_k\|_\infty$; thus, $\|\mathbf{M}\|_\infty \leq \Gamma + \epsilon$. Then $\|[\mathbf{M}_k]_j^T\|_2 \leq \Gamma\sqrt{p}$ and $\|[\mathbf{M}]_j^T\|_2 \leq (\Gamma + \epsilon)\sqrt{p}$.

Define $E = E_1 \cap E_2 \cap E_3 \cap E_4 \cap E_5$. Then, from Corollaries G.1 and G.2, we have that $\mathbb{P}(E^c) \leq \frac{C_1}{(qp)^{10}}$ for large enough constant $C_1 > 0$.

Under E_5 , we have $\varepsilon = \max(\hat{\rho}/\rho, \rho/\hat{\rho}) \leq \left(1 - \sqrt{\frac{20 \log(qp)}{qp\rho}}\right)^{-1}$. Under this choice of ε and using $\rho \geq C \log(qp)/q$, we have that for C large enough, $\varepsilon \leq C$ and $(\varepsilon - 1)^2 \leq C/p$.

Given this setup, under event E , Lemma G.2 leads to the following: for all $j \in [q]$ and with appropriately (re-defined) large enough constant $C > 0$,

$$\begin{aligned} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^2 &\leq C \frac{\sigma^2 p + \rho^2 \epsilon qp}{\rho^2 \sigma_k(\mathbf{M}_k)^2} \left(p\Gamma^2 + \frac{\sigma^2 p}{\rho^2} \right) \\ &\quad + \frac{C\sigma^2 k \log p}{\rho^2} + C(\Gamma + \epsilon)^2 + 2p\epsilon^2. \end{aligned}$$

That is, under event E ,

$$\max_{j \in [q]} \frac{1}{p} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^2 \leq C \frac{p(\sigma^2 + \rho^2 \epsilon q)}{\rho^2 \sigma_k(\mathbf{M}_k)^2} \left(\Gamma^2 + \frac{\sigma^2}{\rho^2} \right) + \frac{C\sigma^2 k \log p}{p\rho^2} + \frac{C(\Gamma + \epsilon)^2}{p} + 2\epsilon^2. \quad (29)$$

For any random variable X and event A , such that under event A , $X \leq B$ and $\mathbb{P}(A^c) \leq \delta$, we have

$$\begin{aligned} \mathbb{E}[X] &= \mathbb{E}[X\mathbf{1}(A)] + \mathbb{E}[X\mathbf{1}(A^c)] \\ &\leq \mathbb{E}[X\mathbf{1}(A)] + \mathbb{E}[X^2]^{\frac{1}{2}} \mathbb{P}(A^c)^{\frac{1}{2}} \\ &\leq B + \mathbb{E}[X^2]^{\frac{1}{2}} \delta^{\frac{1}{2}}. \end{aligned} \quad (30)$$

We shall use this reasoning above to bound $\mathbb{E}[\max_{j \in [q]} \frac{1}{p} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^2]$: let $X = \max_{j \in [q]} \frac{1}{p} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^2$ and $A = E$; B is given by right hand side of (29), $\delta = \frac{C_1}{(qp)^{10}}$; the only missing quantity that remains to be bounded is $\mathbb{E}[X^2]$. We do that next.

To begin with, for any $j \in [q]$,

$$\|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2 \leq \|\widehat{\mathbf{M}}_{j\cdot}^T\|_2 + \|\mathbf{M}_{j\cdot}^T\|_2 \quad (31)$$

by triangle inequality. As stated earlier, $\|[\mathbf{M}]_j^T\|_2 \leq (\Gamma + \epsilon)\sqrt{p}$. Next, we bound $\|\widehat{\mathbf{M}}_{j\cdot}^T\|_2$. From (15), the fact that $\hat{\rho} \geq 1/(qp)$, and Lemma G.1, we have

$$\begin{aligned} \|\widehat{\mathbf{M}}_{j\cdot}^T\|_2 &= \frac{1}{\hat{\rho}} \|\text{HSVT}_{\lambda_k}(\mathbf{Y})_{j\cdot}^T\|_2 \\ &\leq q p \|\phi_{\lambda_k}^{\mathbf{Y}}(\mathbf{Y}_{j\cdot}^T)\|_2 \\ &\leq q p \|\phi_{\lambda_k}^{\mathbf{Y}}\|_2 \|\mathbf{Y}_{j\cdot}^T\|_2 \\ &\leq q p \|\mathbf{Y}_{j\cdot}^T\|_2, \end{aligned} \quad (32)$$

where we used the fact that $\phi_{\lambda_k}^{\mathbf{Y}}$ is a projection operator and hence $\|\phi_{\lambda_k}^{\mathbf{Y}}\|_2 = 1$. Note that $Y_{ij} = B_{ij} \times (M_{ij} + \varepsilon_{ij})$, where B_{ij} is an independent Bernoulli variable with $\mathbb{P}(B_{ij} = 1) = \rho$ representing whether $(M_{ij} + \varepsilon_{ij})$ is observed or not. Therefore, $|Y_{ij}| = |B_{ij}| \times |M_{ij} + \varepsilon_{ij}| \leq (\Gamma + \epsilon) + |\varepsilon_{ij}|$. Therefore, from (31) and (32),

$$\begin{aligned} \max_{j \in [q]} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2 &\leq (\Gamma + \epsilon)\sqrt{p} + qp \left(\max_{j \in [q]} \|\mathbf{Y}_{j\cdot}^T\|_2 \right) \\ &\leq (\Gamma + \epsilon)\sqrt{p} + qp \times \sqrt{p} \left(\max_{i \in [p], j \in [q]} |Y_{ij}| \right) \\ &\leq 2qp^{\frac{3}{2}} \left(\Gamma + \epsilon + \max_{i \in [p], j \in [q]} |\varepsilon_{ij}| \right). \end{aligned} \quad (33)$$

Using $(a + b)^2 \leq 2a^2 + 2b^2$ twice, we have $(a + b)^4 \leq 8(a^4 + b^4)$. Therefore, from (34)

$$\max_{j \in [q]} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^4 \leq 16q^4 p^6 ((\Gamma + \epsilon)^4 + \max_{i \in [p], j \in [q]} |\epsilon_{ij}|^4). \quad (34)$$

Recall $\mathbb{E}[\epsilon_{ij}] = 0$, $\|\epsilon_{ij}\|_{\psi_2} \leq \sigma$ and ϵ_{ij} are independent across i, j . A property of ψ_2 -random variables is that $|\eta_{ij}|^\theta$ is a $\psi_{2/\theta}$ -random variable for $\theta \geq 1$. With choice of $\theta = 4$, we have

$$\mathbb{E}[\max_{ij} |\epsilon_{ij}|^4] \leq C' \sigma^4 \log^2(qp), \quad (35)$$

for some $C' > 0$ by Lemma F.1. From (32), (34), and (35), we have that

$$\left(\mathbb{E}[\max_{j \in [q]} \frac{1}{p^2} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^4] \right)^{\frac{1}{2}} \leq 4q^2 p^2 ((\Gamma + \epsilon)^4 + C' \sigma^4 \log^2(qp))^{\frac{1}{2}}. \quad (36)$$

Finally, using (29), (30) and (36), we conclude

$$\mathbb{E}[\max_{j \in [q]} \frac{1}{p} \|\widehat{\mathbf{M}}_{j\cdot}^T - \mathbf{M}_{j\cdot}^T\|_2^2] \leq \frac{p(C\sigma^2 + \rho^2 \epsilon q)}{\rho^2 \sigma_k(\mathbf{M}_k)^2} \left(\Gamma^2 + \frac{\sigma^2}{\rho^2} \right) + \frac{C\sigma^2 k \log p}{p\rho^2} + \frac{C(\Gamma + \epsilon)^2}{p} + 2\epsilon^2 + \frac{C}{(pq)^2}.$$

This completes the proof of Theorem G.1. \square

H Proof of Theorem 4.1

The proof of Theorem 4.1 will utilize Theorem G.1. To begin with, given N time series with observations over $[T]$, the mSSA algorithm as described in Section 1.1 constructs the $L \times (NT/L)$ stacked page matrix $\text{SP}((X_1, \dots, X_N), T, L)$ with $L = \sqrt{\min(N, T)}$, i.e. $L \leq T$.

As per the model described by (1) and Section 2, it follows that each entry of $\text{SP}((X_1, \dots, X_N), T, L)$ is an independent random variable; it is observed with probability $\rho \in (0, 1]$ independently and when it is observed, its equal to value of the latent time series plus zero-mean sub-Gaussian noise. In particular,

$$\mathbb{E}[\text{SP}((X_1, \dots, X_N), T, L)] = \rho \text{SP}((f_1, \dots, f_N), T, L),$$

where $\text{SP}((f_1, \dots, f_N), T, L) \in \mathbb{R}^{L \times (NT/L)}$ with entry in row $\ell \in [L]$ and column $(n-1) \times T/L + j$ equal to $f_n(\ell + (j-1) \times L)$. Further, when entry in row $\ell \in [L]$ and column $(n-1) \times T/L + j$ in $\text{SP}((X_1, \dots, X_N), T, L)$ is observed, i.e. $X_n(\ell + (j-1) \times L) \neq \star$, it is equal to $f_n(\ell + (j-1) \times L) + \eta_n(\ell + (j-1) \times L)$ where $\eta_n(\cdot)$ are independent, zero-mean sub-Gaussian variables with $\|\eta_n(\cdot)\|_{\psi_2} \leq \gamma$ as per the Property 2.3.

Under Properties 2.1 and 4.1, as a direct implication of Proposition D.1, $\text{SP}((f_1, \dots, f_N), T, L)$ has ϵ' -rank at most $R \times G$ with $\epsilon' = R\Gamma_1\epsilon$. That is, there exist rank $k \leq R \times G$ matrix $\mathbf{M}_k \in \mathbb{R}^{L \times (NT/L)}$ so that

$$\text{SP}((f_1, \dots, f_N), T, L) = \mathbf{M}_k + \mathbf{E}_k,$$

where $\|\mathbf{E}_k\|_\infty \leq \epsilon'$. Due to Property 2.1, it follows that $\|\mathbf{M}_k\|_\infty \leq R\Gamma_1\Gamma_2 + \epsilon'$. Under Property 4.2, we have $\sigma_k(\mathbf{M}_k) \geq c\sqrt{NT}/\sqrt{k}$ for some constant $c > 0$.

Define

$$\Gamma = R\Gamma_1\Gamma_2 + \epsilon' = R\Gamma_1(\Gamma_2 + \epsilon).$$

Recall from Section 1.1, the elements of the imputed multivariate time series are simply the entries of the matrix $\widehat{\text{SP}}((X_1, \dots, X_N), T, L)$ where $\widehat{\text{SP}}((X_1, \dots, X_N), T, L) = \frac{1}{\rho} \text{HSVT}_k(\text{SP}((X_1, \dots, X_N), T, L))$. That is, imputation in mSSA is carried out by applying HSVT to the stacked page matrix $\text{SP}((X_1, \dots, X_N), T, L)$.

All in all, the above description precisely meets the setup of Theorem G.1. To apply Theorem G.1, we require $\rho \geq C \log(NT)/\sqrt{NT}$ for $C > 0$ large enough. Note that the number of columns in $\widehat{\text{SP}}((X_1, \dots, X_N), T, L)$ is equal to NT/L for $L = \sqrt{\min(N, T)T}$ – for this choice of L , note that $NT/L \geq L$. Using $\sigma_k^2(\mathbf{M}_k) \geq cNT/k$, for some absolute constant $c \geq 0$, and using Theorem G.1, we obtain

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{(NT/L)} \|\widehat{\text{SP}}((X_1, \dots, X_N), T, L) - \text{SP}((f_1, \dots, f_N), T, L)\|_{2, \infty}^2 \right] \\ & \leq \frac{k(NT/L)(C\gamma^2 + \rho^2\epsilon' L)}{\rho^2 c^2 NT} \left(\Gamma^2 + \frac{\gamma^2}{\rho^2} \right) + \frac{C\gamma^2 k \log NT}{(NT/L)\rho^2} + \frac{C(\Gamma + \epsilon')^2}{(NT/L)} + 2(\epsilon')^2 + \frac{C}{(NT)^2} \end{aligned} \quad (37)$$

Recall that $k \leq R \times G$, $\epsilon' = R\Gamma_1\epsilon$, and $\Gamma = R\Gamma_1(\Gamma_2 + \epsilon)$. Hence, simplifying (37), we obtain that

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{(NT/L)} \|\widehat{\text{SP}}((X_1, \dots, X_N), T, L) - \text{SP}((f_1, \dots, f_N), T, L)\|_{2, \infty}^2 \right] \\ & \leq \tilde{C} \left(\frac{RG(1 + \rho^2 R\epsilon L)}{\rho^2 L} \left(R^2(1 + \epsilon^2) + \frac{1}{\rho^2} \right) + \frac{RG \log NT}{(NT/L)\rho^2} + \frac{(R(1 + \epsilon))^2}{(NT/L)} + (R\epsilon)^2 \right) \\ & \leq \tilde{C} \left(\frac{R^3 G \log NT}{\rho^4 L} + \frac{R^4 G(\epsilon + \epsilon^2 + \epsilon^3)}{\rho^2} \right), \end{aligned} \quad (38)$$

where $\tilde{C} = C(c, \Gamma_1, \Gamma_2, \gamma)$ is a positive constant dependent on model parameters including $\Gamma_1, \Gamma_2, \gamma$.

It can be easily verified that for any matrix, $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\frac{1}{mn} \|\mathbf{A}\|_F^2 \leq \frac{1}{n} \|\mathbf{A}\|_{\infty, 2}^2. \quad (39)$$

Further, there is a one-to-one mapping of $\hat{f}_n(\cdot)$ (resp. $f_n(\cdot)$) to the entries of $\widehat{\text{SP}}((X_1, \dots, X_N), T, L)$ (resp. $\text{SP}((f_1, \dots, f_N), T, L)$). Hence,

$$\text{ImpErr}(N, T) = \mathbb{E} \left[\frac{1}{NT} \|\widehat{\text{SP}}((X_1, \dots, X_N), T, L) - \text{SP}((f_1, \dots, f_N), T, L)\|_F^2 \right] \quad (40)$$

Therefore, from (38), (39), and (40) it follows that

$$\text{ImpErr}(N, T) \leq C(c, \Gamma_1, \Gamma_2, \gamma) \left(\frac{R^3 G \log NT}{\rho^4 L} + \frac{R^4 G(\epsilon + \epsilon^2 + \epsilon^3)}{\rho^2} \right)$$

This completes the proof of Theorem 4.1.

I Proof of Theorem 4.2

The forecasting algorithm, as described in Section 1.1, computes a linear model between the recent past and immediate future to forecast. We shall bound the forecasting error, $\text{ForErr}(N, T, L)$ as defined in (6). We start with some setup and notations, followed by a key proposition that establishes the existence of a linear model under the setup of Theorem 4.2, and then conclude with a detailed analysis of noisy, mis-specified least-squares.

Setup, Notations. For $L \geq 1, k \geq 1$, for ease of notations, we define

- $\text{SP}(X) = \text{SP}((X_1, \dots, X_N), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}(f) = \text{SP}((f_1, \dots, f_N), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}'(X) \in \mathbb{R}^{(L-1) \times (NT/L)}$ as the top $L - 1$ rows of $\text{SP}((X_1, \dots, X_N), T, L)$,
- $\text{SP}'(f) \in \mathbb{R}^{(L-1) \times (NT/L)}$ as the top $L - 1$ rows of $\text{SP}((f_1, \dots, f_N), T, L)$.

It is worth noting that $\mathbb{E}[\text{SP}(X)] = \rho \text{SP}(f)$ and hence

$$\text{SP}_{L\cdot}(X)^T = \rho \text{SP}_{L\cdot}(f)^T + \eta, \quad (41)$$

where $\eta \in \mathbb{R}^{(NT)/L}$ is a random vector with each component being independent, zero-mean with its distribution given as: it is 0 with probability $1 - \rho$ and with probability ρ , due to Property 2.3, it equals a zero-mean sub-Gaussian random variable with $\|\cdot\|_{\psi_2} \leq \gamma$. Therefore, using arguments in [2, 3], each component of η is an independent, zero-mean random variable with $\|\cdot\|_{\psi_2}$ bounded above by $C'(\gamma^2 + R\Gamma_1\Gamma_2)$ for some absolute constant $C' > 0$. Let $K = C'(\gamma^2 + R\Gamma_1\Gamma_2)$ and hence each component of η has $\|\cdot\|_{\psi_2}$ bounded by K .

Now, recall that for forecasting, we first apply the imputation algorithm (i.e. HSVT) to $\text{SP}((X_1, \dots, X_N), T, L)$ by replacing \star s, i.e. missing observations by 0 as well as setting all the entries in the last row equal to 0. Equivalently, the imputation algorithm is applied to $\text{SP}'(X)$ after setting all missing values to 0. Let $\widehat{\text{SP}}' \in \mathbb{R}^{L-1 \times (NT/L)}$ be the estimate produced from the imputation algorithm applied to $\text{SP}'(X)$. Under the setup of Theorem 4.1, by following arguments identical to that of Theorems G.1 and 4.1—in particular, refer to (38)—it follows that by selecting the right choice of $k \leq R \times G$, we have

$$\mathbb{E}\left[\frac{1}{(NT/L)} \|\widehat{\text{SP}}' - \text{SP}'(f)\|_{2,\infty}^2\right] \leq \tilde{C} \left(\frac{R^3 G \log NT}{\rho^4 L} + \frac{R^4 G(\epsilon + \epsilon^2 + \epsilon^3)}{\rho^2} \right), \quad (42)$$

where $\tilde{C} = C(c, \Gamma_1, \Gamma_2, \gamma) > 0$ is a constant dependent on $c, \Gamma_1, \Gamma_2, \gamma$.

Now, the mSSA forecasting algorithm finds $\hat{\beta} = \hat{\beta}((X_1, \dots, X_N), TL; k)$, by solving the following Ordinary Least Squares (OLS):

$$\hat{\beta} \in \text{minimize} \quad \left\| \frac{1}{\rho} \text{SP}(X)_{L\cdot} - \widehat{\text{SP}}'^T \beta \right\|_2^2 \quad \text{over} \quad \beta \in \mathbb{R}^{L-1}. \quad (43)$$

And subsequently, $\widehat{\text{SP}}'^T \hat{\beta}$ is used as the estimate for $\text{SP}(f)_{L\cdot} \in \mathbb{R}^{NT/L}$, the L th row of the latent $\text{SP}(f)$. The goal is to bound the forecasting error $\text{ForErr}(N, T, L)$, which is given by

$$\text{ForErr}(N, T, L) = \mathbb{E}\left[\frac{1}{(NT/L)} \|\text{SP}(f)_{L\cdot} - \widehat{\text{SP}}'^T \hat{\beta}\|_2^2\right].$$

Therefore, our interest is in bounding $\mathbb{E}[\|\text{SP}_{L\cdot}(f) - \widehat{\text{SP}}'^T \hat{\beta}\|_2^2]$.

Now, we recall from Proposition 4.2 that there exists $\beta^* \in \mathbb{R}^{L-1}$, such that

$$\|\text{SP}(f)_{L\cdot}^T - \text{SP}'(f)^T \beta^*\|_\infty \leq C_2 \epsilon,$$

where $C_2 := R\Gamma_1(1 + \|\beta^*\|_1)$.

Bounding $\mathbb{E}[\|\text{SP}_{L\cdot}(f) - \widehat{\text{SP}}'^T \hat{\beta}\|_2^2]$. By (43) and (41)

$$\begin{aligned} \left\| \frac{1}{\rho} \text{SP}(X)_{L\cdot} - \widehat{\text{SP}}'^T \hat{\beta} \right\|_2^2 &\leq \left\| \frac{1}{\rho} \text{SP}(X)_{L\cdot} - \widehat{\text{SP}}'^T \beta^* \right\|_2^2 \\ &= \left\| \frac{\rho}{\rho} \text{SP}(f)_{L\cdot} + \eta - \widehat{\text{SP}}'^T \beta^* \right\|_2^2 \\ &= \left\| \frac{\rho}{\rho} \text{SP}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^* \right\|_2^2 + \|\eta\|_2^2 + 2\eta^T \left(\frac{\rho}{\rho} \text{SP}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^* \right) \end{aligned} \quad (44)$$

Also,

$$\begin{aligned} \left\| \frac{1}{\rho} \text{SP}(X)_{L\cdot} - \widehat{\text{SP}}'^T \hat{\beta} \right\|_2^2 &= \left\| \frac{\rho}{\rho} \text{SP}(f)_{L\cdot} + \eta - \widehat{\text{SP}}'^T \hat{\beta} \right\|_2^2 \\ &= \left\| \frac{\rho}{\rho} \text{SP}(f)_{L\cdot} - \widehat{\text{SP}}'^T \hat{\beta} \right\|_2^2 + \|\eta\|_2^2 + 2\eta^T \left(\frac{\rho}{\rho} \text{SP}(f)_{L\cdot} - \widehat{\text{SP}}'^T \hat{\beta} \right). \end{aligned} \quad (45)$$

From (44) and (45)

$$\begin{aligned} & \mathbb{E}[\|\frac{\rho}{\widehat{\rho}}\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \widehat{\beta}\|_2^2] \\ & \leq \mathbb{E}[\|\frac{\rho}{\widehat{\rho}}\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^*\|_2^2] + 2\mathbb{E}[\eta^T \widehat{\text{SP}}'^T (\beta^* - \widehat{\beta})] \end{aligned} \quad (46)$$

η is independent of $\widehat{\text{SP}}'$, β^* , and $\widehat{\rho}$; $\mathbb{E}[\eta] = \mathbf{0}$; thus, we have that

$$\mathbb{E}[\eta^T \widehat{\text{SP}}'^T \beta^*] = 0. \quad (47)$$

By (43), we have $\widehat{\beta} = \widehat{\text{SP}}'^{T,\dagger} \frac{1}{\widehat{\rho}} \widehat{\text{SP}}(X)_{L\cdot}$, where $\widehat{\text{SP}}'^{T,\dagger}$ is pseudo-inverse of $\widehat{\text{SP}}'^T$. That is,

$$\widehat{\beta} = \widehat{\text{SP}}'^{T,\dagger} \frac{\rho}{\widehat{\rho}} \widehat{\text{SP}}(f)_{L\cdot} + \frac{1}{\widehat{\rho}} \widehat{\text{SP}}'^{T,\dagger} \eta. \quad (48)$$

Using cyclic and linearity of Trace operator; the independence properties of η ; and (48); we have

$$\begin{aligned} \mathbb{E}[\eta^T \widehat{\text{SP}}'^T \widehat{\beta}] &= \mathbb{E}[\eta^T \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger} \frac{\rho}{\widehat{\rho}} \widehat{\text{SP}}(f)_{L\cdot}] + \mathbb{E}[\frac{1}{\widehat{\rho}} \eta^T \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger} \eta] \\ &= \mathbb{E}[\eta]^T \mathbb{E}[\widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger} \frac{\rho}{\widehat{\rho}} \widehat{\text{SP}}(f)_{L\cdot}] + \mathbb{E}[\frac{1}{\widehat{\rho}} \text{Tr}(\eta^T \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger} \eta)] \\ &= \mathbb{E}[\frac{1}{\widehat{\rho}} \text{Tr}(\widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger} \eta \eta^T)] \\ &= \text{Tr}(\mathbb{E}[\frac{1}{\widehat{\rho}} \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger}] \mathbb{E}[\eta \eta^T]) \\ &\leq C(\gamma)k/\rho, \end{aligned} \quad (49)$$

where $C(\gamma)$ is a function only of γ . To see the last inequality, we use various facts. First, by the definition of the HSVT algorithm $\widehat{\text{SP}}'^T$ has rank at most k . Second, let $\widehat{\text{SP}}'^T = \mathbf{U} \mathbf{S} \mathbf{V}^T$ be the singular value decomposition of $\widehat{\text{SP}}'^T$, we have

$$\begin{aligned} \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger} &= \mathbf{U} \mathbf{S} \mathbf{V}^T \mathbf{V} \mathbf{S}^\dagger \mathbf{U}^T \\ &= \mathbf{U} \tilde{\mathbf{U}} \mathbf{U}^T, \end{aligned}$$

That is, $\frac{1}{\widehat{\rho}} \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger}$ is a positive semi-definite matrix and $\text{Tr}(\frac{1}{\widehat{\rho}} \widehat{\text{SP}}'^T \widehat{\text{SP}}'^{T,\dagger}) \leq k/\widehat{\rho}$. The matrix $\mathbb{E}[\eta \eta^T]$ is diagonal with all the non-zero entries on diagonal (variance of components of η) bounded above by a constant that depends on γ . For a positive semi-definite matrix A and positive semi-definite diagonal matrix B , $\text{Tr}(AB) \leq \|B\|_2 \text{Tr}(A)$. For $\rho \geq C \log(NT)/\sqrt{NT}$ for large enough C , one can verify that $\mathbb{E}[1/\widehat{\rho}] \leq 2/\rho$. This completes the justification of the last step of (49).

Now consider the term $\|\frac{\rho}{\widehat{\rho}}\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^*\|_2^2$. Note,

$$\begin{aligned} \|\frac{\rho}{\widehat{\rho}}\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^*\|_2^2 &= \|(\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^*) + (\frac{\rho - \widehat{\rho}}{\widehat{\rho}}) \widehat{\text{SP}}(f)_{L\cdot}\|_2^2 \\ &\leq 2\|(\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^*)\|_2^2 + 2\|\frac{\rho - \widehat{\rho}}{\widehat{\rho}} \widehat{\text{SP}}(f)_{L\cdot}\|_2^2. \end{aligned} \quad (50)$$

We will bound the two terms on the r.h.s of (50) separately. We now consider the first term.

$$\|\widehat{\text{SP}}(f)_{L\cdot} - \widehat{\text{SP}}'^T \beta^*\|_2^2 \leq 2\|\widehat{\text{SP}}(f)_{L\cdot} - \text{SP}'(f)^T \beta^*\|_2^2 + 2\|\text{SP}'(f)^T \beta^* - \widehat{\text{SP}}'^T \beta^*\|_2^2. \quad (51)$$

By Proposition 4.2

$$\|\widehat{\text{SP}}(f)_{L\cdot} - \text{SP}'(f)^T \beta^*\|_2 \leq \|\widehat{\text{SP}}(f)_{L\cdot} - \text{SP}'(f)^T \beta^*\|_\infty \sqrt{NT/L} \leq C_2 \epsilon \sqrt{NT/L}, \quad (52)$$

where we used the fact that for any $v \in \mathbb{R}^p$, $\|v\|_2 \leq \|v\|_\infty \sqrt{p}$. And,

$$\|\text{SP}'(f)^T \beta^* - \widehat{\text{SP}'}^T \beta^*\|_2 = \|(\text{SP}'(f) - \widehat{\text{SP}'})^T \beta^*\|_2 \leq \|\text{SP}'(f) - \widehat{\text{SP}'}\|_{2,\infty} \|\beta^*\|_1, \quad (53)$$

where we used the fact that for any $A \in \mathbb{R}^{q \times p}$, $v \in \mathbb{R}^p$, $\|Av\|_2 \leq \|A\|_{2,\infty} \|v\|_1$. Finally, note that

$$\|\text{SP}(f)_{L\cdot} - \widehat{\text{SP}'}^T \widehat{\beta}\|_2^2 \leq 2\|\frac{\rho}{\widehat{\rho}} \text{SP}(f)_{L\cdot} - \widehat{\text{SP}'}^T \widehat{\beta}\|_2^2 + 2\|\frac{\rho - \widehat{\rho}}{\widehat{\rho}} \text{SP}(f)_{L\cdot}\|_2^2. \quad (54)$$

Using (46), (47), (49), (50), (51), (52), (53), and the bound in (54), we obtain

$$\begin{aligned} & \mathbb{E}[\|\text{SP}(f)_{L\cdot} - \widehat{\text{SP}'}^T \widehat{\beta}\|_2^2] \\ & \leq 4C(\gamma)k/\rho + 6\mathbb{E}[\|\frac{\rho - \widehat{\rho}}{\widehat{\rho}} \text{SP}(f)_{L\cdot}\|_2^2] + 2C_2\epsilon^2(NT/L) + 2\|\beta^*\|_1^2 \|\text{SP}'(f) - \widehat{\text{SP}'}\|_{2,\infty}^2. \end{aligned} \quad (55)$$

Note that $\|\text{SP}(f)\|_\infty \leq R\Gamma_1\Gamma_2$. Hence, $\|\text{SP}(f)_{L\cdot}\|_2^2 \leq C(\Gamma_1, \Gamma_2)R^2(NT/L)$, for large enough constant $C(\Gamma_1, \Gamma_2)$ that may depend on Γ_1, Γ_2 . Using the bounds derived in Lemma G.3, one can verify that $\mathbb{E}[(\frac{\rho - \widehat{\rho}}{\widehat{\rho}})^2] \leq C/(NT/L)$ for large enough positive constant C . Therefore, we have that

$$6\mathbb{E}[\|\frac{\rho - \widehat{\rho}}{\widehat{\rho}} \text{SP}(f)_{L\cdot}\|_2^2] \leq C(\Gamma_1, \Gamma_2)R^2 \quad (56)$$

Using (42), (56), and the bound in (55); diving by $1/(NT/L)$ on both sides; and noting $k \leq R \times G$, we obtain

$$\begin{aligned} & \mathbb{E}[\frac{1}{(NT/L)} \|\text{SP}(f)_{L\cdot} - \widehat{\text{SP}'}^T \widehat{\beta}\|_2^2] \\ & \leq C(c, \gamma, \Gamma_1, \Gamma_2) \left(\frac{RG}{\rho(NT/L)} + \frac{R^2}{(NT/L)} + R(1 + \|\beta^*\|_1)\epsilon^2 + \|\beta^*\|_1^2 \left(\frac{R^3 G \log NT}{\rho^4 L} + \frac{R^4 G(\epsilon + \epsilon^2 + \epsilon^3)}{\rho^2} \right) \right) \\ & \leq C(c, \gamma, \Gamma_1, \Gamma_2) \left(\max(1, \|\beta^*\|_1, \|\beta^*\|_1^2) \left(\frac{R^3 G \log NT}{\rho^4 L} + \frac{R^4 G(\epsilon + \epsilon^2 + \epsilon^3)}{\rho^2} \right) \right) \end{aligned} \quad (57)$$

Letting $L = \sqrt{\min(N, T)T}$, using (57), and noting that

$$\text{ForErr}(N, T, L) = \mathbb{E}[\frac{1}{(NT/L)} \|\text{SP}(f)_{L\cdot} - \widehat{\text{SP}'}^T \widehat{\beta}\|_2^2]$$

completes the proof of Theorem 4.2.

I.1 Proof of Proposition 4.2

For this proof, we utilize a modified version of the stacked Hankel matrix defined in Appendix D. Define the modified Hankel matrix for time series f_n , for $n \in [N]$, as $\widetilde{\text{H}}(n) \in \mathbb{R}^{T \times 2T}$, where for $i \in [T]$, $j \in [2T]$, we have

$$\widetilde{\text{H}}(n)_{ij} = f_n(i + j - 1 - T).$$

Define $\widetilde{\text{SH}} \in \mathbb{R}^{T \times NT}$ as the column wise concatenation of the matrices $\widetilde{\text{H}}(n)$ for $n \in [N]$, i.e., $\widetilde{\text{SH}} := [\widetilde{\text{H}}(1), \dots, \widetilde{\text{H}}(N)]$. By a straightforward modification of the proof of Proposition D.1, we have $\widetilde{\text{SH}}$ has ϵ' -rank bounded by $R \times G$ with $\epsilon' = R\Gamma_1\epsilon$. That is, there exists a matrix $\text{M} \in \mathbb{R}^{T \times NT}$ such that,

$$\text{rank}(\text{M}) \leq RG, \quad \|\widetilde{\text{SH}} - \text{M}\|_\infty \leq \epsilon'$$

Since $\text{rank}(\text{M}) \leq RG$, it must be the case that within the last RG rows of M , there exists at least one row, which we denote as r^* , that can be written as a linear combination of at most RG rows above

it, which we denote as r_1, \dots, r_{RG} . Specifically there exists a vector $\theta := (\theta_1, \dots, \theta_{RG}) \in \mathbb{R}^{RG}$ such that

$$\mathbf{M}_{r^*, \cdot} = \sum_{\ell=1}^{RG} \theta_{\ell} \mathbf{M}_{r_{\ell}, \cdot}.$$

Hence for $j \in [2T]$,

$$\begin{aligned} & \left| \widetilde{\mathbf{S}}\mathbf{H}_{r^*, j} - \sum_{\ell=1}^{RG} \theta_{\ell} \widetilde{\mathbf{S}}\mathbf{H}_{r_{\ell}, j} \right| \\ &= \left| \widetilde{\mathbf{S}}\mathbf{H}_{r^*, j} \pm \mathbf{M}_{r^*, j} - \sum_{\ell=1}^{RG} \theta_{\ell} \widetilde{\mathbf{S}}\mathbf{H}_{r_{\ell}, j} \pm \sum_{\ell=1}^{RG} \theta_{\ell} \mathbf{M}_{r_{\ell}, t} \right| \\ &\leq \left| \widetilde{\mathbf{S}}\mathbf{H}_{r^*, j} - \mathbf{M}_{r^*, j} \right| + \left| \sum_{\ell=1}^{RG} \theta_{\ell} \widetilde{\mathbf{S}}\mathbf{H}_{r_{\ell}, j} - \sum_{\ell=1}^{RG} \theta_{\ell} \mathbf{M}_{r_{\ell}, t} \right| + \left| \mathbf{M}_{r^*, j} - \sum_{\ell=1}^{RG} \theta_{\ell} \mathbf{M}_{r_{\ell}, t} \right| \\ &= \left| \widetilde{\mathbf{S}}\mathbf{H}_{r^*, j} - \mathbf{M}_{r^*, j} \right| + \left| \sum_{\ell=1}^{RG} \theta_{\ell} (\widetilde{\mathbf{S}}\mathbf{H}_{r_{\ell}, j} - \mathbf{M}_{r_{\ell}, t}) \right| \\ &\leq \epsilon' + \|\theta\|_1 \|\widetilde{\mathbf{S}}\mathbf{H}_{r_{\ell}, j} - \mathbf{M}_{r_{\ell}, t}\|_{\infty} \\ &\leq R\Gamma_1(1 + \|\theta\|_1)\epsilon. \end{aligned} \tag{58}$$

Observe that every entry of $\text{SP}(f)_L$ appears within $\widetilde{\mathbf{S}}\mathbf{H}_{r^*, \cdot}$; this can be seen by noting that $\widetilde{\mathbf{S}}\mathbf{H}$ is skew-symmetric and thus every entry in the last row of $\widetilde{\mathbf{S}}\mathbf{H}$ appears along the appropriate diagonal. Using this skew-symmetric property of $\widetilde{\mathbf{S}}\mathbf{H}$ and (58), it implies that by appropriately selecting entries in $\widetilde{\mathbf{S}}\mathbf{H}$, there exists $\beta^* \in \mathbb{R}^{L-1}$,

$$\|\text{SP}(f)_L^T - \text{SP}'(f)^T \beta^*\|_{\infty} \leq R\Gamma_1(1 + \|\beta\|_1)\epsilon,$$

where the non-zero entries in β^* correspond to the entries of θ . Noting that $\theta \in \mathbb{R}^{RG}$ implies $\|\beta^*\|_0 \leq RG$. This completes the proof.

J Proof of Theorem 6.1

Setup, Notations. For $L \geq 1, k \geq 1$, for ease of notations, we define

- $\text{SP}(X) = \text{SP}((X_1, \dots, X_N), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}(X^2) = \text{SP}((X_1^2, \dots, X_N^2), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}(f) = \text{SP}((f_1, \dots, f_N), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}(f^2) = \text{SP}((f_1^2, \dots, f_N^2), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}(\sigma^2) = \text{SP}((\sigma_1^2, \dots, \sigma_N^2), T, L) \in \mathbb{R}^{L \times (NT/L)}$,
- $\text{SP}(f^2 + \sigma^2) = \text{SP}(f^2) + \text{SP}(\sigma^2)$.

Recalling that $\rho = 1$, we note that

$$\mathbb{E}[\text{SP}(X)] = \text{SP}(f), \quad \mathbb{E}[\text{SP}(X^2)] = \text{SP}(f^2 + \sigma^2).$$

Further, from the definition of the variance estimation algorithm, we recall

$$\begin{aligned} \widehat{\text{SP}}(f) &:= \widehat{\text{SP}}((X_1, \dots, X_N), T, L) = \frac{1}{\widehat{\rho}} \text{HSVT}_k(\text{SP}((X_1, \dots, X_N), T, L)) \\ \widehat{\text{SP}}(f^2 + \sigma^2) &:= \widehat{\text{SP}}((X_1^2, \dots, X_N^2), T, L) = \frac{1}{\widehat{\rho}} \text{HSVT}_k(\text{SP}((X_1^2, \dots, X_N^2), T, L)) \end{aligned}$$

We denote

- $\widehat{\text{SP}}(f^2) = \widehat{\text{SP}}(f) \circ \widehat{\text{SP}}(f)$
- $\widehat{\text{SP}}(\sigma^2) = \max\left(\widehat{\text{SP}}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2), \mathbf{0}\right),$

where $\mathbf{0} \in \mathbb{R}^{L \times (NT/L)}$ is a matrix of all zeroes, and we apply the $\max(\cdot)$ above entry-wise. We remind the reader the output of the variance estimation algorithm is $\widehat{\text{SP}}(\sigma^2)$. Thus, we have

$$\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T (\sigma_n(t)^2 - \hat{\sigma}_n^2(t))^2 = \frac{1}{NT} \|\text{SP}(\sigma^2) - \widehat{\text{SP}}(\sigma^2)\|_F^2.$$

Initial Decomposition. Note that since $\sigma_n^2(t) \geq 0$ for $n \in [N]$ and $t \in [T]$, we have that

$$\begin{aligned} & \frac{1}{NT} \|\text{SP}(\sigma^2) - \widehat{\text{SP}}(\sigma^2)\|_F^2 \\ & \leq \frac{1}{NT} \|\text{SP}(\sigma^2) - (\widehat{\text{SP}}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2))\|_F^2 \\ & = \frac{1}{NT} \|\text{SP}(f^2 + \sigma^2) - \text{SP}(f^2) - (\widehat{\text{SP}}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2))\|_F^2 \\ & \leq \frac{2}{NT} \|\text{SP}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2 + \sigma^2)\|_F^2 + \frac{2}{NT} \|\text{SP}(f^2) - \widehat{\text{SP}}(f^2)\|_F^2 \end{aligned} \quad (59)$$

We bound the two terms on the r.h.s of (59) separately.

Bounding $\mathbb{E}[\|\text{SP}(f^2) - \widehat{\text{SP}}(f^2)\|_F^2]$.

$$\begin{aligned} \|\text{SP}(f^2) - \widehat{\text{SP}}(f^2)\|_F^2 &= \sum_{n=1}^N \sum_{t=1}^T \left(f_n^2(t) - \hat{f}_n^2(t)\right)^2 \\ &= \sum_{n=1}^N \sum_{t=1}^T \left(f_n(t) - \hat{f}_n(t)\right)^2 \left(f_n(t) + \hat{f}_n(t)\right)^2 \\ &\leq \left[\max_{n \in [N], t \in [T]} \left(f_n(t) + \hat{f}_n(t)\right)^2 \right] \left[\sum_{n=1}^N \sum_{t=1}^T \left(f_n(t) - \hat{f}_n(t)\right)^2 \right] \\ &\stackrel{(a)}{\leq} C(\Gamma_1, \Gamma_2, \Gamma_3) R^2 \left[\sum_{n=1}^N \sum_{t=1}^T \left(f_n(t) - \hat{f}_n(t)\right)^2 \right] \\ &= C(\Gamma_1, \Gamma_2, \Gamma_3) R^2 \|\text{SP}(f) - \widehat{\text{SP}}(f)\|_F^2 \end{aligned} \quad (60)$$

Bounding $\|\text{SP}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2 + \sigma^2)\|_F^2$. To bound $\|\text{SP}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2 + \sigma^2)\|_F^2$, we modify the proof of Theorem 4.1 in a straightforward manner. The need for the modification is that Theorem 4.1 was proven for the case where the coordinate wise noise, $\eta_n(t) = X_n(t) - f_n(t)$ are independent sub-gaussian random variables, and $\|\eta\|_{\psi_2} \leq \gamma$. However, one can verify that $X_n^2(t) - f_n^2(t) - \sigma_n^2(t)$ is a sub-exponential random variable with $\|\cdot\|_{\psi_1}$ norm bounded as

$$\begin{aligned} \|X_n^2(t) - f_n^2(t) - \sigma_n^2(t)\|_{\psi_1} &\leq \|X_n^2(t)\|_{\psi_1} \\ &= \|f_n^2(t) + 2f_n(t)\eta_n(t) + \eta_n^2(t)\|_{\psi_1} \\ &\leq 2\|f_n^2(t)\|_{\psi_1} + 2\|\eta_n^2(t)\|_{\psi_1} \\ &= 2\|f_n(t)\|_{\psi_2}^2 + 2\|\eta_n(t)\|_{\psi_2}^2 \\ &\leq C(\Gamma_1, \Gamma_2) R^2 + 2\gamma^2 \\ &\leq C(\Gamma_1, \Gamma_2, \gamma) R^2, \end{aligned}$$

where we have use the standard facts that for a random variable A , $\|A - \mathbb{E}[A]\|_{\psi_1} \leq \|A\|_{\psi_1}$ and $\|A^2\|_{\psi_1} = \|A\|_{\psi_2}^2$.

Further, note that by using Properties 2.1, 2.2, 6.1, and 6.2, and a straightforward modification of Proposition D.1, we have

$$\begin{aligned}\text{rank}(\text{SP}(f^2 + \sigma^2)) &\leq \text{rank}(\text{SP}(f^2)) + \text{rank}(\text{SP}(\sigma^2)) \\ &\leq (RG)^2 + (R'G'),\end{aligned}$$

where we have used that for any two matrices \mathbf{A}, \mathbf{B} , we have $\text{rank}(\mathbf{A} \circ \mathbf{A}) \leq \text{rank}(\mathbf{A})^2$, where \circ denotes Hadamard product, and $\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B})$. We define $\tilde{k} := (RG)^2 + (R'G')$.

Modified Theorem 4.1. Below, we state the modified version of Theorem 4.1 to get our desired result.

Lemma J.1 (Imputation Error). *Let the conditions of Theorem 6.1 hold. Then,*

$$\begin{aligned}\mathbb{E}\left[\max_{j \in [L]} \frac{1}{(NT/L)} \|\text{SP}(f^2 + \sigma^2)_{L,\cdot}^T - \widehat{\text{SP}}(f^2 + \sigma^2)_{L,\cdot}^T\|_2^2\right] \\ \leq C(\Gamma_1, \Gamma_2, \Gamma'_1, \Gamma'_2, \gamma, R, R') \left(\frac{(G^2 + G') \log^2 NT}{L} \right),\end{aligned}$$

where $C(\Gamma_1, \Gamma_2, \Gamma'_1, \Gamma'_2, \gamma, R, R')$ is a term that depends only polynomially on $\Gamma_1, \Gamma_2, \Gamma'_1, \Gamma'_2, \gamma, R, R'$.

Proof. To reduce redundancy, we provide an overview of the argument needed for this proof, focusing only the parts of the arguments made in Theorem 4.1 that need to be modified. For ease of exposition, we let $\tilde{C} = C(\Gamma_1, \Gamma_2, \Gamma'_1, \Gamma'_2, \gamma, R, R')$. We begin by matching notation with that used in Theorem 4.1; in particular with respect to $\rho, k, \epsilon, \Gamma$. Under the setup of Theorem 6.1, we have $\rho = 1, k = \tilde{k}, \epsilon = 0, \Gamma \leq \tilde{C}$. Further, recall the definition of $\mathbf{Y}, \mathbf{M}, p, q, \sigma$ from Appendix G.1. We will now use $\mathbf{Y} = \text{SP}(X^2)$, and $\mathbf{M} = \text{SP}(f^2 + \sigma^2)$, $\sigma = \gamma, p = (NT/L), q = L$. One can verify that there is only required change to the proof of Theorem 4.1; in particular, in the argument made to prove Theorem G.1, we need to re-define events E_2, E_3, E_4 in (25), (26), (27) for the case where $(\mathbf{Y} - \mathbf{M})_{ij}$ is mean-zero sub-exponential. Using the result from [2, 3], which bounds the operator norm of a matrix with sub-exponential mean-zero entries, we have with probability at least $1 - 1/((NT)^{10})$

$$\|\mathbf{Y} - \mathbf{M}\|_2 \leq \tilde{C} \sqrt{(NT/L)} \log^2 NT \quad (61)$$

As a result (61), and standard concentration inequalities for sub-exponential random variables, we have the modified events, $\tilde{E}_2, \tilde{E}_3, \tilde{E}_4$.

$$\begin{aligned}\tilde{E}_2 &:= \left\{ \|\mathbf{Y} - \rho \mathbf{M}\|_2 \leq \tilde{C} \sqrt{(NT/L)} \log^2 NT \right\}, \\ \tilde{E}_3 &:= \left\{ \|\mathbf{Y} - \rho \mathbf{M}\|_{\infty, 2}, \|\mathbf{Y} - \rho \mathbf{M}\|_{2, \infty} \leq \tilde{C} \sqrt{(NT/L)} \log^2 NT \right\}, \\ \tilde{E}_4 &:= \left\{ \max_{j \in [q]} \|\varphi_{\sigma_k(\mathbf{B})}^{\mathbf{B}} (\mathbf{Y}_{j,\cdot}^T - \rho \mathbf{M}_{j,\cdot}^T)\|_2^2 \leq \tilde{C} \tilde{k} \log^2 (NT/L) \right\},\end{aligned}$$

Using these modified events in the proofs of Theorem G.1 and Theorem 4.1, and appropriately simplifying leads to the desired result. \square

By Lemma J.1 and (39), we have that

$$\begin{aligned}\frac{1}{NT} \mathbb{E}[\|\text{SP}(f^2 + \sigma^2) - \widehat{\text{SP}}(f^2 + \sigma^2)\|_F^2] &\leq \mathbb{E}\left[\max_{j \in [L]} \frac{1}{(NT/L)} \|\text{SP}(f^2 + \sigma^2)_{L,\cdot}^T - \widehat{\text{SP}}(f^2 + \sigma^2)_{L,\cdot}^T\|_2^2\right] \\ &\leq C(\Gamma_1, \Gamma_2, \Gamma'_1, \Gamma'_2, \gamma, R, R') \left(\frac{(G^2 + G') \log^2 NT}{L} \right).\end{aligned} \quad (62)$$

Completing proof. Substituting (60) and (62) into (59) and letting $L = \sqrt{\min(N, T)T}$

$$\frac{1}{NT} \|\text{SP}(\sigma^2) - \widehat{\text{SP}}(\sigma^2)\|_F^2 \leq C(\Gamma_1, \Gamma_2, \Gamma_3, \Gamma'_1, \Gamma'_2, \gamma, R, R') \left(\frac{(G^2 + G') \log^2 NT}{\sqrt{\min(N, T)T}} \right).$$

This completes the proof.

K tSSA Proofs

K.1 Proof of Proposition 6.1

Consider $n \in [N]$, $\ell \in [L]$, $s \in [T/L]$. By Property 2.1,

$$\begin{aligned}\mathbf{T}_{n\ell s} &= f_n((s-1) \times L + \ell) \\ &= \sum_{r=1}^R U_{nr} W_{r((s-1) \times L + \ell)}.\end{aligned}\tag{63}$$

The Hankel matrix induced by time series W_r has rank at most G as per Property 2.2. The Page matrix associated with it is of dimension $L \times T/L$ with entry in its ℓ -th row and s -th column equal to $W_{r((s-1) \times L + \ell)}$. Since this Page matrix can be viewed as a sub-matrix of the Hankel matrix, it has rank at most G as well. That is, there exists vectors $w_{\ell}^r, v_s^r \in \mathbb{R}^G$ such that

$$W_{r((s-1) \times L + \ell)} = \sum_{g=1}^G w_{\ell g}^r v_{sg}^r.\tag{64}$$

From (63) and (64), it follows that

$$\begin{aligned}\mathbf{T}_{n\ell s} &= \sum_{r=1}^R U_{nr} \left(\sum_{g=1}^G w_{\ell g}^r v_{sg}^r \right) \\ &= \sum_{r \in [R], g \in [G]} U_{nr} w_{\ell g}^r v_{sg}^r \\ &= \sum_{r \in [R], g \in [G]} a_{n(r,g)} b_{\ell(r,g)} c_{s(r,g)},\end{aligned}\tag{65}$$

where $a_{n(r,g)} = U_{nr}$, $b_{\ell(r,g)} = w_{\ell g}^r$ and $c_{s(r,g)} = v_{sg}^r$. Thus (65) implies that \mathbf{T} has CP-rank at most $R \times G$.

By the setup and model definition, it follows $\mathbb{T}_{n\ell s} = X_n((s-1) \times L + \ell)$. And $X_n((s-1) \times L + \ell) = \star$ with probability $1 - \rho$ and $f_n((s-1) \times L + \ell) + \eta_n((s-1) \times L + \ell)$ with probability ρ , where $\eta_n((s-1) \times L + \ell)$ are independent and zero-mean. Therefore, it follows that the entries of \mathbb{T} are independent and

$$\begin{aligned}\mathbb{E}[\mathbb{T}_{n\ell s}] &= \mathbb{E}[X_n((s-1) \times L + \ell)] \\ &= \rho f_n((s-1) \times L + \ell) \\ &= \rho \mathbf{T}_{n\ell s}.\end{aligned}$$

That is, $\mathbb{E}[\mathbb{T}] = \rho \mathbf{T}$. This concludes the proof.

K.2 Proof of Proposition 6.2

From Property 6.4, and our choice of parameter L for mSSA ($L = \sqrt{\min(N, T)T}$) and tSSA ($L = \sqrt{T}$), we have that

$$\text{ImpErr}(N, T; \text{tSSA}) = \tilde{\Theta} \left(\frac{1}{\min(N, \sqrt{T})^2} \right) = \tilde{\Theta} \left(\frac{1}{\min(N^2, T)} \right),\tag{66}$$

$$\text{ImpErr}(N, T; \text{mSSA}) = \tilde{\Theta} \left(\frac{1}{\sqrt{\min(N, T)T}} \right),\tag{67}$$

$$\text{ImpErr}(N, T; \text{ME}) = \tilde{\Theta} \left(\frac{1}{\min(N, T)} \right).\tag{68}$$

We proceed in cases.

Case 1: $T = o(N)$. In this case, from (66), (67), and (68), we have

$$\text{ImpErr}(N, T; \text{tSSA}), \text{ImpErr}(N, T; \text{mSSA}), \text{ImpErr}(N, T; \text{ME}) = \tilde{\Theta}\left(\frac{1}{T}\right)$$

Case 2: $N = o(T)$. In this case, from (66), (67), and (68), we have

$$\text{ImpErr}(N, T; \text{tSSA}) = \tilde{\Theta}\left(\frac{1}{N^2}\right), \quad (69)$$

$$\text{ImpErr}(N, T; \text{mSSA}) = \tilde{\Theta}\left(\frac{1}{\sqrt{NT}}\right), \quad (70)$$

$$\text{ImpErr}(N, T; \text{ME}) = \tilde{\Theta}\left(\frac{1}{N}\right).$$

In this case, we have

$$\text{ImpErr}(N, T; \text{tSSA}), \text{ImpErr}(N, T; \text{mSSA}) = \tilde{o}(\text{ImpErr}(N, T; \text{ME})).$$

It remains to compare the relative performance of tSSA and mSSA for the regime $N = o(T)$. Towards this, note from (69) and (70) that

$$\begin{aligned} \text{ImpErr}(N, T; \text{tSSA}) &= \tilde{o}(\text{ImpErr}(N, T; \text{mSSA})) \\ \iff \frac{1}{N^2} &= \tilde{o}\left(\frac{1}{\sqrt{NT}}\right) \\ \iff T^{1/3} &= o(N) \end{aligned}$$

This completes the proof.

K.3 Proof of Proposition C.1

Proposition K.1. *Let Properties C.1, 2.2, and 2.3 hold. Then, for any $1 \leq L \leq \sqrt{T}$, \mathbf{HT} has CP-rank at most $R \times G$. Further, all entries of \mathbb{HT} are independent random variables with each entry observed with probability $\rho \in (0, 1]$, and $\mathbb{E}[\mathbb{HT}] = \rho \mathbf{HT}$.*

Consider $n_1, \dots, n_d \in [N_1] \times \dots \times [N_d]$, $\ell \in [L]$, $s \in [T/L]$. By Property C.1,

$$\begin{aligned} \mathbf{HT}_{n_1, \dots, n_d, \ell, s} &= f_{n_1, \dots, n_d}((s-1) \times L + \ell) \\ &= \sum_{r=1}^R U_{n_1, r} \dots U_{n_d, r} W_{r, ((s-1) \times L + \ell)}, \end{aligned}$$

The rest of the proof follows in a similar fashion to that of Proposition 6.1.