

# SARMA: A COMPUTATIONALLY SCALABLE HIGH-DIMENSIONAL TIME SERIES MODEL

BY FEIQING HUANG<sup>1</sup>, YAO ZHENG<sup>2,‡</sup>, KEXIN LU<sup>1,\*</sup> AND GUODONG LI<sup>1,†</sup>

<sup>1</sup>*Department of Statistics, University of Hong Kong, [amieehuang@hku.hk](mailto:amieehuang@hku.hk); [neithen@hku.hk](mailto:neithen@hku.hk); [gqli@hku.hk](mailto:gqli@hku.hk)*

<sup>2</sup>*Department of Statistics, University of Connecticut, [yao.zheng@uconn.edu](mailto:yao.zheng@uconn.edu)*

This paper introduces a novel parametric infinite-order vector autoregressive model. As a variant of the vector autoregressive moving average (ARMA) model, it not only inherits desirable properties such as parsimony and rich temporal dependence structures, but also avoids two well-known drawbacks of the former: (i) non-identifiability and (ii) computational intractability even for moderate-dimensional data. Moreover, its parameter estimation is scalable with respect to the complexity of temporal dependence, namely the number of decay patterns constituting the autoregressive structure; hence it is called the scalable ARMA (SARMA) model. In the high-dimensional setup, we further impose a low-Tucker-rank assumption on the coefficient tensor of the proposed model. The resulting model has the form of a regression with embedded dynamic factors and hence can be especially suited for financial and economic data. Non-asymptotic error bounds for the proposed estimator are derived, and a tractable alternating least squares algorithm is developed. Theoretical and computational properties of the proposed method are verified by simulation studies, and the advantages over existing methods are illustrated in real applications.

**1. Introduction.** The advent of big data era has spurred a rapid growth of interest in high-dimensional time series modeling in recent years. The availability of large temporal datasets enables us to better capture the dependence structure of multivariate time series along both time and variable dimensions. Applications can be found in many areas such as economics, finance, environmental science, biology, and neuroscience [17, 18, 14]. Arguably, as the most widely used time series model, the vector autoregressive (AR) model has been a

---

*MSC2020 subject classifications:* Primary 62M10; secondary 62H12, 60G10.

*Keywords and phrases:* high-dimensional time series, identifiability, reduced-rank regression, scalability, tensor decomposition, vector  $AR(\infty)$ , vector ARMA.

primary workhorse for high-dimensional tasks thanks to dimension reduction techniques such as sparsity-inducing methods (see a recent survey by [7]), parameter constraints (see [41] and the references therein), and reduced-rank methods (see, e.g., [29] and [6]). However, the flexibility of the vector AR model can be seriously limited by the finite AR order in practice, since information further in the past is often needed for fully capturing the complex temporal dependence. This will result in poorer forecasting accuracy when the time series is longer. Instead, it is more realistic to assume that the data follow a general infinite-order vector AR (i.e., vector  $\text{AR}(\infty)$ ) process; see the literature on multivariate time series, e.g., [26] and [34]. Then the infinite-order process can be estimated consistently by fitting a finite-order vector AR model with a large AR order that grows with the time series length. While this method may improve the forecasting performance, it can make the fitted model unnecessarily complicated and hard to interpret owing to the excessive number of AR coefficient matrices.

Alternatively, as a parsimonious parameterization of vector  $\text{AR}(\infty)$  processes, the vector autoregressive moving average (ARMA) model [26, 34] is widely studied:

$$(1.1) \quad \mathbf{y}_t = \sum_{i=1}^p \Phi_i \mathbf{y}_{t-i} + \boldsymbol{\varepsilon}_t - \sum_{j=1}^q \Theta_j \boldsymbol{\varepsilon}_{t-j}, \quad t \in \mathbb{Z}$$

where  $\mathbf{y}_t \in \mathbb{R}^N$  is the observed  $N$ -dimensional time series, and  $\boldsymbol{\varepsilon}_t \in \mathbb{R}^N$  is the innovation term. When  $q \geq 1$ , this model can be written as a vector  $\text{AR}(\infty)$  process, where the AR coefficient matrices can have a wide range of complex decay patterns even when  $p$  and  $q$  are small. Thus, information in the infinite past can be incorporated with great flexibility and parsimony, allowing for better interpretability and forecasting performance in practice; see [3] and [10] for empirical evidence. However, the popularity of the vector ARMA model has long been limited mainly due to two challenges, non-identifiability and computational intractability, both caused by the presence of the MA coefficient matrices  $\Theta_j$ 's; see Section 2.2 for details. There have only been a few existing methods for large vector ARMA models in the literature, such as the EM algorithm in [28], the Bayesian method in [10], the iterative ordinary least squares algorithm in [13], and the optimization-based method in [39]. All these methods tackle the non-identifiability by imposing sophisticated identification constraints, which can render the algorithm time-consuming or inevitably multi-stage.

In this paper, we first investigate the root cause of the aforementioned challenges of the vector ARMA model. As an illustration, consider model (1.1) with orders  $p = q = 1$ . Once we write it in the vector  $\text{AR}(\infty)$  form with AR coefficient matrices  $\mathbf{A}_j = \Theta^{j-1}(\Phi - \Theta)$

for  $j \geq 1$ , where  $\Phi = \Phi_1$  and  $\Theta = \Theta_1$ , it is clear that the MA coefficient matrix  $\Theta$  plays a key role in the corresponding vector ARMA model. On one hand, the multiplicative factor  $\Theta^{j-1}$  makes  $A_j$  decay exponentially fast as  $j \rightarrow \infty$ , which agrees with the intuition that, for stationary time series, the dependence between  $y_t$  and  $y_{t-j}$  will diminish quickly as the time lag  $j$  increases. The decay property allows the vector ARMA model to capture information in the infinite past while maintaining stationarity. On the other hand, the existence of  $\Theta^{j-1}$  also directly causes the two challenges. In fact, the exponential decay pattern is merely driven by the eigenvalues of  $\Theta$ . However, the eigenvectors of  $\Theta$  will interact with  $\Phi - \Theta$ , and this ultimately leads to the non-identifiability of the coefficient matrices  $\Phi$  and  $\Theta$ . Moreover, the loss function for the parameter estimation will involve a very high-order polynomial in the coefficient matrix  $\Theta$ , making the optimization intractable even for a moderately large  $N$ . The situation for general vector ARMA models is similar except that  $\Theta$  is generalized to the MA companion matrix; see Section 2.2 for more details. However, for brevity, we will stick to the notation  $\Theta$  in this section.

The above investigation motivates us to design a new parametric vector AR( $\infty$ ) model, which not only inherits the main advantages of vector ARMA models but avoids the two drawbacks of the latter. Briefly speaking, without loss of much generality, we first assume that  $\Theta$  has distinct eigenvalues, then extract these eigenvalues to generate the temporal decay of  $A_j$ 's, and finally merge all remaining ingredients of  $A_j$  into an  $N \times N \times d$  coefficient tensor  $\mathcal{G}$ , where  $d$  is a small number. This leads to a vector AR( $\infty$ ) model that avoids the problem of non-identifiability. In this model,  $A_j$ 's are parameterized by  $\mathcal{G}$  and a low-dimensional parameter vector  $\omega$ , and hence the parsimony of the vector ARMA model is also preserved. Meanwhile, the polynomial in  $\Theta$  is decomposed as a weighted sum of one or two-dimensional polynomials, where each polynomial represents a particular decay pattern, and the weights are matrices controlling the relative strength of the decay patterns. As a direct consequence of this decomposable parametric structure, the estimation will be scalable with respect to the number of decay patterns, namely the complexity of temporal dependence. Thus, we call this model the scalable ARMA (SARMA) model.

Similar to the vector ARMA model, the number of parameters in the SARMA model is in the order of  $N^2$ , as is determined by the coefficient tensor  $\mathcal{G}$ . Therefore, when fitting the model to high-dimensional time series, it is necessary to restrict the parameter space of

$\mathcal{G}$  to a reasonable number of degrees of freedom. While following the literature on high-dimensional regression, one may consider sparsity-inducing methods such as the Lasso, this paper advocates a different approach for the following two reasons. First, unlike the independent case, the Lasso may fail to identify the true sparsity patterns of the parameters for time series data, due to the presence of strong temporal dependence (see, e.g., [30]) and the fact that the stationarity of the time series will force the entries of all  $N \times N$  coefficient matrices to shrink to zero as  $N \rightarrow \infty$  (see Remark 1 in [38] and Theorem 1 in Section 2.3). Second, in economic and financial time series, there is typically strong cross-sectional dependence among the  $N$  variables. Meanwhile, for such applications  $N$  and  $T$  are often both large in the high-dimensional setup, which is also known as the data-rich environment in econometrics; see, e.g., [27] and the references therein. This is different from, for example, time-course gene expression data, where usually the dependence among the  $N$  genes is believed to be sparse, and  $T$  is relatively small compared to  $N$  [25]. The strong (i.e., non-sparse) cross-sectional dependence can be better characterized by assuming that the  $N$  series are driven by a small number of common latent factors. Indeed, in the econometric literature, high-dimensional data are usually analyzed by factor models, where the estimated factors and loadings can provide insights into the interactive mechanism of large economic and financial systems [4, 24, 5]. Nonetheless, unlike AR and ARMA-type time series models, the factor model by itself cannot be directly used for forecasting.

Motivated by the above considerations, this paper employs a low-Tucker-rank approach that can not only further reduce the dimension of  $\mathcal{G}$  but extract common latent factors across the  $N$  variables, thereby exploiting the advantages of factor modeling. More importantly, unlike factor models in econometrics, the proposed model can be directly used for forecasting since the factors are embedded in the vector  $\text{AR}(\infty)$  framework. Specifically, let  $\mathcal{A} \in \mathbb{R}^{N \times N \times \infty}$  be the tensor formed by stacking all the AR coefficient matrices  $\mathbf{A}_j$ 's. The parametric structure of the proposed vector  $\text{AR}(\infty)$  model satisfies the tensor factorization,  $\mathcal{A} = \mathcal{G} \times_3 \mathbf{L}(\boldsymbol{\omega})$ , where  $\mathcal{G} \in \mathbb{R}^{N \times N \times d}$ , and  $\mathbf{L}$  is a known function such that  $\mathbf{L}(\boldsymbol{\omega}) \in \mathbb{R}^{\infty \times d}$  is of full rank. This directly implies that  $\mathcal{A}$  has Tucker rank  $d$  at the third mode. In addition, note that the Tucker ranks of  $\mathcal{A}$  at the first and second mode correspond to the dimensions of the column and row spaces of  $\mathbf{A}_j$ 's, respectively. Similar to the multivariate reduced-rank regression [36], it is then natural to assume that  $\mathcal{G}$  has low Tucker rank  $\mathcal{R}_i$  at mode  $i$  for

$i = 1$  and  $2$ , so that  $\mathcal{A}$  has Tucker ranks  $(\mathcal{R}_1, \mathcal{R}_2, d)$  across three modes. As a result, the degrees of freedom of  $\mathcal{G}$  can be restricted simultaneously at its first two modes, and  $\mathcal{A}$  can be further factorized as  $\mathcal{A} = \mathcal{S} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{L}(\boldsymbol{\omega})$ , where  $\mathcal{S} \in \mathbb{R}^{\mathcal{R}_1 \times \mathcal{R}_2 \times d}$ , and  $\mathbf{U}_i \in \mathbb{R}^{N \times \mathcal{R}_i}$  for  $i = 1, 2$ . Interestingly, the low-Tucker-rankness of  $\mathcal{A}$  implies different factor structures in three directions: response variables', predictor variables', and predictor time lags' directions; that is, the concept of factor modeling is embedded in the proposed model. Moreover, we can show that the  $\mathcal{R}_1$ -dimensional space for the response  $\mathbf{y}_t$  recovers exactly the latent factor space of the factor model [24, 5]. However, the proposed low-Tucker-rank framework provides more flexibility in characterizing the cross-sectional dependence by allowing for a different,  $\mathcal{R}_2$ -dimensional space for the predictors  $\mathbf{y}_{t-j}$ 's with  $j \geq 1$ . This flexibility is essential for forecasting frameworks, since although  $\mathbf{y}_t$  and  $\mathbf{y}_{t-j}$ 's correspond to the same  $N$  variables, the former is the output of the dynamic system, whereas the latter is the input. For a more detailed discussion on the connection with factor models, see Section 2.4.

This paper has three main contributions:

- (i) A new multivariate time series model, the SARMA model, is proposed, which preserves the parsimony and exponentially decaying temporal dependence patterns of the vector ARMA model, but avoids the challenges in identification and computation. A sufficient condition for the stationarity of the model is established.
- (ii) For high-dimensional modeling and forecasting, a low-Tucker-rank estimation method is introduced to embed factor-based dimensionality reduction in the proposed model, and the corresponding algorithm is scalable with respect to the complexity of temporal dependence. Additionally, a procedure for consistent model selection is developed.
- (iii) The nonasymptotic analysis of the proposed estimator involves new techniques for handling nonlinear and nonconvex parametric structure for high-dimensional time series models, and may be of independent interest.

The remainder of this paper is organized as follows. In Section 2, a reparameterization of the vector ARMA model is first derived, which leads to the proposed model, and then the high-dimensional low-Tucker-rank model is introduced. Section 3 presents the proposed estimator, algorithm, and model selection procedure, together with theoretical properties. Simulation studies and empirical examples are provided in Sections 4 and 5, respectively. Section 6 concludes with a discussion on some interesting extensions. We collect the most

important theoretical analysis in the Appendix, while all remaining proofs and some further details for Sections 5 and 6 are given in a separate supplementary file.

## 2. Scalable autoregressive moving average models .

**2.1. Tensor algebra and notations.** Tensors, a.k.a. multidimensional arrays, are natural high-order extensions of matrices; see [23] for a review of basic tensor algebra. The order of a tensor is known as the dimension, way or mode. This paper will focus on the third-order tensor  $\mathcal{X} \in \mathbb{R}^{d_1 \times d_1 \times d_3}$ . Its mode-1 matricization  $\mathcal{X}_{(1)}$  is the  $d_1$ -by- $(d_2 d_3)$  matrix whose  $\{i, (k-1)d_2 + j\}$ -th entry is  $\mathcal{X}_{ijk}$  for  $1 \leq i \leq d_1, 1 \leq j \leq d_2$  and  $1 \leq k \leq d_3$ . The mode-2 and mode-3 matricizations can be defined similarly. For any sequence of  $d_1 \times d_2$  matrices  $\{\mathbf{X}_i, 1 \leq i \leq d_3\}$ , where  $d_3$  may be infinity, we define the stacking operation  $\text{stack}(\cdot)$  such that  $\mathcal{X} = \text{stack}(\{\mathbf{X}_i\})$  is the  $d_1 \times d_2 \times d_3$  tensor obtained by stacking all the matrices along the third mode. Note that  $\mathbf{X}_i$ 's are called the frontal slices of  $\mathcal{X}$ , and  $\mathcal{X}_{(1)} = (\mathbf{X}_1, \dots, \mathbf{X}_{d_3})$ . Furthermore, the mode-1 multiplication of a tensor  $\mathcal{X} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  by a matrix  $\mathbf{Y} \in \mathbb{R}^{q_1 \times d_1}$  leads to a  $q_1 \times d_2 \times d_3$  tensor,  $\mathcal{X} \times_1 \mathbf{Y} = (\sum_{i=1}^{d_1} \mathcal{X}_{ijk} \mathbf{Y}_{si})_{1 \leq s \leq q_1, 1 \leq j \leq d_2, 1 \leq k \leq d_3}$ ; the mode-2 and mode-3 multiplications can be defined similarly. The Tucker rank at mode  $i$  of  $\mathcal{X}$  is defined as the rank of  $\mathcal{X}_{(i)}$  for  $1 \leq i \leq 3$  [35, 12]. Unlike row and column ranks of a matrix,  $r_1, r_2$  and  $r_3$  in general are not identical. Suppose that  $\mathcal{X}$  has Tucker ranks  $r_i = \text{rank}(\mathcal{X}_{(i)})$  for  $1 \leq i \leq 3$ . Then there exists a Tucker decomposition,  $\mathcal{X} = \mathcal{Y} \times_1 \mathbf{Y}_1 \times_2 \mathbf{Y}_2 \times_3 \mathbf{Y}_3$ , where  $\mathcal{Y} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$  is the core tensor, and  $\mathbf{Y}_i \in \mathbb{R}^{d_i \times r_i}$  with  $1 \leq i \leq 3$  are factor matrices.

Throughout the paper, unless otherwise specified, we will denote scalars by lowercase letters  $x, y, \dots$ , vectors by boldface lowercase letters  $\mathbf{x}, \mathbf{y}, \dots$ , matrices by boldface capital letters  $\mathbf{X}, \mathbf{Y}, \dots$ , and tensors by calligraphic capital letters  $\mathcal{X}, \mathcal{Y}, \dots$ . For any  $a, b \in \mathbb{R}$ , denote  $a \vee b = \max\{a, b\}$  and  $a \wedge b = \min\{a, b\}$ . For any vector  $\mathbf{x}$ , denote its  $\ell_2$  norm by  $\|\mathbf{x}\|_2$ . For any matrix  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$ , let  $\sigma_1(\mathbf{X}) \geq \sigma_2(\mathbf{X}) \geq \dots \geq \sigma_{d_1 \wedge d_2}(\mathbf{X}) \geq 0$  be its singular values in descending order. Let  $\mathbf{X}'$ ,  $\sigma_{\max}(\mathbf{X})$  (or  $\sigma_{\min}(\mathbf{X})$ ),  $\lambda_{\max}(\mathbf{X})$  (or  $\lambda_{\min}(\mathbf{X})$ ),  $\text{colsp}(\mathbf{X})$ , and  $\text{rank}(\mathbf{X})$  denote its transpose, largest (or smallest) singular value, largest (or smallest) eigenvalue, column space, and rank, respectively. Its vectorization  $\text{vec}(\mathbf{X})$  is the long vector obtained by stacking all its columns. In addition, its operator norm, Frobenius norm, and nuclear norm are  $\|\mathbf{X}\|_{\text{op}} = \sigma_{\max}(\mathbf{X})$ ,  $\|\mathbf{X}\|_{\text{F}} = \sqrt{\sum_{i,j} \mathbf{X}_{ij}^2} = \sqrt{\sum_{k=1}^{d_1 \wedge d_2} \sigma_k^2(\mathbf{X})}$ , and  $\|\mathbf{X}\|_* = \sum_{k=1}^{d_1 \wedge d_2} \sigma_k(\mathbf{X})$ , respectively. The Frobenius norm of a tensor  $\mathcal{X}$  is  $\|\mathcal{X}\|_{\text{F}} = \sqrt{\sum_{i,j,k} \mathcal{X}_{ijk}^2}$ .

For any real-valued function  $f$ , let  $\nabla f(\mathbf{x}) \in \mathbb{R}^m$  and  $\nabla^2 f(\mathbf{x}) \in \mathbb{R}^{m \times m}$  be its first- and second-order derivatives with respect to a scalar or vector  $\mathbf{x} \in \mathbb{R}^m$ , respectively. For any scalar-, vector-, or matrix-valued function  $\mathbf{f}(\boldsymbol{\eta})$  with  $\boldsymbol{\eta} = (\gamma, \theta)' \in \mathbb{R}^2$ , let  $\nabla_\gamma \mathbf{f}(\boldsymbol{\eta}) = \partial \mathbf{f}(\boldsymbol{\eta}) / \partial \gamma$  and  $\nabla_\theta \mathbf{f}(\boldsymbol{\eta}) = \partial \mathbf{f}(\boldsymbol{\eta}) / \partial \theta$  be its first-order partial derivatives, and they both have the same size as  $\mathbf{f}(\boldsymbol{\eta})$ . For any two sequences  $x_n$  and  $y_n$ , denote  $x_n \lesssim y_n$  (or  $x_n \gtrsim y_n$ ) if there exists an absolute constant  $C > 0$  such that  $x_n \leq C y_n$  (or  $x_n \geq C y_n$ ). Write  $x_n \asymp y_n$  if  $x_n \lesssim y_n$  and  $x_n \gtrsim y_n$ . Let  $\mathbb{I}_{\{\cdot\}}$  be the indicator function taking value one when the condition is true and zero otherwise. The capital letters  $C, C_g, \dots$  and lowercase letters  $c, c_g, \dots$  represent generic large and small positive absolute constants, respectively, whose values may vary from place to place.

**2.2. Reparameterization of vector ARMA models.** Consider a simple stationary vector ARMA(1, 1) model,  $\mathbf{y}_t = \boldsymbol{\Phi} \mathbf{y}_{t-1} + \boldsymbol{\varepsilon}_t - \boldsymbol{\Theta} \boldsymbol{\varepsilon}_{t-1}$ , for an  $N$ -dimensional time series  $\{\mathbf{y}_t\}_{t=1}^T$ . We assume that it is invertible, i.e., all eigenvalues of  $\boldsymbol{\Theta}$  are less than one in absolute value, and then it can be written into the following vector AR( $\infty$ ) form,

$$(2.1) \quad \mathbf{y}_t = \sum_{j=1}^{\infty} \underbrace{\boldsymbol{\Theta}^{j-1}(\boldsymbol{\Phi} - \boldsymbol{\Theta})}_{\mathbf{A}_j} \mathbf{y}_{t-j} + \boldsymbol{\varepsilon}_t,$$

where  $\mathbf{A}_j = \mathbf{A}_j(\boldsymbol{\Phi}, \boldsymbol{\Theta}) = \boldsymbol{\Theta}^{j-1}(\boldsymbol{\Phi} - \boldsymbol{\Theta})$  with  $j \geq 1$  are specified in matrix multiplicative forms. Note that the common factor  $\boldsymbol{\Phi} - \boldsymbol{\Theta}$  in  $\mathbf{A}_j$ 's can be regarded as a shifted version of the AR coefficient matrix  $\boldsymbol{\Phi}$  in the vector ARMA model. It is well known that when fitting the vector ARMA model to real-world datasets, two challenges will be encountered:

- (i) *Non-identifiability*: There generally exist many different combinations of  $(\boldsymbol{\Theta}, \boldsymbol{\Phi})$  that lead to the same values for  $\{\mathbf{A}_1, \mathbf{A}_2, \dots\}$  and hence the same data generating process.
- (ii) *Computational intractability*: When estimating  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Phi}$  in a sample of size  $T$ , for instance, by the least squares method, due to the form in (2.1), the loss function will be a matrix polynomial in  $\boldsymbol{\Theta}$  of degree as high as  $2(T-1)$ . This will make the optimization intractable even for moderate dimension  $N$ ; see Section 4.2 for numerical evidence.

To overcome these challenges, we aim to block-diagonalize the MA coefficient matrix  $\boldsymbol{\Theta}$  such that the high-order term  $\boldsymbol{\Theta}^{j-1}$  in (2.1) can be more easily computed.

**LEMMA 2.1** (Theorem 1 in [20]). *Real matrices with  $\mathcal{R}$  distinct nonzero eigenvalues are dense in the set of all  $N \times N$  real matrices with rank at most  $\mathcal{R}$ , where  $0 < \mathcal{R} \leq N$ .*

By Lemma 2.1, without loss of much generality, we can assume that all of the  $\mathcal{R}$  nonzero eigenvalues of  $\Theta$  are distinct, where  $\mathcal{R} = r + 2s \leq N$  is the rank of  $\Theta$ . Specifically, suppose that  $\Theta$  has  $r$  distinct nonzero real eigenvalues,  $\lambda_1, \dots, \lambda_r$ , and  $s$  distinct conjugate pairs of nonzero complex eigenvalues,  $\lambda_{r+1}, \dots, \lambda_{r+2s}$ , where  $|\lambda_j| \in (0, 1)$  for  $1 \leq j \leq r$ ,  $(\lambda_{r+2k-1}, \lambda_{r+2k}) = (\gamma_k e^{i\theta_k}, \gamma_k e^{-i\theta_k})$  with  $\gamma_k \in (0, 1)$  and  $\theta_k \in (-\pi/2, \pi/2)$  for  $1 \leq k \leq s$ , and  $i$  represents the imaginary unit. Then, the Jordan decomposition  $\Theta = BJB^{-1}$  exists such that  $B \in \mathbb{R}^{N \times N}$  is invertible and  $J \in \mathbb{R}^{N \times N}$  is the real Jordan form; that is,  $J = \text{diag}\{\lambda_1, \dots, \lambda_r, C_1, \dots, C_s, \mathbf{0}\}$  is a real block diagonal matrix, with

$$(2.2) \quad C_k = \gamma_k \cdot \begin{pmatrix} \cos(\theta_k) & \sin(\theta_k) \\ -\sin(\theta_k) & \cos(\theta_k) \end{pmatrix} \in \mathbb{R}^{2 \times 2}, \quad 1 \leq k \leq s;$$

see [22, Chap. 3]. We make use of the block-diagonal structure of  $J$  to reparameterize the coefficient matrices of the vector AR( $\infty$ ) form in (2.1) into

$$(2.3) \quad \begin{aligned} A_1 &= \Phi - \Theta := G_1 \quad \text{and} \\ A_{1+j} &= BJ^j B^{-1}(\Phi - \Theta) \\ &:= \sum_{k=1}^r \lambda_k^j G_k^I + \sum_{k=1}^s \gamma_k^j \left\{ \cos(j\theta_k) G_k^{II,1} + \sin(j\theta_k) G_k^{II,2} \right\}, \quad j \geq 1, \end{aligned}$$

where  $G_k^I$ 's,  $G_k^{II,1}$ 's and  $G_k^{II,2}$ 's are  $N \times N$  real matrices determined jointly by  $B$  and  $B^{-1}(\Phi - \Theta)$  such that

$$(2.4) \quad \text{rank}(G_j^I) \leq 1, \quad \text{rank}(G_k^{II,1}) \leq 2, \quad \text{rank}(G_k^{II,2}) \leq 2, \quad \forall 1 \leq j \leq r, 1 \leq k \leq s,$$

and, for each  $1 \leq k \leq s$ ,  $G_k^{II,1}$  and  $G_k^{II,2}$  have the same row and column spaces; see the proof of Proposition 2.2 for details. We can view (2.3) as a new parametric vector AR( $\infty$ ) model, whose unknown parameters are  $G_1$ ,  $G_k^I$ 's,  $G_k^{II,1}$ 's,  $G_k^{II,2}$ 's,  $\lambda_k$ 's,  $\gamma_k$ 's and  $\theta_k$ 's.

The above reparameterization circumvents the challenges of the vector ARMA model mentioned in (i) and (ii). To better understand how this is possible, consider the MA coefficient matrix  $\Theta$ . Since  $\Theta = BJB^{-1}$ , it plays two different roles through  $J$  and  $B$ : while  $J$  induces the temporal decay of  $A_j$  via the eigenvalues,  $B$  is essentially a redundant factor irrelevant to the temporal dynamic. However,  $B$  will interact with the shifted AR coefficient matrix  $\Phi - \Theta$  through the matrix multiplicative form of  $A_j$ 's, leading to the identification problem of the vector ARMA model. To avoid this pitfall, we first extract the ingredients of  $\Theta$  that generate the temporal decay of  $A_{1+j}$ , i.e., the scalars  $\lambda_k^j$ 's,  $\gamma_k^j$ 's and  $j\theta_k$ 's. Then we



merge all other ingredients of  $\mathbf{A}_{1+j}$  into the  $\mathbf{G}$ -matrices. Consequently,  $\mathbf{A}_{1+j}$ 's in (2.3) are simply linear combinations of the  $\mathbf{G}$ -matrices. By getting rid of the matrix multiplications, the identification problem will be avoided. As a by-product, since there is no longer any high-order matrix polynomial involved, the computational challenge is also overcome. The cost of these desirable properties is the addition of a small number of unknown parameters: the vector ARMA(1,1) model without any constraint on  $\Phi$  has  $N^2 + 2N(r + 2s) - (r + 2s)^2$  parameters, while model (2.3) has  $N^2 + (2N - 1)(r + 2s)$  parameters with the low-rank properties of the  $\mathbf{G}$ -matrices in (2.4) taken into account.

In general, for a stationary and invertible vector ARMA( $p, q$ ) model in the form of (1.1), the vector AR( $\infty$ ) representation can be written as

$$(2.5) \quad \mathbf{y}_t = \sum_{j=1}^{\infty} \underbrace{\left( \sum_{i=0}^{p \wedge j} \mathbf{P} \underline{\Theta}^{j-i} \mathbf{P}' \Phi_i \right)}_{\mathbf{A}_j} \mathbf{y}_{t-j} + \varepsilon_t, \quad \underline{\Theta} = \begin{pmatrix} \Theta_1 & \Theta_2 & \cdots & \Theta_{q-1} & \Theta_q \\ \mathbf{I} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{I} & \mathbf{0} \end{pmatrix},$$

where  $\Phi_0 = -\mathbf{I}$  and  $\mathbf{P} = (\mathbf{I}_N, \mathbf{0}_{N \times N(q-1)})$  are known matrices, and  $\underline{\Theta}$  is called the MA companion matrix [26]. Similarly, to block-diagonalize  $\underline{\Theta}$ , we consider its Jordan decomposition:  $\underline{\Theta} = \mathbf{B} \mathbf{J} \mathbf{B}^{-1}$ , where  $\mathbf{B} \in \mathbb{R}^{Nq \times Nq}$  is invertible, and  $\mathbf{J}$  is the real Jordan form of  $\underline{\Theta}$ . Let  $\tilde{\mathbf{B}} = \mathbf{P} \mathbf{B}$  and  $\tilde{\mathbf{B}} = \mathbf{B}^{-1} (\sum_{i=0}^p \underline{\Theta}^{p-i} \mathbf{P}' \Phi_i)$ . The following proposition gives the general reparameterization for (2.5) without requiring that all nonzero eigenvalues of  $\underline{\Theta}$  are distinct.

**PROPOSITION 2.2.** Suppose that there are  $r$  distinct nonzero real eigenvalues of  $\underline{\Theta}$ ,  $\lambda_j$  for  $1 \leq j \leq r$ , and  $s$  distinct conjugate pairs of nonzero complex eigenvalues of  $\underline{\Theta}$ ,  $(\lambda_{r+2k-1}, \lambda_{r+2k}) = (\gamma_k e^{i\theta_k}, \gamma_k e^{-i\theta_k})$  with  $\gamma_k \in (0, 1)$  and  $\theta_k \in (-\pi/2, \pi/2)$  for  $1 \leq k \leq s$ . Moreover, the algebraic multiplicity of  $\lambda_j$  is  $n_j$  for  $1 \leq j \leq r$ , and that of  $(\lambda_{r+2k-1}, \lambda_{r+2k})$  is  $m_k$  for  $1 \leq k \leq s$ , so there are  $R + 2S$  nonzero eigenvalues of  $\underline{\Theta}$  in total, where  $R = \sum_{k=1}^r n_k$  and  $S = \sum_{k=1}^s m_k$ . Assume that the geometric multiplicities of all nonzero eigenvalues are

one. Then for all  $j \geq 1$ , we have

$$\begin{aligned}
 \mathbf{A}_j &= \sum_{k=1}^p \mathbb{I}_{\{j=k\}} \mathbf{G}_k + \sum_{k=1}^r \sum_{i=1}^{n_k} \mathbb{I}_{\{j \geq p+(i-1) \vee 1\}} \lambda_k^{j-p-i+1} \binom{j-p}{i-1} \mathbf{G}_{k,i}^I \\
 (2.6) \quad &+ \sum_{k=1}^s \sum_{i=1}^{m_k} \mathbb{I}_{\{j \geq p+(i-1) \vee 1\}} \gamma_k^{j-p-i+1} \binom{j-p}{i-1} \\
 &\cdot \left[ \cos\{(j-p-i+1)\theta_k\} \mathbf{G}_{k,i}^{II,1} + \sin\{(j-p-i+1)\theta_k\} \mathbf{G}_{k,i}^{II,2} \right],
 \end{aligned}$$

where the first term is suppressed if  $p = 0$ , and  $\mathbf{G}_{k,i}^I$ 's,  $\mathbf{G}_{k,i}^{II,1}$ 's and  $\mathbf{G}_{k,i}^{II,2}$ 's are all determined jointly by  $\tilde{\mathbf{B}}$  and  $\check{\mathbf{B}}$ . Moreover, for any fixed  $k$  and  $i$ ,  $\mathbf{G}_{k,i}^{II,h}$  for  $h = 1, 2$  have the same row and column spaces, and  $\text{rank}(\mathbf{G}_{j,l}^I) \leq n_k$  and  $\text{rank}(\mathbf{G}_{k,i}^{II,h}) \leq 2m_k$  for all  $1 \leq j \leq r$ ,  $1 \leq k \leq s$ ,  $1 \leq l \leq n_k$ ,  $1 \leq i \leq m_k$ , and  $h = 1, 2$ .

By Proposition 2.2, if all nonzero eigenvalues of  $\underline{\Theta}$  are further assumed to be distinct, i.e.,  $n_1 = \dots = n_r = m_1 = \dots = m_s = 1$ , as its generality is supported by Lemma 2.1, then (2.6) for the vector AR( $\infty$ ) form in (2.5) can be reparameterized into

$$\begin{aligned}
 \mathbf{A}_j &= \sum_{k=1}^p \mathbb{I}_{\{j=k\}} \mathbf{G}_k + \sum_{k=1}^r \mathbb{I}_{\{j \geq p+1\}} \lambda_k^{j-p} \mathbf{G}_k^I \\
 (2.7) \quad &+ \sum_{k=1}^s \mathbb{I}_{\{j \geq p+1\}} \gamma_k^{j-p} \left[ \cos\{(j-p)\theta_k\} \mathbf{G}_k^{II,1} + \sin\{(j-p)\theta_k\} \mathbf{G}_k^{II,2} \right]
 \end{aligned}$$

for all  $j \geq 1$ , where the first term is suppressed if  $p = 0$ .

**2.3. Proposed SARMA model.** We first use tensor notations to obtain a more succinct form of (2.7). Let  $d = p + r + 2s$  and  $\boldsymbol{\omega} = (\lambda_1, \dots, \lambda_r, \boldsymbol{\eta}'_1, \dots, \boldsymbol{\eta}'_s)'$ , where  $\boldsymbol{\eta}_k = (\gamma_k, \theta_k)'$  for  $1 \leq k \leq s$ . Define a matrix-valued function of  $\boldsymbol{\omega}$  as follows:

$$\mathbf{L}(\boldsymbol{\omega}) = (\ell_{j,k}(\boldsymbol{\omega}))_{j \geq 1, 1 \leq k \leq d} = \begin{pmatrix} \mathbf{I}_p \\ \boldsymbol{\ell}^I(\lambda_1) \cdots \boldsymbol{\ell}^I(\lambda_r) \boldsymbol{\ell}^{II}(\boldsymbol{\eta}_1) \cdots \boldsymbol{\ell}^{II}(\boldsymbol{\eta}_s) \end{pmatrix} \in \mathbb{R}^{\infty \times d},$$

where, for any  $\lambda$  and  $\boldsymbol{\eta} = (\gamma, \theta)'$ , the functions  $\boldsymbol{\ell}^I(\cdot)$  and  $\boldsymbol{\ell}^{II}(\cdot)$  are defined as

$$\boldsymbol{\ell}^I(\lambda) = (\lambda, \lambda^2, \lambda^3, \dots)' \quad \text{and} \quad \boldsymbol{\ell}^{II}(\boldsymbol{\eta}) = \begin{pmatrix} \gamma \cos(\theta) & \gamma^2 \cos(2\theta) & \gamma^3 \cos(3\theta) & \cdots \\ \gamma \sin(\theta) & \gamma^2 \sin(2\theta) & \gamma^3 \sin(3\theta) & \cdots \end{pmatrix}'.$$

For simplicity, we relabel  $\mathbf{G}_{p+j} = \mathbf{G}_j^I$  for  $1 \leq j \leq r$ , and  $\mathbf{G}_{p+r+2k-1} = \mathbf{G}_k^{II,1}$  and  $\mathbf{G}_{p+r+2k} = \mathbf{G}_k^{II,2}$  for  $1 \leq k \leq s$ ; the two sets of notations may be used interchangeably when there is no confusion. By stacking all  $\mathbf{G}_k$ 's into the tensor  $\mathfrak{G} = \text{stack}(\{\mathbf{G}_j, 1 \leq j \leq$

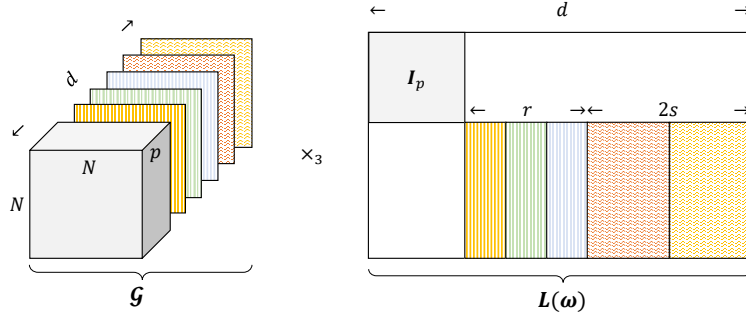


FIGURE 1. Reparameterization of  $\mathcal{A}$  for vector ARMA( $p, q$ ) models when all nonzero eigenvalues of  $\underline{\Theta}$  are distinct,  $r = 3$  and  $s = 2$ . In  $\mathbf{L}(\omega)$ , rectangles filled with vertical stripes are blocks in Type I parametric form  $\ell^I(\cdot)$ , those filled with waves are blocks in Type II parametric form  $\ell^{II}(\cdot)$ , different colors indicate dependence on different subsets of parameters, and empty spaces represent zeros.

$d\}} \in \mathbb{R}^{N \times N \times d}$  and all  $\mathbf{A}_j$ 's into the tensor  $\mathcal{A} = \text{stack}(\{\mathbf{A}_j, j \geq 1\}) \in \mathbb{R}^{N \times N \times \infty}$ , (2.7) can be written in the tensor form:

$$(2.8) \quad \mathcal{A} := \text{stack}(\{\mathbf{A}_j, j \geq 1\}) = \mathcal{G} \times_3 \mathbf{L}(\omega);$$

see Figure 1 for an illustration. In matrix form, this is equivalent to  $\mathbf{A}_j = \sum_{k=1}^d \ell_{j,k}(\omega) \mathbf{G}_k$  for  $j \geq 1$ , which is exactly (2.7).

Based on the above notation, this paper introduces a new parametric vector AR( $\infty$ ) model:

$$(2.9) \quad \mathbf{y}_t = \sum_{j=1}^{\infty} \mathbf{A}_j(\omega, \mathcal{G}) \mathbf{y}_{t-j} + \varepsilon_t, \quad \text{or} \quad \mathbf{y}_t = \sum_{k=1}^d \mathbf{G}_k \sum_{j=1}^{\infty} \ell_{j,k}(\omega) \mathbf{y}_{t-j} + \varepsilon_t,$$

where  $\varepsilon_t \in \mathbb{R}^N$  are the innovations, the coefficient tensor  $\mathcal{A}(\omega, \mathcal{G}) = \text{stack}(\{\mathbf{A}_j(\omega, \mathcal{G}), j \geq 1\})$  satisfies (2.8) (i.e., (2.7)),  $\boldsymbol{\eta}_k \in \boldsymbol{\Pi} = [0, 1) \times (-\pi/2, \pi/2)$  for  $1 \leq k \leq s$ , and  $\omega \in (-1, 1)^r \times \boldsymbol{\Pi}^s \subset \mathbb{R}^{r+2s}$ . Model (2.9) does not impose any low-rank constraints on  $\mathcal{G}$ , and it contains  $N^2d + r + 2s$  parameters in total. Note that in the matrix  $\mathbf{L}(\omega)$ , the diagonal block  $\mathbf{I}_p$  implies  $\mathbf{A}_j = \mathbf{G}_j$  for  $1 \leq j \leq p$ , while  $\ell^I(\cdot)$  and  $\ell^{II}(\cdot)$  introduce two types of decay patterns for  $\mathbf{A}_j$ 's with  $j \geq p + 1$ : the exponential decay, or a pair of exponentially damped cosine and sine waves. As a result, the magnitude of  $\ell_{j,k}(\omega)$  will decay exponentially fast as the row index  $j \rightarrow \infty$  for all  $k$ , and hence so will  $\mathbf{A}_j$ .

To better understand the temporal structure induced by the two types of decay patterns, we consider two special vector ARMA models whose MA coefficient matrices are scaled identity matrices: the vector ARMA(1, 1) with  $\Phi = \mathbf{G}^I + \lambda \mathbf{I}_N$  and  $\Theta = \lambda \mathbf{I}_N$ , and the vector ARMA(2, 2) with  $\Phi_1 = \mathbf{G}^{II,1} + 2\gamma \cos \theta \mathbf{I}_N$ ,  $\Phi_2 = \gamma \sin \theta \mathbf{G}^{II,2} - \gamma \cos \theta \mathbf{G}^{II,1} - \gamma^2 \mathbf{I}_N$ ,  $\Theta_1 = 2\gamma \cos \theta \mathbf{I}_N$ , and  $\Theta_2 = -\gamma^2 \mathbf{I}_N$ . Interestingly, their vector AR( $\infty$ ) representations are

$\mathbf{y}_t = \sum_{j=1}^{\infty} \lambda^{j-1} \mathbf{G}^I \mathbf{y}_{t-j} + \varepsilon_t$  (Type I model) and

$$\mathbf{y}_t = \sum_{j=1}^{\infty} \gamma^{j-1} [\cos\{(j-1)\theta\} \mathbf{G}^{II,1} + \sin\{(j-1)\theta\} \mathbf{G}^{II,2}] \mathbf{y}_{t-j} + \varepsilon_t \quad (\text{Type II model}),$$

which resemble the two types of decay patterns. As a result, in light of (2.7), model (2.9) in the case of  $p = 1$  can be viewed as the superimposition of  $r$  Type I  $\text{AR}(\infty)$  models,  $s$  Type II  $\text{AR}(\infty)$  models, and a vector  $\text{AR}(1)$  model,  $\mathbf{y}_t = (\mathbf{G}_1 - \sum_{k=1}^r \mathbf{G}_k^I - \sum_{k=1}^s \mathbf{G}_k^{II,1}) \mathbf{y}_{t-1} + \varepsilon_t$ ; note that the subtracted terms are because  $j$  starts with 1 rather than  $p + 1$  in the above vector  $\text{AR}(\infty)$  models. The richness of the temporal dynamic in the proposed model is thus mainly controlled by  $r$  and  $s$ , and similar results hold for the model with a general order of  $p$ .

Another important finding related to the above discussion is that the proposed model in general admits an additive decomposition in the form of

$$\mathbf{y}_t = \sum_{k=1}^p \mathbf{G}_k \mathbf{y}_{t-k} + \sum_{k=1}^r f^I(\mathbf{x}_t; \lambda_k) + \sum_{k=1}^s f^{II}(\mathbf{x}_t; \boldsymbol{\eta}_k) + \varepsilon_t,$$

where  $\mathbf{x}_t = (\mathbf{y}'_{t-1}, \mathbf{y}'_{t-2}, \dots)'$ ,  $f^I(\mathbf{x}_t; \lambda_k) = \sum_{j=p+1}^{\infty} \lambda^{j-p} \mathbf{G}_{p+k} \mathbf{y}_{t-j}$  and  $f^{II}(\mathbf{x}_t; \boldsymbol{\eta}_k) = \sum_{j=p+1}^{\infty} \gamma^{j-p} [\cos\{(j-p)\theta\} \mathbf{G}_{p+r+2k-1} + \sin\{(j-p)\theta\} \mathbf{G}_{p+r+2k}] \mathbf{y}_{t-j}$ . Here the dependence on  $\mathbf{G}_{p+1}, \dots, \mathbf{G}_d$  is suppressed for succinctness. It is noteworthy that the parameters pertaining to different decay patterns are separable; that is, each summand  $f^I(\cdot)$  or  $f^{II}(\cdot)$  involves only a particular  $\lambda_k$  or  $\boldsymbol{\eta}_k$ . As a result, the parameter estimation can be implemented efficiently via an alternating minimization algorithm, where we iterate through all  $\lambda_k$ 's and  $\boldsymbol{\eta}_k$ 's, updating one at a time while keep other terms constant. Since  $\lambda_k$ 's and  $\boldsymbol{\eta}_k$ 's are only one- or two-dimensional, the computational cost associated with any increase in  $r$  or  $s$  will be low. In other words, the estimation will be scalable with respect to the complexity of temporal dependence; see the algorithm in Section 3.2 and numerical evidence in Section 4.2. For this reason, we call model (2.9) the scalable ARMA model of orders  $(p, r, s)$ , or the  $\text{SARMA}(p, r, s)$  model. The following theorem provides a sufficient condition for the existence of a stationary causal solution of this model.

**THEOREM 2.3.** *Suppose that  $\{\varepsilon_t\}$  is a strictly stationary sequence with  $E(\|\varepsilon_t\|_2) < \infty$ . If there exists  $0 < \rho < 1$  such that*

$$(2.10) \quad \max\{|\lambda_1|, \dots, |\lambda_r|, \gamma_1, \dots, \gamma_s\} \leq \rho \quad \text{and} \quad \sum_{k=1}^p \|\mathbf{G}_k\|_{\text{op}} + \frac{\rho}{1-\rho} \sum_{k=p+1}^d \|\mathbf{G}_k\|_{\text{op}} < 1,$$

then there exists a unique strictly stationary solution to model (2.9), and it has the form of  $\mathbf{y}_t = \boldsymbol{\varepsilon}_t + \sum_{j=1}^{\infty} \boldsymbol{\Psi}_j \boldsymbol{\varepsilon}_{t-j}$ , where  $\boldsymbol{\Psi}_j = \sum_{k=1}^{\infty} \sum_{j_1+\dots+j_k=j} \mathbf{A}_{j_1} \cdots \mathbf{A}_{j_k}$  and  $\mathbf{A}_j = \sum_{k=1}^d \ell_{j,k}(\boldsymbol{\omega}) \mathbf{G}_k$  for all  $j \geq 1$ .

In addition,  $\mathcal{G}$  can be divided into two subtensors  $\mathcal{G}^{\text{AR}} := \text{stack}(\{\mathbf{G}_k, 1 \leq k \leq p\}) \in \mathbb{R}^{N \times N \times p}$  and  $\mathcal{G}^{\text{MA}} := \text{stack}(\{\mathbf{G}_k, p+1 \leq k \leq d\}) \in \mathbb{R}^{N \times N \times (r+2s)}$ , corresponding to the AR and MA parts of the model, respectively. Then we can write (2.8) in the partitioned form:

$$(2.11) \quad \mathcal{A} = \mathcal{G} \times_3 \mathbf{L}(\boldsymbol{\omega}) = \text{stack}(\mathcal{G}^{\text{AR}}, \mathcal{G}^{\text{MA}}) \times_3 \begin{pmatrix} \mathbf{I}_p \\ \mathbf{L}^{\text{MA}}(\boldsymbol{\omega}) \end{pmatrix},$$

where  $\mathbf{L}^{\text{MA}}(\boldsymbol{\omega}) = (\ell^I(\lambda_1), \dots, \ell^I(\lambda_r), \ell^{II}(\boldsymbol{\eta}_1), \dots, \ell^{II}(\boldsymbol{\eta}_s))$ . In Figure 1, the gray box represents  $\mathcal{G}^{\text{AR}}$ , the colored slices behind it represent  $\mathcal{G}^{\text{MA}}$ , and the colored rectangles in  $\mathbf{L}(\boldsymbol{\omega})$  constitute  $\mathbf{L}^{\text{MA}}(\boldsymbol{\omega})$ .

EXAMPLE 1 (SMA model). When  $p = 0$ , (2.9) reduces to the SMA( $r, s$ ) model. The AR parts in  $\mathcal{G}$  and  $\mathbf{L}(\boldsymbol{\omega})$  no longer exist, and (2.11) reduces to  $\mathcal{A} = \mathcal{G}^{\text{MA}} \times_3 \mathbf{L}^{\text{MA}}(\boldsymbol{\omega})$ . This model can be regarded as an analogue of the vector MA model.

EXAMPLE 2 (SAR model). When  $r = s = 0$ , (2.9) reduces to  $\mathbf{y}_t = \sum_{j=1}^p \mathbf{G}_j \mathbf{y}_{t-j} + \boldsymbol{\varepsilon}_t$  with  $\mathcal{A} = \mathcal{G}^{\text{AR}}$ , and the MA parts are absent in (2.11). Hence, the SAR( $p$ ) model is exactly the vector AR( $p$ ) model.

2.4. *Low-Tucker-rank model for high-dimensional time series.* Similar to vector ARMA models, the proposed SARMA model in (2.9) contains  $O(N^2)$  unknown parameters and thus requires further dimension reduction when  $N$  is large. As shown in Section 2.2, the frontal slices of  $\mathcal{G}^{\text{MA}}$  are low-rank matrices if the data follow a vector ARMA process. In light of Example 2 and the literature on high-dimensional low-rank vector AR models [36, 11, 6], it is natural to assume that the frontal slices of  $\mathcal{G}^{\text{AR}}$  are also of low rank when  $N$  is large. However, this will result in many matrix constraints, making the model hard to interpret.

For better interpretability, we consider the low-Tucker-rank assumption on  $\mathcal{G} \in \mathbb{R}^{N \times N \times d}$ . This will automatically induce low-rankness uniformly for all its frontal slices. Note that the order  $d$  is usually small. Thus, it is only necessary to assume that  $\mathcal{G}$  has low Tucker ranks at

its first two modes. We propose the following low-Tucker-rank SARMA model:

$$(2.12) \quad \mathbf{y}_t = \sum_{j=1}^{\infty} \mathbf{A}_j(\boldsymbol{\omega}, \mathcal{G}) \mathbf{y}_{t-j} + \boldsymbol{\varepsilon}_t, \quad \text{with} \quad \text{rank}(\mathcal{G}_{(i)}) \leq \mathcal{R}_i, \quad i = 1, 2,$$

for some predetermined Tucker ranks  $0 < \mathcal{R}_1, \mathcal{R}_2 \leq N$ . As a result,  $\text{rank}(\mathbf{G}_k) \leq \mathcal{R}_1 \wedge \mathcal{R}_2$  for all  $1 \leq k \leq d$ . Moreover, there exists a Tucker decomposition,  $\mathcal{G} = \mathcal{S} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2$ , where  $\mathcal{S} \in \mathbb{R}^{\mathcal{R}_1 \times \mathcal{R}_2 \times d}$  is the core tensor, and  $\mathbf{U}_i \in \mathbb{R}^{N \times \mathcal{R}_i}$  with  $i = 1$  and  $2$  are factor matrices. Let  $\mathcal{S}_k \in \mathbb{R}^{\mathcal{R}_1 \times \mathcal{R}_2}$  be the  $k$ th frontal slice of  $\mathcal{S}$  for  $1 \leq k \leq d$ , and then it holds

$$(2.13) \quad \mathcal{A}(\boldsymbol{\omega}, \mathcal{G}) = \mathcal{S} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{L}(\boldsymbol{\omega}), \quad \text{or} \quad \mathbf{A}_j(\boldsymbol{\omega}, \mathcal{G}) = \sum_{k=1}^d \ell_{j,k}(\boldsymbol{\omega}) \mathbf{U}_1 \mathcal{S}_k \mathbf{U}_2', \quad j \geq 1.$$

Note that the Tucker decomposition is not unique, since  $\mathcal{G} = \mathcal{S} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 = (\mathcal{S} \times_1 \mathbf{O}_1 \times_2 \mathbf{O}_2) \times_1 (\mathbf{U}_1 \mathbf{O}_1^{-1}) \times_2 (\mathbf{U}_2 \mathbf{O}_2^{-1})$  for any invertible matrices  $\mathbf{O}_i \in \mathbb{R}^{\mathcal{R}_i \times \mathcal{R}_i}$  with  $i = 1$  and  $2$ . We can consider the higher-order singular value decomposition (HOSVD) of  $\mathcal{G}$ , namely the special Tucker decomposition uniquely defined by choosing  $\mathbf{U}_i$  as the tall matrix consisting of the top  $\mathcal{R}_i$  left singular vectors of  $\mathcal{G}_{(i)}$  and then setting  $\mathcal{S} = \mathcal{G} \times_1 \mathbf{U}_1' \times_2 \mathbf{U}_2'$ . Note that  $\mathbf{U}_1$  is orthonormal, i.e.,  $\mathbf{U}_1' \mathbf{U}_1 = \mathbf{I}_{\mathcal{R}_1}$ . Then (2.12) implies the following dynamic factor structure:

$$(2.14) \quad \underbrace{\mathbf{U}_1' \mathbf{y}_t}_{\text{response factor}} = \sum_{k=1}^d \mathcal{S}_k \underbrace{\left\{ \sum_{j=1}^{\infty} \ell_{j,k}(\boldsymbol{\omega}) \underbrace{\mathbf{U}_2' \mathbf{y}_{t-j}}_{\text{lag-}j \text{ predictor factor}} \right\}}_{k\text{th predictor-lag factor}} + \mathbf{U}_1' \boldsymbol{\varepsilon}_t.$$

The dynamic factor regression form in (2.14) enables us to interpret model (2.12) from the viewpoint of factor modeling [32, 4, 24, 5]. Specifically,  $\mathbf{U}_1' \mathbf{y}_t$  can be viewed as an  $\mathcal{R}_1$ -dimensional response factor obtained by projecting  $\mathbf{y}_t$  onto a low-dimensional space. Similarly,  $\mathbf{U}_2' \mathbf{y}_{t-j}$  represents the  $\mathcal{R}_2$ -dimensional predictor factor at lag  $j \geq 1$ . While there are infinitely many lags,  $\mathbf{U}_2' \mathbf{y}_{t-j}$ 's are subsequently reweighted by  $\ell_{j,k}(\boldsymbol{\omega})$ 's, so that the  $k$ th temporal decay pattern is incorporated for  $1 \leq k \leq d$ . This eventually results in  $d$  predictor-lag factors, with associated weighting matrices  $\mathcal{S}_k$ 's. From (2.14), it is noteworthy that the HOSVD of  $\mathcal{G}$  at its first two modes allows for different response and predictor subspaces, since  $\mathcal{R}_1 \neq \mathcal{R}_2$  in general and the factor loadings  $\mathbf{U}_1$  and  $\mathbf{U}_2$  can be different even if  $\mathcal{R}_1 = \mathcal{R}_2$ . The difference between the two subspaces is intuitive and plausible in practice, since the response factor is the output of the dynamic system, whereas the predictor factors are the inputs.

To understand the connection between our model and the conventional factor model, similarly to (2.14), we can rewrite model (2.12) as

$$(2.15) \quad \mathbf{y}_t = \mathbf{U}_1 \sum_{k=1}^d \mathbf{S}_k \sum_{j=1}^{\infty} \ell_{j,k}(\boldsymbol{\omega}) \mathbf{U}_2' \mathbf{y}_{t-j} + \boldsymbol{\varepsilon}_t = \boldsymbol{\Lambda} \mathbf{f}_t + \boldsymbol{\varepsilon}_t,$$

where the right hand side has the form of a factor model with latent factor series  $\mathbf{f}_t \in \mathbb{R}^{\mathcal{R}_1}$  and loading matrix  $\boldsymbol{\Lambda} \in \mathbb{R}^{N \times \mathcal{R}_1}$ . There are different versions of normalization restrictions for the identification of  $\boldsymbol{\Lambda}$  and  $\mathbf{f}_t$  in the literature on factor models; see, e.g., [24] and [5]. Nevertheless, it always holds  $\text{colsp}(\boldsymbol{\Lambda}) = \text{colsp}(\mathbf{U}_1)$ . This reveals that the  $\mathcal{R}_1$ -dimensional response factor space of our model recovers the latent factor space of the factor model. In particular, if the normalization restriction in [24] is used,  $\boldsymbol{\Lambda}$  and  $\mathbf{f}_t$  are identical to  $\mathbf{U}_1$  and  $\sum_{k=1}^d \mathbf{S}_k \sum_{j=1}^{\infty} \ell_{j,k}(\boldsymbol{\omega}) \mathbf{U}_2' \mathbf{y}_{t-j}$ , respectively, up to an orthogonal rotation. Then, suppose that the estimate  $\hat{\boldsymbol{\Lambda}}$  of  $\boldsymbol{\Lambda}$  is obtained by fitting the factor model in the form of the right hand side of (2.15). The factor series  $\mathbf{f}_t$  will be estimated by  $\hat{\mathbf{f}}_t = \hat{\boldsymbol{\Lambda}}' \mathbf{y}_t$ ; see [24] for details. Moreover, by the theory for the factor model, without actually fitting our model we can get another estimator of  $\mathbf{U}_1$ , i.e.,  $\hat{\boldsymbol{\Lambda}}$ , which is consistent up to an orthogonal rotation.

Our model is also related to dynamic factor models in the literature [32, 33]. Consider the special case of model (2.12) where  $\mathbf{U}_1 = \mathbf{U}_2$ , i.e., the response and predictor factors are identical. Then it can be seen directly from (2.14) that the response (or predictor) factor series  $\{\mathbf{U}_1' \mathbf{y}_t\}$  follows an  $\mathcal{R}_1$ -dimensional SARMA model without any low-Tucker-rank structure. Thus, for this special case, a two-step dynamic factor modeling procedure is possible: first obtain the fitted factor series  $\hat{\mathbf{f}}_t$  by the factor modeling method in [24], and then fit the low-dimensional SARMA model to  $\{\hat{\mathbf{f}}_t\}$ . However, it is worth noting that this method will be inconsistent when  $\text{colsp}(\mathbf{U}_1) \not\supset \text{colsp}(\mathbf{U}_2)$ .

### 3. High-dimensional estimation.

3.1. *Rank-constrained estimation and its nonasymptotic properties.* Suppose that the observed time series  $\{\mathbf{y}_t\}_{t=1}^T$  is generated by model (2.12) with Tucker ranks  $(\mathcal{R}_1, \mathcal{R}_2)$  and model orders  $(p, r, s)$ . Let  $\boldsymbol{\Omega} \subset (-1, 1)^r \times \boldsymbol{\Pi}^s$  be the parameter space of  $\boldsymbol{\omega} = (\boldsymbol{\lambda}', \boldsymbol{\eta}')'$ , where  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_r)'$ ,  $\boldsymbol{\eta} = (\boldsymbol{\eta}'_1, \dots, \boldsymbol{\eta}'_s)'$ ,  $\boldsymbol{\eta}_k = (\gamma_k, \theta_k)'$  for  $1 \leq k \leq s$ , and  $\boldsymbol{\Pi} = [0, 1) \times (-\pi/2, \pi/2)$ . Denote by  $\boldsymbol{\omega}^*$ ,  $\mathcal{G}^*$ , and  $\mathcal{A}^*$  the true values of  $\boldsymbol{\omega}$ ,  $\mathcal{G}$ , and  $\mathcal{A}$ , respectively. Similarly,  $\lambda_j^*$ 's,  $\boldsymbol{\eta}_j^*$ 's and  $\mathbf{G}_k^*$ 's denote the true values of the corresponding parameters.

The squared error loss function for the proposed model can be written as  $\mathbb{L}_T(\omega, \mathcal{G}) = \sum_{t=1}^T \|\mathbf{y}_t - \mathcal{A}_{(1)} \mathbf{x}_t\|_2^2$ , where  $\mathbf{x}_t = (\mathbf{y}'_{t-1}, \mathbf{y}'_{t-2}, \dots)'$ . Since the loss function depends on observations in the infinite past, initial values for  $\{\mathbf{y}_s, s \leq 0\}$  will be needed in practice. We set them to zero for simplicity, and the corresponding effect will be accounted for in our theoretical analysis; see the discussion below Theorem 3.2. Let  $\tilde{\mathbf{x}}_t = (\mathbf{y}'_{t-1}, \dots, \mathbf{y}'_1, 0, 0, \dots)'$  be the initialized version of  $\mathbf{x}_t$ , and consider the rank-constrained least squares estimator:

$$(3.1) \quad \hat{\mathcal{A}} = \hat{\mathcal{G}} \times_3 \mathbf{L}(\hat{\omega}), \quad \text{with} \quad (\hat{\omega}, \hat{\mathcal{G}}) = \arg \min_{\omega \in \Omega, \mathcal{G} \in \Gamma(\mathcal{R}_1, \mathcal{R}_2)} \tilde{\mathbb{L}}_T(\omega, \mathcal{G}),$$

where  $\Gamma(\mathcal{R}_1, \mathcal{R}_2) = \{\mathcal{G} \in \mathbb{R}^{N \times N \times d} \mid \text{rank}(\mathcal{G}_{(1)}) \leq \mathcal{R}_1, \text{rank}(\mathcal{G}_{(2)}) \leq \mathcal{R}_2\}$ , and

$$\tilde{\mathbb{L}}_T(\omega, \mathcal{G}) = \sum_{t=1}^T \|\mathbf{y}_t - \mathcal{A}_{(1)} \tilde{\mathbf{x}}_t\|_2^2 = \sum_{t=1}^T \left\| \mathbf{y}_t - \sum_{j=1}^{t-1} \mathbf{A}_j(\omega, \mathcal{G}) \mathbf{y}_{t-j} \right\|_2^2.$$

After conducting the estimation in (3.1), we can further obtain the estimated orthonormal factor matrices by the HOSVD,  $\hat{\mathcal{G}} = \hat{\mathbf{S}} \times_1 \hat{\mathbf{U}}_1 \times_2 \hat{\mathbf{U}}_2$ . Then the fitted model can be interpreted from the perspective of factor modeling as in (2.14).

**ASSUMPTION 1 (Parameter space and stationarity).** (i) There exists an absolute constant  $0 < \bar{\rho} < 1$  such that for all  $\omega \in \Omega$ ,  $|\lambda_1|, \dots, |\lambda_r|, \gamma_1, \dots, \gamma_s \in \Lambda$ , where  $\Lambda$  is a compact subset of  $(0, \bar{\rho})$ ; and (ii) the model corresponding to the true values  $\omega^*$  and  $\mathcal{G}^*$  is stationary, with  $\|\mathbf{G}_k^*\|_{\text{op}} \leq C_{\mathcal{G}}$  for all  $1 \leq k \leq d$  and some absolute constant  $C_{\mathcal{G}} > 0$ .

**ASSUMPTION 2 (Identifiability).** (i) The elements of  $\omega^*$  are nonzero, where  $\lambda_1^*, \dots, \lambda_r^*$  are distinct, and  $\eta_1^*, \dots, \eta_s^*$  are distinct; and (ii) there exist absolute constants  $0 < c_{\mathcal{G}} < C_{\mathcal{G}}$  such that  $c_{\mathcal{G}} \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \leq \|\mathbf{G}_k^*\|_{\text{F}} \leq C_{\mathcal{G}} \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2}$  for  $p+1 \leq k \leq d$ .

Theorem 2.3 in the previous section provides a sufficient condition for the strict stationarity. Under Assumption 1(ii), there exists a unique strictly stationary and causal solution of the true model, given by  $\mathbf{y}_t = \Psi_*(B) \varepsilon_t$ , where  $B$  is the backshift operator,  $\Psi_*(B) = \mathbf{I}_N + \sum_{j=1}^{\infty} \Psi_j^* B^j$ , and  $\{\Psi_j^*, j \geq 1\}$  are determined by  $\omega^*$  and  $\mathcal{G}^*$ ; see Theorem 2.3. Denote

$$\mu_{\min}(\Psi_*) = \min_{|z|=1} \lambda_{\min}(\Psi_*(z) \Psi_*^{\text{H}}(z)) \quad \text{and} \quad \mu_{\max}(\Psi_*) = \max_{|z|=1} \lambda_{\max}(\Psi_*(z) \Psi_*^{\text{H}}(z)),$$

where  $\Psi_*^{\text{H}}(z)$  is the conjugate transpose of  $\Psi_*(z)$  for  $z \in \mathbb{C}$ . It can be verified that  $\mu_{\min}(\Psi_*) > 0$ ; see also [8]. Assumption 2(i) requires that the elements of  $\omega^*$  are distinguishable, which is necessary for the identifiability of the SARMA model. For an illustration,



consider model (2.9) with  $\lambda_1^* = \lambda_2^*$  and, without loss of generality, assume that  $p = 0$ . It holds  $\ell_{j,1}(\omega^*) = \ell_{j,2}(\omega^*)$  for all  $j \geq 1$ , and then the data generating process has the form of

$$\mathbf{y}_t = \sum_{k=1}^d \mathbf{G}_k^* \sum_{j=1}^{\infty} \ell_{j,k}(\omega^*) \mathbf{y}_{t-j} + \varepsilon_t = \bar{\mathbf{G}}_1^* \sum_{j=1}^{\infty} \ell_{j,1}(\omega^*) \mathbf{y}_{t-j} + \sum_{k=3}^d \mathbf{G}_k^* \sum_{j=1}^{\infty} \ell_{j,k}(\omega^*) \mathbf{y}_{t-j} + \varepsilon_t,$$

where  $\bar{\mathbf{G}}_1^* = \mathbf{G}_1^* + \mathbf{G}_2^*$ . Thus, the model is not identifiable.

To understand Assumption 2(ii), consider the role of  $\mathbf{G}_{p+1}^*, \dots, \mathbf{G}_d^*$  in model (2.9). We can regard  $\sum_{j=1}^{\infty} \ell_{j,k}(\omega^*) \mathbf{y}_{t-j}$  as the  $k$ th signal with a specific temporal decay pattern. Then  $\mathbf{G}_k^*$  can be viewed as the energy of the corresponding signal, and its overall strength can be measured by  $\|\mathbf{G}_k^*\|_F$ . Assumption 2(ii) requires that all of the last  $d - p = r + 2s$  signals are of comparable strength. In fact, we can relax this assumption to  $\|\mathbf{G}_k^*\|_F \asymp \|\mathbf{G}_j^*\|_F$  for all  $p + 1 \leq k \neq j \leq d$  or replace  $\sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2}$  with any smaller rate, and Theorem 3.2 will be unaffected; note that  $\sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2}$  is the largest possible rate implied by Assumption 1(ii) since  $\text{rank}(\mathbf{G}_k^*) \leq \mathcal{R}_1 \wedge \mathcal{R}_2$ . On the contrary, if there exists  $\mathbf{G}_{k_0}^*$  that dominates all other  $\mathbf{G}_k^*$ s with  $p + 1 \leq k \neq k_0 \leq d$ , such that  $\|\mathbf{G}_{k_0}^*\|_F / \|\mathbf{G}_k^*\|_F \rightarrow \infty$  as  $N \rightarrow \infty$ , then only the  $k_0$ -th signal (i.e., the strong signal) will be identified, while all other signals (i.e., weak signals) will be ignored by the estimation procedure. For more details about the necessity of Assumption 2(ii), see the proof of Proposition 3.1. This is similar to the concept of strong vs. weak factors in factor models [24], and variants of this assumption may be considered with a more involved technical analysis.

Note that although the proposed model is linear in  $\mathcal{A}$ , the loss function in (3.1) is non-convex with respect to  $\omega$  and  $\mathcal{G}$  jointly. In fact, it is nonconvex with respect to  $\omega$  even when fixing  $\mathcal{G}$ , since the nonlinear function  $L(\cdot)$  is nonconvex. This leads to a substantial challenge for the nonasymptotic analysis. To overcome it, a key intermediate step is to “linearize” the tensor mapping  $\mathcal{A}(\omega, \mathcal{G})$  within a fixed local neighborhood of  $\omega^*$ ; see (A.3) in the Appendix for details. We achieve this by establishing an equivalence between the perturbation of  $\mathcal{A}^*$  and those of  $\omega^*$  and  $\mathcal{G}^*$ , as shown in the proposition below. It will allow us to cover the parameter space of  $\mathcal{A}$  by covering those of  $\omega$  and  $\mathcal{G}$ , and subsequently obtain Theorem 3.2.

**PROPOSITION 3.1.** Under Assumptions 1(i) and 2, for any  $\mathcal{A} = \mathcal{G} \times_3 L(\omega)$  with  $\mathcal{G} \in \mathbb{R}^{N \times N \times d}$  and  $\omega \in \Omega$ , if  $\|\omega - \omega^*\|_2 \leq c_\omega$ , then

$$\begin{aligned}
c_\Delta \left( \|\mathcal{G} - \mathcal{G}^*\|_F + \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \|\omega - \omega^*\|_2 \right) &\leq \|\mathcal{A} - \mathcal{A}^*\|_F \\
&\leq C_\Delta \left( \|\mathcal{G} - \mathcal{G}^*\|_F + \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \|\omega - \omega^*\|_2 \right),
\end{aligned}$$

where  $c_\omega, c_\Delta, C_\Delta > 0$  are absolute constants defined in (A.2) in the Appendix.

ASSUMPTION 3 (Sub-Gaussian errors). Let  $\varepsilon_t = \Sigma_\varepsilon^{1/2} \xi_t$ , where  $\xi_t$  is a sequence of i.i.d. random vectors with zero mean and  $\text{var}(\xi_t) = \mathbf{I}_N$ , and  $\Sigma_\varepsilon$  is a positive definite covariance matrix. In addition, the coordinates  $(\xi_{it})_{1 \leq i \leq N}$  within  $\xi_t$  are mutually independent and  $\sigma^2$ -sub-Gaussian.

Assumption 3 is weaker than the commonly imposed Gaussian assumption in the literature on high-dimensional time series (see, e.g., [8] and [39]). Let

$$\kappa_1 = \lambda_{\min}(\Sigma_\varepsilon) \mu_{\min}(\Psi_*) \min\{1, c_{\bar{\rho}}^2\} \quad \text{and} \quad \kappa_2 = \lambda_{\max}(\Sigma_\varepsilon) \mu_{\max}(\Psi_*) \max\{1, C_{\bar{\rho}}^2\},$$

where  $c_{\bar{\rho}}, C_{\bar{\rho}} > 0$  are absolute constants defined in Lemma A.2 in the Appendix.

THEOREM 3.2. Let  $d_{\mathcal{M}} = \mathcal{R}_1 \mathcal{R}_2 d + (\mathcal{R}_1 + \mathcal{R}_2)N$ . Suppose that  $\|\hat{\omega} - \omega^*\|_2 \leq c_\omega$  and  $T \gtrsim (\kappa_2/\kappa_1)^2 d_{\mathcal{M}} \log(\kappa_2/\kappa_1)$ . Then under Assumptions 1–3, with probability at least  $1 - 4e^{-cd_{\mathcal{M}} \log(\kappa_2/\kappa_1)} - 8e^{-cN} - \{2 + \sqrt{\kappa_2/\lambda_{\max}(\Sigma_\varepsilon)}\} \sqrt{N/((\mathcal{R}_1 + \mathcal{R}_2)T)}$ , we have the following estimation and prediction error bounds:

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F \lesssim \sqrt{\frac{\kappa_2 \lambda_{\max}(\Sigma_\varepsilon) d_{\mathcal{M}}}{\kappa_1^2 T}} \quad \text{and} \quad \frac{1}{T} \sum_{t=1}^T \|(\hat{\mathcal{A}} - \mathcal{A}^*)_{(1)} \tilde{x}_t\|_2^2 \lesssim \frac{\kappa_2 \lambda_{\max}(\Sigma_\varepsilon) d_{\mathcal{M}}}{\kappa_1 T}.$$

The above theorem, together with Proposition 3.1, implies that, it holds with the same high probability,  $\|\hat{\omega} - \omega^*\|_2 \lesssim \sqrt{\kappa_2 \lambda_{\max}(\Sigma_\varepsilon) d_{\mathcal{M}} / \{\kappa_1^2 (\mathcal{R}_1 \wedge \mathcal{R}_2) T\}}$  and  $\|\hat{\mathcal{G}} - \mathcal{G}^*\|_F \lesssim \sqrt{\kappa_2 \lambda_{\max}(\Sigma_\varepsilon) d_{\mathcal{M}} / (\kappa_1^2 T)}$ . As a result, the consistency of  $\hat{\mathcal{A}}$ ,  $\hat{\omega}$  and  $\hat{\mathcal{G}}$  can be achieved as  $d_{\mathcal{M}}/T \rightarrow 0$ . Moreover, in the proof of Theorem 3.2 in the Appendix, we show that the effect of initial values for  $\{\mathbf{y}_s, s \leq 0\}$  is dominated by other quantities and hence has no contribution to the final estimation and prediction error rates; see the quantities  $|S_i(\hat{\Delta})|$  for  $1 \leq i \leq 3$  in the proof. The nonexponential tail probability  $\{2 + \sqrt{\kappa_2/\lambda_{\max}(\Sigma_\varepsilon)}\} \sqrt{N/((\mathcal{R}_1 + \mathcal{R}_2)T)}$  comes from naively bounding  $|S_i(\hat{\Delta})|$ 's by Markov's inequality for technical simplicity; see the proof of Lemma A.5 for details. However, this tail probability may be sharpened by employing more sophisticated concentration inequalities.

3.2. *Alternating least squares algorithm.* Applying the Tucker decomposition for  $\mathcal{A}$  in (2.13), the loss function in (3.1) can be equivalently written as

$$\tilde{\mathbb{L}}_T(\boldsymbol{\omega}, \mathcal{S}, \mathbf{U}_1, \mathbf{U}_2) = \sum_{t=1}^T \|\mathbf{y}_t - \mathcal{A}_{(1)} \tilde{\mathbf{x}}_t\|_2^2 = \sum_{t=1}^T \|\mathbf{y}_t - \{\mathcal{S} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{L}(\boldsymbol{\omega})\}_{(1)} \tilde{\mathbf{x}}_t\|_2^2,$$

where  $\mathcal{S} \in \mathbb{R}^{\mathcal{R}_1 \times \mathcal{R}_2 \times d}$ , and  $\mathbf{U}_i \in \mathbb{R}^{N \times \mathcal{R}_i}$  for  $i = 1, 2$ . Thus, the estimation can be implemented as  $(\hat{\boldsymbol{\omega}}, \hat{\mathcal{S}}, \hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) = \arg \min \tilde{\mathbb{L}}_T(\boldsymbol{\omega}, \mathcal{S}, \mathbf{U}_1, \mathbf{U}_2)$ , and then the estimator of  $\mathcal{A}$  is  $\hat{\mathcal{A}} = \hat{\mathcal{G}} \times_3 \mathbf{L}(\hat{\boldsymbol{\omega}})$ , where  $\hat{\mathcal{G}} = \hat{\mathcal{S}} \times_1 \hat{\mathbf{U}}_1 \times_2 \hat{\mathbf{U}}_2$ . Note that  $\mathbf{U}_i$ 's need not be subject to any orthonormal constraint in this minimization; i.e., the outputs  $\hat{\mathbf{U}}_i$ 's from the algorithm need not be orthonormal. Instead, to recover the orthonormal factor loading matrices  $\mathbf{U}_i$ 's such that the fitted model can be interpreted via the dynamic factor regression form in (2.14), we will conduct the HOSVD of  $\hat{\mathcal{G}}$  to get orthonormal estimates of  $\mathbf{U}_i$ 's after the optimization procedure.

Let  $\mathbf{z}_t(\boldsymbol{\omega}) = (z'_{t,1}(\boldsymbol{\omega}), \dots, z'_{t,d}(\boldsymbol{\omega}))' = \{\mathbf{L}'(\boldsymbol{\omega}) \otimes \mathbf{I}_N\} \tilde{\mathbf{x}}_t \in \mathbb{R}^{Nd}$ , where each  $N$ -dimensional vector  $z_{t,k}(\boldsymbol{\omega}) = \sum_{j=1}^{t-1} \ell_{j,k}(\boldsymbol{\omega}) \mathbf{y}_{t-j}$ . In addition, denote  $\mathbf{Z}_t(\boldsymbol{\omega}) = (z_{t,1}(\boldsymbol{\omega}), \dots, z_{t,d}(\boldsymbol{\omega})) \in \mathbb{R}^{N \times d}$  and  $\tilde{\mathbf{X}}_t = (\mathbf{y}_{t-1}, \dots, \mathbf{y}_1, \mathbf{0}) \in \mathbb{R}^{N \times (T-1)}$ . It can be verified that

$$\begin{aligned} \mathcal{A}_{(1)} \tilde{\mathbf{x}}_t &= \left[ \{z'_t(\boldsymbol{\omega})(\mathbf{I}_d \otimes \mathbf{U}_2) \mathcal{S}'_{(1)}\} \otimes \mathbf{I}_N \right] \text{vec}(\mathbf{U}_1) = \mathbf{U}_1 \mathcal{S}_{(1)} \{ \mathbf{Z}'_t(\boldsymbol{\omega}) \otimes \mathbf{I}_{R_2} \} \text{vec}(\mathbf{U}'_2) \\ &= \left[ \{z'_t(\boldsymbol{\omega})(\mathbf{I}_d \otimes \mathbf{U}_2)\} \otimes \mathbf{U}_1 \right] \text{vec}(\mathcal{S}_{(1)}) = \mathbf{U}_1 \mathcal{G}_{(1)} (\mathbf{I}_d \otimes (\mathbf{U}'_2 \tilde{\mathbf{X}}_t)) \text{vec}(\mathbf{L}(\boldsymbol{\omega})), \end{aligned}$$

which is linear with respect to each of  $\mathbf{L}(\boldsymbol{\omega})$ ,  $\mathbf{U}_1, \mathbf{U}_2$  and  $\mathcal{S}$ , fixing the other three components. This motivates us to update the four components cyclically by an alternating least squares method.

Since  $\mathbf{L}(\boldsymbol{\omega})$  is nonlinear in  $\boldsymbol{\omega}$ , let us consider the update of  $\boldsymbol{\omega}$  first. Recall that  $\mathbf{G}_k = \mathbf{U}_1 \mathcal{S}_k \mathbf{U}'_2$  for  $1 \leq k \leq d$ , where  $\mathcal{S}_k$  is the  $k$ th frontal slice of  $\mathcal{S}$ ; see (2.13). Let  $f^I(\tilde{\mathbf{x}}_t; \lambda_k) = \sum_{j=p+1}^{t-1} \lambda_k^{j-p} \mathbf{G}_{p+k} \mathbf{y}_{t-j}$  for  $1 \leq k \leq r$ , and

$$f^{II}(\tilde{\mathbf{x}}_t; \boldsymbol{\eta}_k) = \sum_{j=p+1}^{t-1} \gamma_k^{j-p} [\cos\{(j-p)\theta_k\} \mathbf{G}_{p+r+2k-1} + \sin\{(j-p)\theta_k\} \mathbf{G}_{p+r+2k}] \mathbf{y}_{t-j}$$

for  $1 \leq k \leq s$ , where we suppress the dependence on  $\mathbf{G}_{p+1}, \dots, \mathbf{G}_d$  since they will be fixed when updating  $\boldsymbol{\omega}$ . It can be verified that

$$\mathcal{A}_{(1)} \tilde{\mathbf{x}}_t - \sum_{k=1}^p \mathbf{G}_k \mathbf{y}_{t-k} = \sum_{k=1}^r f^I(\tilde{\mathbf{x}}_t; \lambda_k) + \sum_{k=1}^s f^{II}(\tilde{\mathbf{x}}_t; \boldsymbol{\eta}_k),$$

where each  $\lambda_k$  or  $\boldsymbol{\eta}_k$  appears in one summand only. As a result, the loss function will be block-separable with respect to  $\lambda_k$ 's and  $\boldsymbol{\eta}_k$ 's, when  $\mathcal{S}, \mathbf{U}_1$  and  $\mathbf{U}_2$  are fixed. Thus, at each

---

**Algorithm 1:** Alternating least squares algorithm

---

```

1 Input: ranks  $(\mathcal{R}_1, \mathcal{R}_2)$ , model orders  $(p, r, s)$ , initialization  $\omega^{(0)}, U_1^{(0)}, U_2^{(0)}, \mathcal{S}^{(0)}$  and  $\mathcal{G}^{(0)}$ .
2 repeat  $i = 0, 1, 2, \dots$ 
3   for  $k = 1, \dots, r$ :
4      $\lambda_k^{(i+1)} \leftarrow \arg \min_{\lambda \in (-1, 1)} \tilde{\mathbb{L}}_T(\lambda_1^{(i+1)}, \dots, \lambda_{k-1}^{(i+1)}, \lambda, \lambda_{k+1}^{(i)}, \dots, \lambda_r^{(i)}, \eta^{(i)}, \mathcal{G}^{(i)})$ 
5   for  $k = 1, \dots, s$ :
6      $\eta_k^{(i+1)} \leftarrow \arg \min_{\eta \in [0, 1] \times (-\pi/2, \pi/2)} \tilde{\mathbb{L}}_T(\lambda^{(i+1)}, \eta_1^{(i+1)}, \dots, \eta_{k-1}^{(i+1)}, \eta, \eta_{k+1}^{(i)}, \dots, \eta_s^{(i)}, \mathcal{G}^{(i)})$ 
7      $U_1^{(i+1)} \leftarrow \arg \min_{U_1} \sum_{t=1}^T \|y_t - [\{z'_t(\omega^{(i+1)})(I_d \otimes U_2^{(i)} \mathcal{S}_{(1)}^{(i)})' \} \otimes I_N] \text{vec}(U_1)\|_2^2$ 
8      $U_2^{(i+1)} \leftarrow \arg \min_{U_2} \sum_{t=1}^T \|y_t - U_1^{(i+1)} \mathcal{S}_{(1)}^{(i)} \{Z'_t(\omega^{(i+1)}) \otimes I_{R_2}\} \text{vec}(U_2')\|_2^2$ 
9      $\mathcal{S}^{(i+1)} \leftarrow \arg \min_{\mathcal{S}} \sum_{t=1}^T \|y_t - [\{z'_t(\omega^{(i+1)})(I_d \otimes U_2^{(i+1)})\} \otimes U_1^{(i+1)}] \text{vec}(\mathcal{S}_{(1)})\|_2^2$ 
10     $\mathcal{G}^{(i+1)} = \mathcal{S}^{(i+1)} \times_1 U_1^{(i+1)} \times_2 U_2^{(i+1)}$ .
11 until convergence

```

---

$\lambda_k$ - or  $\eta_k$ -update step, the irrelevant summands will be treated as the intercept and absorbed into the response. The corresponding optimization subproblem will be only one- or two-dimensional, which can be solved by the Newton-Raphson method. This feature will significantly alleviate the computational burden, since an increase in the orders  $r$  and  $s$  of the fitted model will have little impact on the cost of  $\omega$ -updates; see Algorithm 1 and the numerical evidence in Section 4.2.

The alternating updates of  $U_1, U_2$  and  $\mathcal{S}$  are straightforward. For lines 7–9 in Algorithm 1, we may adopt the closed-form solutions to the corresponding linear least squares problems. Alternatively, when the dimension  $N$  is large and the computation of closed-form solutions is time-consuming, the gradient descent method can be used to further speed up the algorithm. Note that Algorithm 1 can also be applied to the SARMA model without any low-Tucker-rank constraint on  $\mathcal{G}$ . In this case, lines 7, 8 and 10 in the algorithm will be dropped, and line 9 will be the update of  $\mathcal{G}^{(i+1)}$ , where both  $U_1$  and  $U_2$  in the loss function are set to the  $N \times N$  identity matrix.

Algorithm 1 requires predetermined ranks  $(\mathcal{R}_1, \mathcal{R}_2)$  and model orders  $(p, r, s)$ . In practice, we can adopt the data-driven selection procedure to be introduced in the next subsection prior to the parameter estimation. The initialization for Algorithm 1 can be conducted as follows:

We first obtain the initial value  $\mathcal{A}^{(0)}$  for  $\mathcal{A}$  by the spectral method in [19] or simply set it to  $\hat{\mathcal{A}}^{\text{init}}$ , the initial estimator used for the rank selection; see Remark 1 in Section 3.3. To obtain initial values for  $U_1$  and  $U_2$ , we conduct the HOSVD of  $\mathcal{A}^{(0)}$  for the first two modes,  $\mathcal{A}^{(0)} = \mathcal{H}^{(0)} \times_1 U_1^{(0)} \times_2 U_2^{(0)}$ , where  $\mathcal{H}^{(0)} \in \mathbb{R}^{\mathcal{R}_1 \times \mathcal{R}_2 \times \infty}$ . Then we randomly generate the initial value  $\omega^{(0)}$  from the uniform distribution on the parameter space  $\Omega$ . To get  $\mathcal{S}^{(0)}$ , note that  $\mathcal{A} = \mathcal{H} \times_1 U_1 \times_2 U_2$ , where  $\mathcal{H} = \mathcal{S} \times_3 L(\omega)$ . Since it is infeasible to recover  $\mathcal{S}^{(0)}$  exactly from  $\mathcal{H}^{(0)}$  and  $L(\omega^{(0)})$  based on the latter equation, with a little sacrifice of information, we can first replace  $\mathcal{H}^{(0)}$  with its projection,  $\mathcal{H}^{(0)} \times_3 [L(\omega^{(0)})\{L'(\omega^{(0)})L(\omega^{(0)})\}^{-1}L'(\omega^{(0)})]$ . This allows us to get the feasible solution,  $\mathcal{S}^{(0)} = \mathcal{H}^{(0)} \times_3 [\{L'(\omega^{(0)})L(\omega^{(0)})\}^{-1}L'(\omega^{(0)})]$ . Finally, we set  $\mathcal{G}^{(0)} = \mathcal{S}^{(0)} \times_1 U_1^{(0)} \times_2 U_2^{(0)}$ .

**3.3. Selection of ranks and model orders.** While it is assumed in previous subsections that the Tucker ranks and model orders are correctly specified, they are unknown in practice. In the following, we develop a two-stage procedure to select them.

Denote the true values of the ranks and orders by  $\mathcal{M}^* := (\mathcal{R}_1^*, \mathcal{R}_2^*, p^*, r^*, s^*)$ . Suppose that  $\hat{\mathcal{A}}^{\text{init}}$  is a consistent initial estimator of  $\mathcal{A}^*$ . Since the Tucker ranks are the same for  $\mathcal{A}^*$  and  $\mathcal{G}^*$ , at the first stage, we can select the ranks by the ridge-type ratio estimator [40]:

$$\hat{\mathcal{R}}_i = \arg \min_{1 \leq j \leq N-1} \frac{\sigma_{j+1}(\hat{\mathcal{A}}_{(i)}^{\text{init}}) + \tau}{\sigma_j(\hat{\mathcal{A}}_{(i)}^{\text{init}}) + \tau}, \quad i = 1, 2,$$

where  $\tau$  is a parameter to be chosen such that Assumption 4 below is satisfied; see similar methods in [40] and [38]. Let

$$\zeta_i = \frac{1}{\sigma_{\mathcal{R}_i^*}(\mathcal{A}_{(i)}^*)} \max_{1 \leq j \leq \mathcal{R}_i^* - 1} \frac{\sigma_j(\mathcal{A}_{(i)}^*)}{\sigma_{j+1}(\mathcal{A}_{(i)}^*)}, \quad i = 1, 2.$$

**ASSUMPTION 4.** The parameter  $\tau > 0$  is specified such that (i)  $\|\hat{\mathcal{A}}^{\text{init}} - \mathcal{A}^*\|_F / \tau = o_p(1)$  and (ii)  $\tau \max\{\zeta_1, \zeta_2\} = o(1)$ .

In Assumption 4, Condition (i) requires that the estimation error of  $\hat{\mathcal{A}}^{\text{init}}$  is dominated by  $\tau$ , and Condition (ii) can be regarded as the minimal signal assumption which will simply reduce to  $\tau = o(1)$  if  $\sigma_j(\mathcal{A}_{(i)}^*)$  for  $1 \leq j \leq \mathcal{R}_i^*$  and  $i = 1, 2$  are bounded above and away from zero by some absolute constants; see Remark 12 in [38] for more details.

**REMARK 1.** The initial estimator can be obtained based on the vector AR( $P$ ) approximation of the vector AR( $\infty$ ) model, where  $P$  grows with  $T$  [26]. Let  $\mathcal{A}_{1:P} = \text{stack}\{\mathcal{A}_j, 1 \leq j \leq$

$P\}$ . Then we can employ the regularized estimation,  $\hat{\mathcal{A}}_{1:P}^{\text{init}} = \arg \min_{\mathcal{A} \in \mathbb{R}^{N \times N \times P}} \sum_{t=P+1}^T \|\mathbf{y}_t - \sum_{j=1}^P \mathbf{A}_j \mathbf{y}_{t-j}\|_2^2 / (T - P) + \lambda_{\text{nuc}} \sum_{i=1}^2 \|\mathcal{A}_{(i)}\|_*$ , where  $\lambda_{\text{nuc}} > 0$ , and the nuclear norm regularizer enforces the low-rankness of  $\mathcal{A}_{(i)}$  for  $i = 1, 2$ ; see, e.g., [15] and [31]. Let  $\hat{\mathcal{A}}^{\text{init}} = \text{stack}\{\hat{\mathcal{A}}_{1:P}^{\text{init}}, \mathbf{0}_{N \times N}, \mathbf{0}_{N \times N}, \dots\}$ . Along the lines of the proofs of Theorem 2 in [37] and Proposition 4.2 in [39], under some regularity conditions, it can be shown that the approximation error due to the truncation after lag  $P$  will be negligible if  $P \asymp T^{1/2-\epsilon}$ , where  $\epsilon \in (0, 1/2)$ , and then  $\|\hat{\mathcal{A}}^{\text{init}} - \mathcal{A}^*\|_F = O_p\{\sqrt{(\mathcal{R}_1^* + \mathcal{R}_2^*)NP/(T-P)}\}$ . In our simulation studies, we set  $\tau = \sqrt{NP \log(T-P)/10(T-P)}$  as  $\mathcal{R}_1^*$  and  $\mathcal{R}_2^*$  are small.

Fixing  $(\mathcal{R}_1, \mathcal{R}_2) = (\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2)$ , at the second stage, we obtain  $(\hat{p}, \hat{r}, \hat{s})$  by minimizing the high-dimensional Bayesian information criterion (BIC),

$$\text{BIC}(p, r, s; \mathcal{R}_1, \mathcal{R}_2) = \log \left\{ \frac{1}{T} \sum_{t=1}^T \left\| \mathbf{y}_t - \sum_{j=1}^{t-1} \mathbf{A}_j(\check{\omega}, \check{\mathcal{G}}) \mathbf{y}_{t-j} \right\|_2^2 \right\} + \frac{cd_{\mathcal{M}} \log T}{T},$$

where  $(p, r, s)$  is searched over the range  $0 \leq p \leq p_{\max}$ ,  $0 \leq r \leq r_{\max}$ , and  $0 \leq s \leq s_{\max}$ , for some predetermined upper bounds, and  $c > 0$  is an arbitrary constant. Here  $\check{\omega}$  and  $\check{\mathcal{G}}$  represent the estimates obtained from (3.1) by fitting the model with  $\mathcal{M} = (\mathcal{R}_1, \mathcal{R}_2, p, r, s)$ .

Let  $\mathcal{M} = \{\mathcal{M} = (\mathcal{R}_1^*, \mathcal{R}_2^*, p, r, s) \mid 0 \leq p \leq p_{\max}, 0 \leq r \leq r_{\max}, 0 \leq s \leq s_{\max}\}$ . For any  $\mathcal{M} \in \mathcal{M}$ , let  $\Omega_{\mathcal{M}} \subset (-1, 1)^r \times \Pi^s$  be the parameter space of  $\omega$ , and  $\Gamma_{\mathcal{M}} = \{\mathcal{G} \in \mathbb{R}^{N \times N \times d} \mid \text{rank}(\mathcal{G}_{(1)}) \leq \mathcal{R}_1^*, \text{rank}(\mathcal{G}_{(2)}) \leq \mathcal{R}_2^*\}$ , where  $d = p + r + 2s$ . Then define the population minimizer for (3.1) corresponding to  $\mathcal{M}$  by

$$(\check{\omega}, \check{\mathcal{G}}) = \arg \min_{\omega \in \Omega_{\mathcal{M}}, \mathcal{G} \in \Gamma_{\mathcal{M}}} \mathbb{E} \left\{ \|\mathbf{y}_t - \mathcal{A}_{(1)} \mathbf{x}_t\|_2^2 \mid \mathcal{F}_{t-1} \right\},$$

where  $\mathcal{F}_t$  denotes the  $\sigma$ -field generated by  $\{\varepsilon_s\}_{s \leq t}$ . Note that  $\mathcal{M}^* \in \mathcal{M}$ , and when  $\mathcal{M} = \mathcal{M}^*$ , we simply have  $(\check{\omega}, \check{\mathcal{G}}) = (\omega^*, \mathcal{G}^*)$  and  $(\check{\omega}, \check{\mathcal{G}}) = (\hat{\omega}, \hat{\mathcal{G}})$ . In addition, let  $\mathcal{M}_{\text{mis}} = \{\mathcal{M} \in \mathcal{M} \mid p \neq p^*\} \cup \{\mathcal{M} \in \mathcal{M} \mid r < r^* \text{ or } s < s^*\}$  denote the subset of misspecified models.

**ASSUMPTION 5.** For every  $\mathcal{M} \in \mathcal{M}$ , the following conditions hold: (i)  $|\lambda_1|, \dots, |\lambda_r|, \gamma_1, \dots, \gamma_s \in \Lambda$  for all  $\omega \in \Omega_{\mathcal{M}}$ , where  $\Lambda$  is the compact subset of  $(0, \bar{\rho})$  defined as in Assumption 1; (ii)  $\|\check{\mathcal{G}}_k\|_{\text{op}} \leq C_g$  for  $1 \leq k \leq d$ , and  $c_g \sqrt{\mathcal{R}_1^* \wedge \mathcal{R}_2^*} \leq \|\check{\mathcal{G}}_k\|_F \leq C_g \sqrt{\mathcal{R}_1^* \wedge \mathcal{R}_2^*}$  for  $p+1 \leq k \leq d$ , where  $0 < c_g < C_g$  are absolute constants defined as in Assumption 2; (iii)  $\sigma_{\min}(\mathbf{L}_{\text{stack}}(\omega)) \asymp 1$  and  $\sigma_{\max}(\mathbf{L}_{\text{stack}}(\omega)) \asymp 1$  for all  $\omega \in \Omega_{\mathcal{M}}$ , where  $\mathbf{L}_{\text{stack}}(\omega)$  is defined as in (A.1) in the Appendix; and (iv)  $\|\check{\omega} - \hat{\omega}\|_2 \leq c_{\check{\omega}}$ , where  $c_{\check{\omega}} = \min\{2, c_g(1 - \bar{\rho})\sigma_{\min}(\mathbf{L}_{\text{stack}}(\check{\omega})) / (8\sqrt{2}C_L C_g)\}$  is an absolute constant.

ASSUMPTION 6. For all  $\mathcal{M} \in \mathcal{M}_{\text{mis}}$ ,  $\mathbb{E}\{\|(\mathcal{A}^* - \hat{\mathcal{A}})_{(1)} \mathbf{x}_t\|_2^2\} \geq C_{\text{mis}} N$ , where  $\hat{\mathcal{A}} = \hat{\mathcal{G}} \times_3 \mathbf{L}(\hat{\omega})$ , and  $C_{\text{mis}} > 0$  is a sufficiently large absolute constant.

Assumption 5 is analogous to Assumptions 1 and 2 for all  $\mathcal{M} \in \mathcal{M}$ . Assumption 6 guarantees sufficient departure of the misspecified models from the true model.

THEOREM 3.3. Under Assumption 4,  $\mathbb{P}(\hat{\mathcal{R}}_1 = \mathcal{R}_1^*, \hat{\mathcal{R}}_2 = \mathcal{R}_2^*) \rightarrow 1$  as  $T \rightarrow \infty$ . Moreover, under Assumptions 1–6, if there exist absolute constants  $c, C > 0$  such that  $c \leq \lambda_{\min}(\Sigma_\varepsilon) \leq \lambda_{\max}(\Sigma_\varepsilon) \leq C$  and  $c \leq \mu_{\min}(\Psi_*) \leq \mu_{\max}(\Psi_*) \leq C$ , and  $d_{\mathcal{M}, \max}/T \rightarrow 0$  as  $T \rightarrow \infty$ , then  $\mathbb{P}(\hat{\mathcal{M}} = \mathcal{M}^*) \rightarrow 1$  as  $T \rightarrow \infty$ , where  $d_{\mathcal{M}, \max} = \mathcal{R}_1^* \mathcal{R}_2^* (p_{\max} + r_{\max} + 2s_{\max}) + (\mathcal{R}_1^* + \mathcal{R}_2^*)N$ , and  $\hat{\mathcal{M}} = (\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2, \hat{p}, \hat{r}, \hat{s})$ .

#### 4. Simulation studies.

4.1. *Estimation error rate and model selection consistency.* Since the proposed model is intended to serve as a tractable alternative to the vector ARMA model, we evaluate the performance of our procedure by considering data generated from vector ARMA processes. Specifically, we consider two data generating processes (DGPs),

- DGP1: the vector MA(1) model  $\mathbf{y}_t = \varepsilon_t - \Theta \varepsilon_{t-1}$ , and
- DGP2: the vector ARMA(1, 1) model  $\mathbf{y}_t = \Phi \mathbf{y}_{t-1} + \varepsilon_t - \Theta \varepsilon_{t-1}$ ,

which correspond to  $p = 0$  and 1, respectively. For both DGPs,  $\{\varepsilon_t\}$  are *i.i.d.*  $N(\mathbf{0}, \mathbf{I}_N)$ , and we set  $\Theta = \mathbf{B} \mathbf{J} \mathbf{B}^{-1}$ , where  $\mathbf{B} \in \mathbb{R}^{N \times N}$  is a randomly generated orthogonal matrix,  $\mathbf{J} = \text{diag}\{\lambda_1, \dots, \lambda_r, \mathbf{C}_1, \dots, \mathbf{C}_s, \mathbf{0}\}$  is the real Jordan normal form, with each  $\mathbf{C}_k$  being the  $2 \times 2$  block determined by  $\boldsymbol{\eta}_k = (\gamma_k, \theta_k)'$ , defined as in (2.2). For DGP2, we set  $\Phi = \mathbf{B} \mathbf{K} \mathbf{B}^{-1}$ , where  $\mathbf{K} = \text{diag}\{\delta, 0, \dots, 0\}$ , with the entry  $\delta \neq 0$ . Note that both DGPs can be written as the proposed SARMA( $p, r, s$ ) model, where  $\boldsymbol{\omega} = (\lambda_1, \dots, \lambda_r, \boldsymbol{\eta}'_1, \dots, \boldsymbol{\eta}'_s)'$  corresponds to the parameters in  $\mathbf{J}$ . In addition,  $\mathcal{R}_1 = \mathcal{R}_2 = r + 2s$  for both DGP1 and DGP2.

In the first experiment, we verify the rate of  $\sqrt{d_{\mathcal{M}}/T}$  for the estimation errors  $\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F$ ,  $\|\hat{\omega} - \omega^*\|_2$ , and  $\|\hat{\mathcal{G}} - \mathcal{G}^*\|_F$  derived by combining Theorem 3.2 and Proposition 3.1. We set  $(r, s) = (1, 0)$  and  $\lambda_1 = -0.7$  for both DGPs, and  $\delta = 0.5$  for DGP2. We consider  $N = 10, 20$  or 40, and  $T$  is chosen such that  $d_{\mathcal{M}}/T \in \{0.05, 0.1, 0.15, 0.2, 0.25\}$ . The estimation is conducted using the algorithm in Section 3.2 given the true ranks and model orders, where

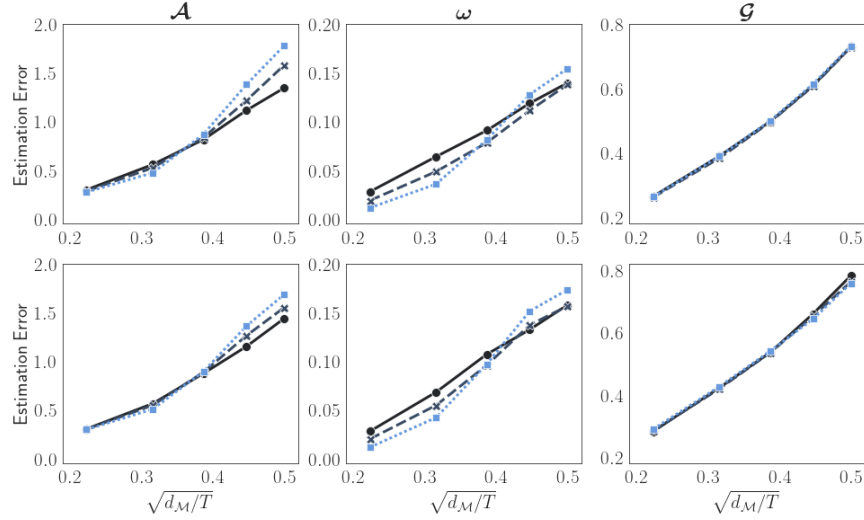


FIGURE 2. Plots of estimation errors  $\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F$  (left panel),  $\|\hat{\omega} - \omega^*\|_2$  (middle panel) and  $\|\hat{\mathcal{G}} - \mathcal{G}^*\|_F$  (right panel) against  $\sqrt{d_{\mathcal{M}}/T}$ , where  $(\mathcal{R}_1, \mathcal{R}_2, p, r, s) = (1, 1, 0, 1, 0)$  (top panel) or  $(\mathcal{R}_1, \mathcal{R}_2, p, r, s) = (1, 1, 1, 1, 0)$  (bottom panel), and  $N = 10$  (—●—),  $20$  (---×---) or  $40$  (····■···).

$U_1, U_2$  and  $\mathcal{S}$  are updated by the alternating least squares algorithm using their corresponding closed-form solutions. Figure 2 plots the estimation errors averaged over 500 replications against  $\sqrt{d_{\mathcal{M}}/T}$ . An approximately linear relationship between the estimation error and the theoretical rate can be observed across all settings, which confirms our theoretical results.

The aim of the second experiment is to verify the consistency of the rank and model order selection procedure. We consider three cases under DGP1:  $(\mathcal{R}_1, \mathcal{R}_2, r, s) = (1, 1, 1, 0)$  (model A),  $(2, 2, 0, 1)$  (model B), and  $(3, 3, 1, 1)$  (model C). The results for DGP2 are similar and hence are omitted for brevity. For models B and C, we set  $\theta_1 = \pi/4$ . Note that  $\mathcal{A}_{(1)}$  and  $\mathcal{A}_{(2)}$  have the same singular values under DGP1. Moreover, when  $r \leq 1$  and  $s \leq 1$ , the magnitude of the nonzero singular values are directly determined by  $|\lambda_1|$  and/or  $\gamma_1$ , which control the signal strength for the rank selection. We consider four levels of signal strength  $\{0.6, 0.65, 0.7, 0.75\}$ , and set  $-\lambda_1$  in model A,  $\gamma_1$  in model B, and  $-\lambda_1 = \gamma_1$  in model C to these values. In addition, we consider  $N = 10$  and  $T \in [100, 600]$ . At the first stage, the initial estimator  $\hat{\mathcal{A}}^{\text{init}}$  is obtained by the method in Remark 1, where  $P = T^{1/3}$ , and  $\lambda_{\text{nuc}}$  is chosen by the cross validation method for time series similar to that in [39]. At the second stage, we minimize the BIC with  $c = 0.1$ ,  $p_{\max} = 1$ , and  $r_{\max} = s_{\max} = 2$ . Figure 3 presents the proportion of correct rank selection,  $\{(\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2) = (\mathcal{R}_1^*, \mathcal{R}_2^*)\}$ , and that of correct rank and model order selection,  $\{\hat{\mathcal{M}} = \mathcal{M}^*\}$ , i.e.,  $\{(\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2, \hat{p}, \hat{r}, \hat{s}) = (\mathcal{R}_1^*, \mathcal{R}_2^*, p^*, r^*, s^*)\}$ , based on the two-stage procedure. It can be clearly seen that both proportions increase to one



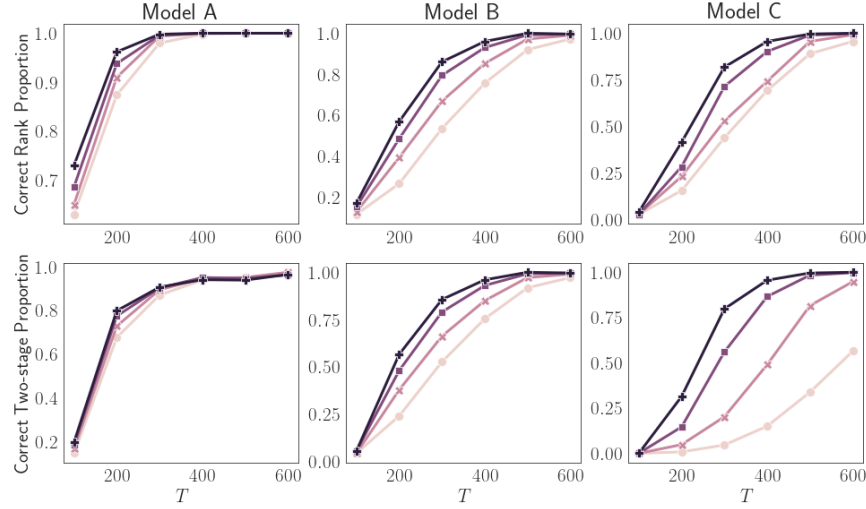


FIGURE 3. Proportion of correct rank selection (top panel) and two-stage selection (bottom panel) for models A (left panel), B (middle panel) and C (right panel), where the signal strength is 0.6 (—○—), 0.65 (—×—), 0.7 (—■—) or 0.75 (—+—).

as  $T$  and the signal strength increase. For models A and B, the proportion that the ranks and model orders are correctly selected simultaneously is fairly close to one when  $T \geq 500$  across all settings. For model C, a larger  $T$  and signal strength is required to achieve similar performance since the model is more complex. Nevertheless, although not reported in the figure, the proportion of correct selection of  $\mathcal{M}^*$  for model C under signal strength 0.6 will be close to one when  $T \geq 900$ .

**4.2. Computational advantage over vector ARMA models.** We conduct two experiments to confirm the computational advantage of the proposed SARMA model over the vector ARMA model. The computation time is reported for a Linux server with Intel(R) Xeon(R) Gold 6246 CPU@3.30GHz and 503GB of RAM, and the algorithms are all implemented by Python.

For a fair comparison, we first consider the case without any low-Tucker-rank structure for  $\mathcal{G}$ , i.e.,  $\mathcal{R}_1 = \mathcal{R}_2 = N$ . The data are generated from DGP2, where  $\Theta$  is defined as in Section 4.1 with  $(r, s) = (1, 1)$ ,  $\lambda_1 = -0.8$ ,  $\gamma_1 = 0.8$ , and  $\theta_1 = \pi/4$ . In addition,  $\Phi = BKB^{-1}$ , where  $K = \text{diag}\{0.8, -0.8, -0.8, 0.8, \dots, 0.8\}$ . We fit both the proposed model and the vector ARMA model to the data assuming that the true model orders are known. The estimation of the vector ARMA model is implemented by the Python function `statsmodels.tsa.statespace.varmax`, which conducts the maximum likelihood

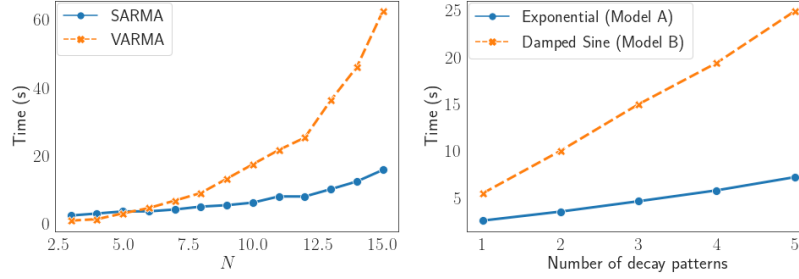


FIGURE 4. Computation time in seconds against  $N$  for two methods (left panel) and that against the number of decay patterns, i.e.,  $r$  or  $s$ , for the proposed method (right panel).

estimation (MLE) via the Kalman filter [26]. This implementation of the MLE employs the limited-memory Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm by default, which is generally faster than the Newton-Raphson method. We consider  $T = 500$  and  $N \in [3, 15]$ , since the MLE of the vector ARMA model will become extremely slow for larger  $N$ . We initialize both algorithms with the true parameter values and record the computation time of 20 iterations, averaged over 500 replications. The left panel of Figure 4 plots the average computation time against  $N$ . It can be observed that the computation time for the MLE increases much more rapidly than that of the proposed method as  $N$  increases. As discussed in previous sections, unlike the vector ARMA model, the tractable parametric structure of the proposed model allows us to avoid the computation of high-order  $N \times N$  matrix polynomials. This results in the significant improvement in the computational speed.

In the next experiment, we verify the computational scalability of the proposed method with respect to increases in  $r$  and  $s$ . Similarly to the second experiment in Section 4.1, we consider two cases under DGPI:  $(\mathcal{R}_1, \mathcal{R}_2, r, s) = (r, r, r, 0)$  (model A), and  $(2s, 2s, 0, s)$  (model B), where  $r, s \in \{1, \dots, 5\}$ . For model A, we set  $-\lambda_1 = 0.8$ ,  $-\lambda_2 = \lambda_5 = 0.5$ ,  $-\lambda_3 = \lambda_4 = 0.2$ . For model B, we set  $\gamma_k = 0.1 + 0.15(k - 1)$  and  $\theta_k = (-1)^{k-1}\pi/4$  for  $1 \leq k \leq 5$ . We fix  $T = 500$  and  $N = 40$ , and run 5 iterations of the proposed algorithm with rank constraints for 500 replications. The right panel of Figure 4 displays the average computation time against  $r$  and  $s$  for models A and B, respectively. It shows that an increase in  $s$  is more time-consuming than that in  $r$ , which is expected since an extra  $\eta_k$  involves two parameters, while another  $\lambda_k$  involves only one. Nevertheless, it can be observed that the computation time scales linearly in both  $r$  and  $s$ . On the other hand, since the MLE is too slow under large  $N$ , we do not report its computation time in the figure. Indeed, in most

cases, it takes more than 17 minutes for the MLE to finish 5 iterations, which is about 40 times the computation time of the proposed method.

## 5. Two empirical examples.

**5.1. Macroeconomic dataset.** This dataset contains observations of 20 quarterly macroeconomic variables from June 1959 to December 2019, with  $T = 243$ , retrieved from FRED-QD [27]. These variables come from four categories: (i) interest rates, (ii) money and credit, (iii) exchange rates, and (iv) stock market. These categories are usually considered in the construction of financial condition index, since they reflect important factors that can affect the stance of monetary policy and aggregate demand conditions [16, 9, 21]. All series are transformed to be stationary, and standardized to have zero mean and unit variance; see Table S.1 in the supplementary file for more details of the variables and their transformations.

We first explore the factor structures of this dataset. As discussed in Section 2.4,  $U_1$  and  $U_2$  capture response factor and predictor factor spaces, respectively. By the rank selection method in Section 3.3, we obtain  $(\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2) = (3, 3)$ . Note that  $U_1$  can be alternatively estimated by  $\hat{\Lambda}$ , up to an orthogonal rotation, by fitting the factor model  $\mathbf{y}_t = \Lambda \mathbf{f}_t + \varepsilon_t$  via the method in [24], where  $\mathbf{f}_t$  is  $\mathcal{R}_1$ -dimensional; see the discussion below (2.15) in Section 2.4. To verify this, we obtain the estimate  $\hat{\Lambda}$  with  $\mathcal{R}_1 = 3$ . Figure 5 displays  $\hat{\Lambda}$  based on the factor model, together with  $\hat{U}_1$  and  $\hat{U}_2$  based on the proposed method in Section 3.1; note that all three matrices are orthonormal. Overall, it can be observed that  $\hat{\Lambda}$  and  $\hat{U}_1$  are mainly influenced by variables in categories (i) and (ii), while the influence from categories (iii) and (iv) is relatively weak. On the other hand, a couple of variables from category (iv) contribute significantly to the first factor of  $\hat{U}_2$ . Moreover, the pattern of  $\hat{\Lambda}$  is much similar to  $\hat{U}_1$  than  $\hat{U}_2$ . Although  $\hat{U}_1$  and  $\hat{\Lambda}$  do not match exactly due to the rotation and estimation error, their subspace distance  $\|\hat{U}_1 \hat{U}_1' - \hat{\Lambda} \hat{\Lambda}'\|_F^2$  is as small as 0.52, whereas that between  $\hat{U}_1$  and  $\hat{U}_2$  is much larger, with  $\|\hat{U}_1 \hat{U}_1' - \hat{U}_2 \hat{U}_2'\|_F^2 = 3.16$ . This confirms empirically the equivalence between  $U_1$  and  $\Lambda$ .

We evaluate the performance of our method based on out-of-sample forecast accuracy. The following rolling forecast procedure is adopted: we first fit the models using historical data with the end point rolling from the fourth quarter of 2015 to the third quarter of 2019, and then conduct one-step-ahead forecasts based on the fitted models. In addition to the proposed

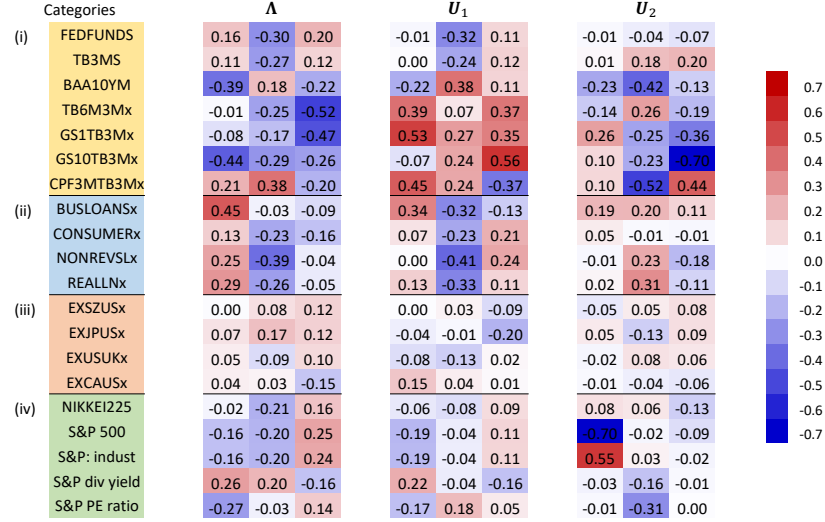


FIGURE 5. Estimates of factor loadings in the factor model ( $\Lambda$ ) and the proposed model ( $U_1$  and  $U_2$ ) for the macroeconomic dataset.

TABLE 1  
Forecast errors for macroeconomic and realized volatility datasets. The smallest numbers in each row are marked in bold.

		Vector AR			Vector ARMA		SARMA
		Lasso	MLR	SHORR	$\ell_1$	HLag	
Macroeconomic	MSFE	2.78	2.77	2.71	2.80	2.79	<b>2.67</b>
	MAFE	9.26	9.27	8.99	9.28	9.24	<b>8.75</b>
Realized Volatility	MSFE	5.17	4.93	4.87	5.19	5.19	<b>4.78</b>
	MAFE	21.58	19.02	18.22	21.70	21.70	<b>16.45</b>

method, we consider five other existing methods, including three based on the vector AR model and two based on the vector ARMA model. Specifically, for the vector AR model, we consider (a) the Lasso method [8] and two methods in [38]: (b) the multilinear low-rank (MLR) method and (c) the sparse higher-order reduced-rank (SHORR) method, which amounts to further imposing sparsity on the factor matrices in (b). For the vector ARMA model, we apply the method in [39] with (d) the  $\ell_1$ -penalty or (e) the HLag penalty. Note that (a) is used as the Phase-I estimator for the estimators in (d) and (e), and the AR order is selected according to [39]. The AR order for (b) and (c) is chosen as in [38]. For the proposed model, the estimated model orders are  $(\hat{p}, \hat{r}, \hat{s}) = (0, 1, 0)$ . Throughout the rolling forecast procedure, the same model orders and ranks are used. Table 1 reports the mean squared forecast error (MSFE) and mean absolute forecast error (MAFE) for all methods. It can be observed that the proposed method achieves the smallest forecast errors among all competing ones. Compared to sparsity-inducing estimation methods such as (a), (d) and

(e), the proposed model can better capture the factor structure which is prominent in this dataset. Meanwhile, the proposed model is more flexible than the finite-order vector AR model, resulting in better forecasting performance than (b) and (c).

As mentioned in Section 2.4, when  $U_1 = U_2$ , the proposed model may also be estimated by a two-step method, where the SARMA model is fitted to the estimated  $\mathcal{R}$ -dimensional factor series  $\hat{\mathbf{f}}_t$  after fitting the factor model, where  $\mathcal{R} = \mathcal{R}_1 = \mathcal{R}_2$ . We have performed the above rolling forecast procedure using this method and found that  $\text{MSFE} = 2.80$  and  $\text{MAFE} = 9.36$ , where  $\mathcal{R} = 3$  is adopted since  $(\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2) = (3, 3)$ . The worse forecasting performance of this method suggests that the predictor and response subspaces may be quite different.

*5.2. Realized volatility.* To further demonstrate the forecasting performance of the proposed method, we study daily realized volatilities for 46 stocks from January 2, 2012 to December 31, 2013, with  $T = 495$ . These are the stocks of top S&P 500 companies ranked by trading volumes on the first day of 2013. Specifically, we obtain the tick-by-tick data from WRDS (<https://wrds-www.wharton.upenn.edu>) and compute the daily realized volatility from five-minute returns [1]. By examining the sample autocorrelation functions, we have confirmed the stationarity of all series. Each series is then standardized to have zero mean and unit variance. More information about the stocks is given in Table S.2 in the supplementary file. We conduct the same rolling forecast procedure as in Section 5.1, where the last 10% of the sample is used as the forecast period. As shown in Table 1, the proposed method considerably outperforms the other ones in terms of the forecast accuracy.

The estimated ranks and model orders of the proposed model are  $(\hat{\mathcal{R}}_1, \hat{\mathcal{R}}_2, \hat{p}, \hat{r}, \hat{s}) = (1, 1, 0, 1, 0)$  for this dataset. Under this model, (2.14) can be simply written as the univariate regression,  $\mathbf{u}'_1 \mathbf{y}_t = \beta \sum_{j=1}^{\infty} \lambda^j \mathbf{u}'_2 \mathbf{y}_{t-j} + \mathbf{u}'_1 \varepsilon_t$ , where the one-dimensional factor series  $\{\mathbf{u}'_1 \mathbf{y}_t\}$  and  $\{\mathbf{u}'_2 \mathbf{y}_t\}$  may be regarded as volatilities of two different market indices, signifying how the market responds to and picks up risks across different stocks, respectively. In addition, it is worth noting that  $\hat{\lambda}$  for this dataset is close to 0.95. This lends further support to the well-established fact that the volatility of asset returns is highly persistent, that is, the AR process of the volatility is nearly unit-root; see, e.g., [2].

**6. Conclusion and discussion.** In this paper, we develop a statistically and computationally attractive alternative to the vector ARMA model. The proposed model mimics the vector ARMA temporal dynamic through a prespecified parametric function  $L(\cdot)$  that depicts the temporal decay pattern of the coefficient matrices  $A_j$ 's in the vector AR( $\infty$ ) form. As a result, the AR structure over infinitely many time lags can be characterized by an  $(r + 2s)$ -dimensional parameter vector  $\omega$ ; or using tensor notations,  $\mathcal{A} \in \mathbb{R}^{N \times N \times \infty}$  is the mode-3 product of the tensor  $\mathcal{G} \in \mathbb{R}^{N \times N \times d}$  and the matrix  $L(\omega) \in \mathbb{R}^{\infty \times d}$ . In the high-dimensional regime, we further impose a low-Tucker-rank structure on  $\mathcal{G}$ , leading to the factorization  $\mathcal{A} = \mathcal{S} \times_1 U_1 \times_2 U_2 \times_3 L(\omega)$ , where  $\mathcal{S} \in \mathbb{R}^{\mathcal{R}_1 \times \mathcal{R}_2 \times d}$ , and  $U_i \in \mathbb{R}^{N \times \mathcal{R}_i}$  for  $i = 1, 2$ . This enables us to embed the concept of factor modeling in the proposed predictive framework, where the orthonormalized matrices  $U_1$  and  $U_2$  can be interpreted as the loadings of the extracted response and predictor factors, respectively.

We discuss some extensions of the proposed method that are worthwhile to explore further. Firstly, when  $N$  is large, the estimated factor loadings may contain many nearly zero entries; see the results of  $\hat{U}_1$  and  $\hat{U}_2$  in Figure 5. To improve the estimation efficiency of the loading matrices, it is natural to further apply sparsity-inducing regularizations to  $U_1$  and  $U_2$ . Following [38], we may assume that these matrices are simultaneously orthonormal and sparse, and apply the vectorized  $\ell_1$ -norm regularization  $\|U_1 \otimes U_2\|_1$  together with orthonormality constraints on each  $U_i$ . Combined with the nonasymptotic techniques in this paper, it can be shown that the estimation error bound for the corresponding estimator of  $\mathcal{A}$  will be sharpened to  $O_p(\sqrt{K \log N/T})$ , where  $K$  denotes the number of nonzero entries in  $U_1 \otimes U_2$ . Thus, a much larger dimension  $N$  can be accommodated. Nevertheless, the non-smoothness and non-convexity of the corresponding optimization problem will pose challenges to the development of an efficient algorithm. Thus, we leave this for future research.

Another interesting direction is to consider other forms of  $L(\cdot)$ . A straightforward generalization is to adopt the more sophisticated decay patterns derived from the general vector ARMA model; see Proposition 2.2, Figure 6 and Section S4 of the supplementary material for details. To allow for even more flexibility in the temporal structure, we can further consider other user-defined decay patterns based on prior knowledge of the data. For instance, economic and financial time series often exhibit seasonality, i.e., similar fluctuations recurring at regular intervals. Denote by  $h$  the length of the regular intervals, a.k.a. the sea-

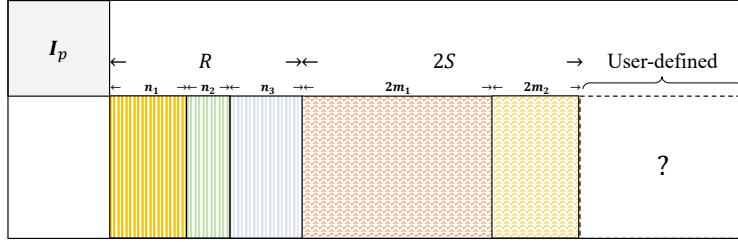


FIGURE 6. General form of  $\mathbf{L}(\omega)$ , where rectangles filled with vertical stripes are blocks in Type I parametric form (exponential decays), those filled with waves are blocks in Type II parametric form (damped sine and cosine waves), different colors indicate dependence on different subsets of parameters  $\lambda_k$  or  $\eta_k$ ; see the supplementary material for detailed expressions. The area with a question mark represents possible extra user-defined parametric blocks.

sonal period; e.g.,  $h = 4$  for quarterly data and 12 for monthly data. The seasonality can be induced by incorporating seasonal versions of  $\ell^I(\lambda_k)$  and  $\ell^{II}(\eta_k)$ . For example, we may define  $\ell_{\text{season}}^I(\lambda_k) = (\lambda_k, 0, \dots, 0, \lambda_k^2, 0, \dots)'$  and

$$\ell_{\text{season}}^{II}(\eta_k) = \begin{pmatrix} \gamma_k \cos(\theta_k) & 0 \cdots 0 & \gamma_k^2 \cos(2\theta_k) & 0 \cdots \\ \gamma_k \sin(\theta_k) & 0 \cdots 0 & \gamma_k^2 \sin(2\theta_k) & 0 \cdots \end{pmatrix}',$$

where the nonzero entries occur once every  $h$  rows in each column of  $\ell_{\text{season}}^I(\lambda_k)$  and  $\ell_{\text{season}}^{II}(\eta_k)$ . In general, both the regular decay patterns,  $\ell^I(\cdot)$  and  $\ell^{II}(\cdot)$ , and the seasonal decay patterns can be included in  $\mathbf{L}(\cdot)$ ; see also Figure 6. Hence, the corresponding seasonal SARMA model will be able to generate various kinds of patterns involving both seasonal and non-seasonal components.

## APPENDIX: TECHNICAL DETAILS

**A.1. Useful properties of  $\mathbf{L}(\omega)$ .** According to the definition of  $\mathbf{L}(\omega)$ , for  $j \geq 1$ , denote the  $j$ th entry of  $\ell^I(\lambda_i)$  by  $\ell_j^I(\lambda_i) = \lambda_i^j$  with  $1 \leq i \leq r$ , and denote the transpose of the  $j$ th row of the  $\infty \times 2$  matrix  $\ell^{II}(\eta_k)$  by  $\ell_j^{II}(\eta_k) := (\ell_j^{II,1}(\eta_k), \ell_j^{II,2}(\eta_k))' = (\gamma_k^j \cos(j\theta_k), \gamma_k^j \sin(j\theta_k))'$  with  $1 \leq k \leq s$ . Let  $\mathbf{L}^I(\boldsymbol{\lambda}) = (\ell^I(\lambda_1), \dots, \ell^I(\lambda_r))$  and  $\mathbf{L}^{II}(\boldsymbol{\eta}) = (\ell^{II}(\eta_1), \dots, \ell^{II}(\eta_s))$ . In addition, define the following matrix by augmenting  $\mathbf{L}(\omega)$  with  $(r + 2s)$  extra columns consisting of first-order derivatives:

$$(A.1) \quad \mathbf{L}_{\text{stack}}(\omega) = \begin{pmatrix} \mathbf{I}_p \\ \mathbf{L}^I(\boldsymbol{\lambda}) \quad \mathbf{L}^{II}(\boldsymbol{\eta}) \quad \nabla \mathbf{L}^I(\boldsymbol{\lambda}) \quad \nabla_{\theta} \mathbf{L}^{II}(\boldsymbol{\eta}) \end{pmatrix},$$

where  $\nabla \mathbf{L}^I(\boldsymbol{\lambda}) = (\nabla \ell^I(\lambda_1), \dots, \nabla \ell^I(\lambda_r))$  and  $\nabla_{\theta} \mathbf{L}^{II}(\boldsymbol{\eta}) = (\nabla_{\theta} \ell^{II}(\eta_1), \dots, \nabla_{\theta} \ell^{II}(\eta_s))$ .

We can similarly define  $\nabla_{\gamma} \mathbf{L}^{II}(\boldsymbol{\eta})$ . Note that from (A.12), it holds  $\text{colsp}\{\nabla_{\gamma} \mathbf{L}^{II}(\boldsymbol{\eta})\} =$

$\text{colsp}\{\nabla_{\theta} \mathbf{L}^{II}(\boldsymbol{\eta})\}$ , which is why  $\nabla_{\gamma} \mathbf{L}^{II}(\boldsymbol{\eta})$  is not included in  $\mathbf{L}_{\text{stack}}(\boldsymbol{\omega})$ . Denote

$$\sigma_{\min,L} = \sigma_{\min}(\mathbf{L}_{\text{stack}}(\boldsymbol{\omega}^*)) \quad \text{and} \quad \sigma_{\max,L} = \sigma_{\max}(\mathbf{L}_{\text{stack}}(\boldsymbol{\omega}^*)),$$

where  $\boldsymbol{\omega}^*$  is the true value of  $\boldsymbol{\omega}$ . Lemma A.1 below gives some exponential decay properties induced by the parametric form of  $\mathbf{L}(\boldsymbol{\omega})$ , which will be used repeatedly in our theoretical analysis. Then, based on Lemma A.1(ii), we can show that  $\sigma_{\min,L} \asymp 1$  and  $\sigma_{\max,L} \asymp 1$  if Assumption 2(i) is further imposed; this is stated in Lemma A.2.

LEMMA A.1. *Suppose that Assumption 1(i) holds. Then (i) there exists an absolute constant  $C_L > 0$  such that for all  $\boldsymbol{\omega} \in \boldsymbol{\Omega}$  and  $j \geq 1$ ,*

$$\max_{1 \leq i \leq r, 1 \leq k \leq s, 1 \leq h \leq 2} \{|\nabla \ell_j^I(\lambda_i)|, \|\nabla \ell_j^{II,h}(\boldsymbol{\eta}_k)\|_2, |\nabla^2 \ell_j^I(\lambda_i)|, \|\nabla^2 \ell_j^{II,h}(\boldsymbol{\eta}_k)\|_F\} \leq C_L \bar{\rho}^j;$$

*and (ii) there exists an absolute constant  $C_* > 0$  such that  $\|\mathbf{A}_j^*\|_{\text{op}} \leq C_* \bar{\rho}^j$  for all  $j \geq 1$  if Assumption 1(ii) further holds.*

LEMMA A.2. *Let  $J = 2(r + 2s)$ . Denote  $x_k^* = \lambda_k^*$  for  $1 \leq k \leq r$  and  $x_{r+2k-1}^* = \gamma_k^* e^{i\theta_k^*}$ ,  $x_{r+2k}^* = \gamma_k^* e^{-i\theta_k^*}$  for  $1 \leq k \leq s$ , and let  $\nu_1^* = \min\{|x_k^*|, 1 \leq k \leq r + 2s\}$  and  $\nu_2^* = \min\{|x_j^* - x_k^*|, 1 \leq j < k \leq r + 2s\}$ . Under Assumptions 1(i) and 2(i), the matrix  $\mathbf{L}_{\text{stack}}(\boldsymbol{\omega}^*)$  has full rank, and its maximum and minimum singular values satisfy*

$$\min\{1, c_{\bar{\rho}}\} \leq \sigma_{\min,L} \leq \sigma_{\max,L} \leq \max\{1, C_{\bar{\rho}}\}.$$

where  $C_{\bar{\rho}} = C_L \bar{\rho} \sqrt{J} (1 - \bar{\rho})^{-1} \asymp 1$  and  $c_{\bar{\rho}} = 0.25^s (\nu_1^*)^{3J/2} (\nu_2^*)^{J(J/2-1)} / C_{\bar{\rho}}^{J-1} \asymp 1$ .

By Lemma A.2, we have  $c_{\omega} \asymp 1$ ,  $c_{\Delta} \asymp 1$ , and  $C_{\Delta} \asymp 1$ , where

$$(A.2) \quad c_{\omega} = \min \left\{ 2, \frac{c_g(1 - \bar{\rho})\sigma_{\min,L}}{8\sqrt{2}C_L C_g} \right\},$$

$$c_{\Delta} = c_l \cdot \sigma_{\min,L}, \quad C_{\Delta} = c_u \cdot \max\{\sigma_{\max,L}, (1 - \bar{\rho})^{-1}\},$$

with  $c_l = 0.25 \min\{1, \sqrt{2}c_g\}$  and  $c_u = 1 + \sqrt{2}C_g (\min_{1 \leq k \leq s} \gamma_k^*)^{-1} + 4\sqrt{2}C_L + 2C_L C_g$ ; for details about how these constants are chosen, see the proof of Proposition 3.1.

**A.2. Notations and main idea: linearization of parametric structure.** For simplicity, denote the perturbations of  $\boldsymbol{\omega}^*$ ,  $\mathcal{G}^*$  and  $\mathcal{A}^*$  by  $\delta_{\omega} = \|\boldsymbol{\omega} - \boldsymbol{\omega}^*\|_2$ ,  $\delta_{\mathcal{G}} = \|\mathcal{G} - \mathcal{G}^*\|_F$  and  $\delta_{\mathcal{A}} = \|\mathcal{A} - \mathcal{A}^*\|_F = \|\Delta\|_F$ , respectively. Let

$$\Upsilon = \{\Delta = \mathcal{A} - \mathcal{A}^* \in \mathbb{R}^{N \times N \times \infty} \mid \mathcal{A} = \mathcal{G} \times_3 \mathbf{L}(\boldsymbol{\omega}), \mathcal{G} \in \Gamma(\mathcal{R}_1, \mathcal{R}_2), \boldsymbol{\omega} \in \boldsymbol{\Omega}, \delta_{\omega} \leq c_{\omega}\},$$



where  $\Gamma(\mathcal{R}_1, \mathcal{R}_2) = \{\mathcal{G} \in \mathbb{R}^{N \times N \times d} \mid \text{rank}(\mathcal{G}_{(1)}) \leq \mathcal{R}_1, \text{rank}(\mathcal{G}_{(2)}) \leq \mathcal{R}_2\}$ . It is noteworthy that under the conditions of Theorem 3.2,  $\hat{\Delta} := \hat{\mathcal{A}} - \mathcal{A}^* \in \Upsilon$ .

A crucial intermediate step for our theoretical analysis is to establish the following linear approximation within a fixed local neighborhood of  $\omega^*$ ,

$$(A.3) \quad \Delta(\omega, \mathcal{G}) = \mathcal{A}(\omega, \mathcal{G}) - \mathcal{A}^* \approx \mathcal{M}(\omega - \omega^*, \mathcal{G} - \mathcal{G}^*) \times_3 \mathbf{L}_{\text{stack}}(\omega^*),$$

where  $\mathcal{M} : \mathbb{R}^{r+2s} \times \mathbb{R}^{N \times N \times d} \rightarrow \mathbb{R}^{N \times N \times (d+r+2s)}$  is a bilinear function defined as follows:

$$(A.4) \quad \mathcal{M}(\mathbf{a}, \mathcal{B}) = \text{stack} \left( \mathcal{B}, \{a_i \mathbf{G}_i^{I*}\}_{1 \leq i \leq r}, \right. \\ \left. \left\{ a_{r+2k} \mathbf{G}_k^{II,1*} - \frac{a_{r+2k-1}}{\gamma_k^*} \mathbf{G}_k^{II,2*}, a_{r+2k-1} \mathbf{G}_k^{II,2*} + \frac{a_{r+2k}}{\gamma_k^*} \mathbf{G}_k^{II,1*} \right\}_{1 \leq k \leq s} \right),$$

for any  $\mathbf{a} = (a_1, \dots, a_{r+2s})' \in \mathbb{R}^{r+2s}$  and  $\mathcal{B} \in \mathbb{R}^{N \times N \times d}$ , with the true values  $\omega^*$  and  $\mathcal{G}^*$  fixed. The linear approximation in (A.3) will be formalized in the proof of Proposition 3.1; in particular, see (A.14) and (A.15) for the linear form and the remainder term, respectively. Note that both the conclusion and proof techniques of Proposition 3.1 are used repeatedly throughout our theoretical analysis, which is detailed in the supplementary material.

In addition, the following notations will be used in the proof of Proposition 3.1. First, for the convenience of notation in the proof, according to the block form of  $\mathbf{L}(\omega)$ , we partition  $\mathcal{G} \in \mathbb{R}^{N \times N \times d}$  as  $\mathcal{G} = \text{stack}(\mathcal{G}^{\text{AR}}, \mathcal{G}^{\text{MA}}) = (\mathcal{G}^{\text{AR}}, \mathcal{G}^I, \mathcal{G}^{II})$ , where  $\mathcal{G}^{\text{AR}} = \text{stack}(\mathbf{G}_1, \dots, \mathbf{G}_p)$ ,  $\mathcal{G}^I = \text{stack}(\mathbf{G}_1^I, \dots, \mathbf{G}_r^I)$ , and  $\mathcal{G}^{II} = \text{stack}(\mathbf{G}_1^{II,1}, \mathbf{G}_1^{II,2}, \dots, \mathbf{G}_s^{II,1}, \mathbf{G}_s^{II,2})$  are  $N \times N \times p$ ,  $N \times N \times r$ , and  $N \times N \times 2s$  tensors, respectively. Here,  $\mathbf{G}_i^I = \mathbf{G}_{p+i}$  for  $1 \leq i \leq r$ , and  $\mathbf{G}_k^{II,1} = \mathbf{G}_{p+r+2k-1}$  and  $\mathbf{G}_k^{II,2} = \mathbf{G}_{p+r+2k}$  for  $1 \leq k \leq s$  according to our relabeling at the beginning of Section 2.3. Then, for any  $\mathcal{A} = \mathcal{G} \times_3 \mathbf{L}(\omega)$ , we have  $\mathbf{A}_k = \mathbf{G}_k$  for  $1 \leq k \leq p$ , and

$$(A.5) \quad \mathbf{A}_{p+j} = \sum_{i=1}^r \ell_j^I(\lambda_i) \mathbf{G}_i^I + \sum_{k=1}^s \left\{ \ell_j^{II,1}(\eta_k) \mathbf{G}_k^{II,1} + \ell_j^{II,2}(\eta_k) \mathbf{G}_k^{II,2} \right\}, \text{ for } j \geq 1.$$

Moreover, for simplicity, let

$$\mathcal{G}_{\text{stack}} = \mathcal{M}(\omega - \omega^*, \mathcal{G} - \mathcal{G}^*).$$

Equivalently, we can express  $\mathcal{G}_{\text{stack}} = \text{stack}(\mathcal{G} - \mathcal{G}^*, \mathcal{D}(\omega))$  as the  $N \times N \times (d+r+2s)$  tensor formed by augmenting  $\mathcal{G} - \mathcal{G}^*$  with the  $N \times N \times (r+2s)$  tensor

$$\mathcal{D}(\omega) = \text{stack} \left( \{(\lambda_i - \lambda_i^*) \mathbf{G}_i^{I*}\}_{1 \leq i \leq r}, \right.$$

$$\left\{ (\theta_k - \theta_k^*) \mathbf{G}_k^{II,1*} - \frac{\gamma_k - \gamma_k^*}{\gamma_k^*} \mathbf{G}_k^{II,2*}, (\theta_k - \theta_k^*) \mathbf{G}_k^{II,2*} + \frac{\gamma_k - \gamma_k^*}{\gamma_k^*} \mathbf{G}_k^{II,1*} \right\}_{1 \leq k \leq s}.$$

Lastly, note that for every  $\Delta(\omega, \mathcal{G}) \in \Upsilon$ , its corresponding  $\mathcal{G}_{\text{stack}} \in \Xi$ , where

$$(A.6) \quad \Xi = \left\{ \mathcal{M}(\mathbf{a}, \mathcal{B}) \in \mathbb{R}^{N \times N \times (d+r+2s)} \mid \mathbf{a} \in \mathbb{R}^{r+2s}, \mathcal{B} \in \Gamma(2\mathcal{R}_1, 2\mathcal{R}_2) \right\}.$$

**A.3. Proof of Proposition 3.1.** Let  $\Delta = \mathcal{A} - \mathcal{A}^* = \mathcal{G} \times_3 \mathbf{L}(\omega) - \mathcal{G}^* \times_3 \mathbf{L}(\omega^*)$ . Denote by  $\Delta_j$  with  $j \geq 1$  the frontal slices of  $\Delta$ , i.e.  $\Delta_{(1)} = (\Delta_1, \Delta_2, \dots)$ . Then  $\Delta_j = \mathbf{G}_j - \mathbf{G}_j^*$  for  $1 \leq j \leq p$ . For  $j \geq 1$ , by (A.5) and the Taylor expansion,

$$\begin{aligned} \Delta_{p+j} &= \mathbf{A}_{p+j} - \mathbf{A}_{p+j}^* \\ &= \sum_{k=1}^r \left\{ \ell_j^I(\lambda_k^*) + \nabla \ell_j^I(\lambda_k^*)(\lambda_k - \lambda_k^*) + \frac{1}{2} \nabla^2 \ell_j^I(\tilde{\lambda}_k)(\lambda_k - \lambda_k^*)^2 \right\} \mathbf{G}_k^I \\ &\quad + \sum_{k=1}^s \left\{ \ell_j^{II,1}(\eta_k^*) + (\eta_k - \eta_k^*)' \nabla \ell_j^{II,1}(\eta_k^*) \right. \\ &\quad \quad \left. + \frac{1}{2} (\eta_k - \eta_k^*)' \nabla^2 \ell_j^{II,1}(\tilde{\eta}_k)(\eta_k - \eta_k^*) \right\} \mathbf{G}_k^{II,1} \\ &\quad + \sum_{k=1}^s \left\{ \ell_j^{II,2}(\eta_k^*) + (\eta_k - \eta_k^*)' \nabla \ell_j^{II,2}(\eta_k^*) \right. \\ &\quad \quad \left. + \frac{1}{2} (\eta_k - \eta_k^*)' \nabla^2 \ell_j^{II,2}(\tilde{\eta}_k)(\eta_k - \eta_k^*) \right\} \mathbf{G}_k^{II,2} - \mathbf{A}_{p+j}^* \end{aligned}$$

$$(A.7) \quad := \mathbf{H}_j + \mathbf{R}_j,$$

where  $\tilde{\lambda}_k$  lies between  $\lambda_k^*$  and  $\lambda_k$  for  $1 \leq k \leq r$ ,  $\tilde{\eta}_k$  lies between  $\eta_k^*$  and  $\eta_k$  for  $1 \leq k \leq s$ ,

$$\begin{aligned} (A.8) \quad \mathbf{H}_j &= \sum_{k=1}^r \ell_j^I(\lambda_k^*)(\mathbf{G}_k^I - \mathbf{G}_k^{I*}) + \sum_{k=1}^s \sum_{h=1}^2 \ell_j^{II,h}(\eta_k^*)(\mathbf{G}_k^{II,h} - \mathbf{G}_k^{II,h*}) \\ &\quad + \sum_{k=1}^r (\lambda_k - \lambda_k^*) \nabla \ell_j^I(\lambda_k^*) \mathbf{G}_k^{I*} + \sum_{k=1}^s \sum_{h=1}^2 (\eta_k - \eta_k^*)' \nabla \ell_j^{II,h}(\eta_k^*) \mathbf{G}_k^{II,h*}, \end{aligned}$$

and

$$\begin{aligned}
 \mathbf{R}_j = & \sum_{i=1}^r \nabla \ell_j^I(\lambda_k^*)(\lambda_k - \lambda_k^*)(\mathbf{G}_k^I - \mathbf{G}_k^{I*}) \\
 & + \sum_{k=1}^s \sum_{h=1}^2 (\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*)' \nabla \ell_j^{II,h}(\boldsymbol{\eta}_k^*)(\mathbf{G}_k^{II,h} - \mathbf{G}_k^{II,h*}) \\
 & + \frac{1}{2} \sum_{k=1}^r \nabla^2 \ell_j^I(\tilde{\lambda}_k)(\lambda_k - \lambda_k^*)^2 \mathbf{G}_k^I \\
 & + \frac{1}{2} \sum_{k=1}^s \sum_{h=1}^2 (\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*)' \nabla^2 \ell_j^{II,h}(\tilde{\boldsymbol{\eta}}_k)(\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*) \mathbf{G}_k^{II,h}.
 \end{aligned}
 \tag{A.9}$$

We first handle the terms in  $\mathbf{R}_j$ , and denote  $\mathbf{R}_j = \mathbf{R}_{1j} + \mathbf{R}_{2j} + \mathbf{R}_{3j}$ , where

$$\begin{aligned}
 \mathbf{R}_{1j} = & \sum_{k=1}^r \nabla \ell_j^I(\lambda_k^*)(\lambda_k - \lambda_k^*)(\mathbf{G}_k^I - \mathbf{G}_k^{I*}) \\
 & + \sum_{k=1}^s \sum_{h=1}^2 (\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*)' \nabla \ell_j^{II,h}(\boldsymbol{\eta}_k^*)(\mathbf{G}_k^{II,h} - \mathbf{G}_k^{II,h*}), \\
 \mathbf{R}_{2j} = & \frac{1}{2} \sum_{k=1}^r \nabla^2 \ell_j^I(\tilde{\lambda}_k)(\lambda_k - \lambda_k^*)^2 (\mathbf{G}_k^I - \mathbf{G}_k^{I*}) \\
 & + \frac{1}{2} \sum_{k=1}^s \sum_{h=1}^2 (\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*)' \nabla^2 \ell_j^{II,h}(\tilde{\boldsymbol{\eta}}_k)(\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*) (\mathbf{G}_k^{II,h} - \mathbf{G}_k^{II,h*}), \\
 \mathbf{R}_{3j} = & \frac{1}{2} \sum_{k=1}^r \nabla^2 \ell_j^I(\tilde{\lambda}_k)(\lambda_k - \lambda_k^*)^2 \mathbf{G}_k^{I*} \\
 & + \frac{1}{2} \sum_{k=1}^s \sum_{h=1}^2 (\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*)' \nabla^2 \ell_j^{II,h}(\tilde{\boldsymbol{\eta}}_k)(\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*) \mathbf{G}_k^{II,h*}.
 \end{aligned}
 \tag{A.10}$$

Note that for any matrix  $\mathbf{Y} = \sum_{k=1}^d a_k \mathbf{X}_k$ , it holds

$$\|\mathbf{Y}\|_{\text{op}} \leq \|\mathbf{Y}\|_{\text{F}} \leq \left( \sum_{k=1}^d \|\mathbf{X}_k\|_{\text{F}}^2 \right)^{1/2} \left( \sum_{k=1}^d a_k^2 \right)^{1/2} = \|\mathcal{X}\|_{\text{F}} \|\mathbf{a}\|_2,$$

and  $\sum_{k=1}^d a_k^4 \leq (\sum_{k=1}^d a_k^2)^2$ , where  $\mathbf{a} = (a_1, \dots, a_d)' \in \mathbb{R}^d$ , and  $\mathcal{X}$  is a tensor with frontal slices  $\mathbf{X}_k$ 's such that  $\mathcal{X}_{(1)} = (\mathbf{X}_1, \dots, \mathbf{X}_d)$ . Then, by Lemma A.1(i),

$$\begin{aligned}
 \|\mathbf{R}_{1j}\|_{\text{F}} & \leq C_L \bar{\rho}^j \sqrt{\|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2 + 2\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|_2^2} \cdot \|\mathcal{G}^{\text{MA}} - \mathcal{G}^{\text{MA}*}\|_{\text{F}} \\
 & \leq \sqrt{2} C_L \bar{\rho}^j \delta_{\boldsymbol{\omega}} \cdot \|\mathcal{G}^{\text{MA}} - \mathcal{G}^{\text{MA}*}\|_{\text{F}} \leq \sqrt{2} C_L \bar{\rho}^j \delta_{\boldsymbol{\omega}} \delta_{\mathcal{G}},
 \end{aligned}$$

and

$$\|\mathbf{R}_{2j}\|_{\text{F}} \leq \frac{\sqrt{2}}{2} C_L \bar{\rho}^j \delta_{\boldsymbol{\omega}}^2 \sqrt{\sum_{k=1}^r \|\mathbf{G}_k^I - \mathbf{G}_k^{I*}\|_{\text{F}}^2 + \sum_{k=1}^s \sum_{h=1}^2 \|\mathbf{G}_k^{II,h} - \mathbf{G}_k^{II,h*}\|_{\text{F}}^2}$$

$$\leq \frac{\sqrt{2}}{2} C_L \bar{\rho}^j \delta_{\omega}^2 \cdot \|\mathcal{G}^{\text{MA}} - \mathcal{G}^{\text{MA}*}\|_{\text{F}} \leq \frac{\sqrt{2}}{2} C_L \bar{\rho}^j \delta_{\omega}^2 \delta_{\mathcal{G}}.$$

Moreover, by Assumption 2(ii) and Lemma A.1(i), we can show that

$$\|\mathbf{R}_{3j}\|_{\text{F}} \leq C_L C_{\mathcal{G}} \bar{\rho}^j \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega}^2.$$

As a result,

$$\begin{aligned} \|\mathbf{R}_j\|_{\text{F}} &\leq \|\mathbf{R}_{1j}\|_{\text{F}} + \|\mathbf{R}_{2j}\|_{\text{F}} + \|\mathbf{R}_{3j}\|_{\text{F}} \\ (A.11) \quad &\leq C_L \bar{\rho}^j \delta_{\omega} \left( \sqrt{2} \delta_{\mathcal{G}} + \frac{\sqrt{2}}{2} \delta_{\omega} \delta_{\mathcal{G}} + C_{\mathcal{G}} \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \right). \end{aligned}$$

Now consider  $\mathbf{H}_j$  in (A.8). Notice that for any  $j \geq 1$  and  $1 \leq k \leq s$ ,

$$\begin{aligned} \nabla_{\gamma} \ell_j^{II,1}(\boldsymbol{\eta}_k) &= j \gamma_k^{j-1} \cos(j\theta_k) = \frac{1}{\gamma_k} \nabla_{\theta} \ell_j^{II,2}(\boldsymbol{\eta}_k), \\ (A.12) \quad \nabla_{\gamma} \ell_j^{II,2}(\boldsymbol{\eta}_k) &= j \gamma_k^{j-1} \sin(j\theta_k) = -\frac{1}{\gamma_k} \nabla_{\theta} \ell_j^{II,1}(\boldsymbol{\eta}_k). \end{aligned}$$

Thus, the last term on the right side of (A.8) can be simplified to

$$\begin{aligned} &\sum_{k=1}^s \sum_{h=1}^2 (\boldsymbol{\eta}_k - \boldsymbol{\eta}_k^*)' \nabla \ell_j^{II,h}(\boldsymbol{\eta}_k^*) \mathbf{G}_k^{II,h*} \\ (A.13) \quad &= \sum_{k=1}^s \left[ (\theta_k - \theta_k^*) \mathbf{G}_k^{II,1*} - \frac{1}{\gamma_k^*} (\gamma_k - \gamma_k^*) \mathbf{G}_k^{II,2*} \right] \nabla_{\theta} \ell_j^{II,1}(\boldsymbol{\eta}_k^*) \\ &\quad + \sum_{k=1}^s \left[ (\theta_k - \theta_k^*) \mathbf{G}_k^{II,2*} + \frac{1}{\gamma_k^*} (\gamma_k - \gamma_k^*) \mathbf{G}_k^{II,1*} \right] \nabla_{\theta} \ell_j^{II,2}(\boldsymbol{\eta}_k^*). \end{aligned}$$

Let  $\mathcal{H} = \text{stack}(\mathbf{H}_1, \mathbf{H}_2, \dots)$  and  $\mathcal{R} = \text{stack}(\mathbf{R}_1, \mathbf{R}_2, \dots)$ . Then by (A.8) and (A.13), it can be verified that

$$\begin{aligned} \widetilde{\mathcal{H}} &:= \text{stack}(\mathbf{G}_1 - \mathbf{G}_1^*, \dots, \mathbf{G}_p - \mathbf{G}_p^*, \mathcal{H}) \\ &= (\mathcal{G} - \mathcal{G}^*) \times_3 \mathbf{L}(\boldsymbol{\omega}^*) + \mathcal{D}(\boldsymbol{\omega}) \times_3 (\nabla \mathbf{L}^I(\boldsymbol{\lambda}^*), \nabla_{\theta} \mathbf{L}^{II}(\boldsymbol{\eta}^*)) \\ (A.14) \quad &= \mathcal{G}_{\text{stack}} \times_3 \mathbf{L}_{\text{stack}}(\boldsymbol{\omega}^*), \end{aligned}$$

where  $\mathcal{D}(\boldsymbol{\omega}) \in \mathbb{R}^{N \times N \times (r+2s)}$  and  $\mathcal{G}_{\text{stack}} \in \mathbb{R}^{N \times N \times (d+r+2s)}$  are defined in Appendix A.2.

Note that

$$(A.15) \quad \boldsymbol{\Delta} = \widetilde{\mathcal{H}} + \text{stack}(\mathbf{0}_{N \times N \times p}, \mathcal{R}).$$

Moreover,

$$\|\mathcal{D}(\boldsymbol{\omega})\|_{\text{F}}^2 = \sum_{i=1}^r (\lambda_i - \lambda_i^*)^2 \|\mathbf{G}_i^{I*}\|_{\text{F}}^2 + \sum_{k=1}^s \left\| (\theta_k - \theta_k^*) \mathbf{G}_k^{II,1*} - \frac{\gamma_k - \gamma_k^*}{\gamma_k^*} \mathbf{G}_k^{II,2*} \right\|_{\text{F}}^2$$

$$\begin{aligned}
& + \sum_{k=1}^s \left\| (\theta_k - \theta_k^*) \mathbf{G}_k^{II,2*} + \frac{\gamma_k - \gamma_k^*}{\gamma_k^*} \mathbf{G}_k^{II,1*} \right\|_{\mathbf{F}}^2 \\
& = \sum_{i=1}^r (\lambda_i - \lambda_i^*)^2 \|\mathbf{G}_i^{I*}\|_{\mathbf{F}}^2 + \sum_{k=1}^s (\theta_k - \theta_k^*)^2 (\|\mathbf{G}_k^{II,1*}\|_{\mathbf{F}}^2 + \|\mathbf{G}_k^{II,2*}\|_{\mathbf{F}}^2) \\
& + \sum_{k=1}^s \frac{(\gamma_k - \gamma_k^*)^2}{\gamma_k^{*2}} (\|\mathbf{G}_k^{II,1*}\|_{\mathbf{F}}^2 + \|\mathbf{G}_k^{II,2*}\|_{\mathbf{F}}^2),
\end{aligned} \tag{A.16}$$

which, together with Assumption 2(ii), leads to

$$\sqrt{2}c_g \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \leq \|\mathcal{D}(\omega)\|_{\mathbf{F}} \leq \frac{\sqrt{2}C_g}{\min_{1 \leq k \leq s} \gamma_k^*} \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega}. \tag{A.17}$$

By the simple inequalities  $(|x| + |y|)/2 \leq \sqrt{x^2 + y^2} \leq |x| + |y|$ , we have  $0.5(\delta_g + \|\mathcal{D}(\omega)\|_{\mathbf{F}}) \leq \|\mathcal{G}_{\text{stack}}\|_{\mathbf{F}} \leq \delta_g + \|\mathcal{D}(\omega)\|_{\mathbf{F}}$ , and thus in view of (A.17) we further have

$$0.5 \left( \delta_g + \sqrt{2}c_g \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \right) \leq \|\mathcal{G}_{\text{stack}}\|_{\mathbf{F}} \leq \delta_g + \frac{\sqrt{2}C_g}{\min_{1 \leq k \leq s} \gamma_k^*} \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \tag{A.18}$$

where  $\delta_g = \|\mathcal{G} - \mathcal{G}^*\|_{\mathbf{F}}$ . By Lemma A.2,  $\sigma_{\min,L} = \sigma_{\min}(\mathbf{L}_{\text{stack}}(\omega^*)) > 0$ . Then it follows from (A.18) that

$$0.5\sigma_{\min,L} \left( \delta_g + \sqrt{2}c_g \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \right) \leq \|\widetilde{\mathcal{H}}\|_{\mathbf{F}} \leq \sigma_{\max,L} \left( \delta_g + \frac{\sqrt{2}C_g \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega}}{\min_{1 \leq k \leq s} \gamma_k^*} \right).$$

Combining this with (A.11), (A.15), (A.17), as well as the fact that  $\|\mathcal{G}^{\text{MA}} - \mathcal{G}^{\text{MA}*}\|_{\mathbf{F}} \leq \delta_g$ , we have

$$\begin{aligned}
\|\Delta\|_{\mathbf{F}} & \leq \|\widetilde{\mathcal{H}}\|_{\mathbf{F}} + \|\mathcal{R}\|_{\mathbf{F}} \\
& \leq \left\{ \sigma_{\max,L} + \frac{\sqrt{2}C_L}{1 - \bar{\rho}} \left( \delta_{\omega} + \frac{\delta_{\omega}^2}{2} \right) \right\} \delta_g + \left( \frac{\sqrt{2}C_g \sigma_{\max,L}}{\min_{1 \leq k \leq s} \gamma_k^*} + \frac{C_L C_g}{1 - \bar{\rho}} \delta_{\omega} \right) \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega}
\end{aligned}$$

and

$$\begin{aligned}
\|\Delta\|_{\mathbf{F}} & \geq \|\widetilde{\mathcal{H}}\|_{\mathbf{F}} - \|\mathcal{R}\|_{\mathbf{F}} \\
& \geq \left\{ 0.5\sigma_{\min,L} - \frac{\sqrt{2}C_L}{1 - \bar{\rho}} \left( \delta_{\omega} + \frac{\delta_{\omega}^2}{2} \right) \right\} \delta_g + \left( \frac{c_g \sigma_{\min,L}}{\sqrt{2}} - \frac{C_L C_g}{1 - \bar{\rho}} \delta_{\omega} \right) \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega}.
\end{aligned}$$

Thus, by taking

$$\delta_{\omega} \leq c_{\omega} = \min \left\{ 2, \frac{c_g(1 - \bar{\rho})\sigma_{\min,L}}{8\sqrt{2}C_L C_g} \right\},$$

we can show that

$$c_{\Delta} \left( \delta_g + \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \right) \leq \|\Delta\|_{\mathbf{F}} \leq C_{\Delta} \left( \delta_g + \sqrt{\mathcal{R}_1 \wedge \mathcal{R}_2} \delta_{\omega} \right),$$

where

$$c_{\Delta} = c_l \cdot \sigma_{\min,L} \asymp 1 \quad \text{and} \quad C_{\Delta} = c_u \cdot \max \{ \sigma_{\max,L}, (1 - \bar{\rho})^{-1} \} \asymp 1,$$

with  $c_l = 0.25 \min \{ 1, \sqrt{2}c_g \}$  and  $c_u = 1 + \sqrt{2}C_g (\min_{1 \leq k \leq s} \gamma_k^*)^{-1} + 4\sqrt{2}C_L + 2C_L C_g$ .

**A.4. Proof of Theorem 3.2.** Let  $\mathbf{x}_t = (\mathbf{y}'_{t-1}, \mathbf{y}'_{t-2}, \dots)'$ . Denote by  $\Delta_j$  with  $j \geq 1$  the frontal slices of  $\Delta$ , i.e.,  $\Delta_{(1)} = (\Delta_1, \Delta_2, \dots)$ . Denote

$$\begin{aligned}
 S_1(\Delta) &= \frac{2}{T} \sum_{t=1}^T \left\langle \sum_{j=1}^{\infty} \Delta_j \mathbf{y}_{t-j}, \sum_{k=t}^{\infty} \Delta_k \mathbf{y}_{t-k} \right\rangle, \\
 S_2(\Delta) &= \frac{2}{T} \sum_{t=1}^T \left\langle \sum_{j=t}^{\infty} \mathbf{A}_j^* \mathbf{y}_{t-j}, \sum_{k=1}^{t-1} \Delta_k \mathbf{y}_{t-k} \right\rangle, \\
 S_3(\Delta) &= \frac{2}{T} \sum_{t=1}^T \left\langle \varepsilon_t, \sum_{j=t}^{\infty} \Delta_j \mathbf{y}_{t-j} \right\rangle.
 \end{aligned}
 \tag{A.19}$$

The following three lemmas are sufficient for the proof of Theorem 3.2. We state them here and provide detailed proofs in the separate supplementary file.

**LEMMA A.3** (Strong convexity and smoothness properties). *Under Assumptions 1–3, if  $T \gtrsim (\kappa_2/\kappa_1)^2 d_{\mathcal{M}} \log(\kappa_2/\kappa_1)$ , then with probability at least  $1 - 2e^{-cd_{\mathcal{M}} \log(\kappa_2/\kappa_1)} - 3e^{-cN}$ ,*

$$\kappa_1 \|\Delta\|_{\text{F}}^2 \lesssim \frac{1}{T} \sum_{t=1}^T \|\Delta_{(1)} \mathbf{x}_t\|_2^2 \lesssim \kappa_2 \|\Delta\|_{\text{F}}^2, \quad \forall \Delta \in \Upsilon.$$

**LEMMA A.4** (Deviation bound). *Under the conditions of Lemma A.3, with probability at least  $1 - 2e^{-cd_{\mathcal{M}} \log(\kappa_2/\kappa_1)} - 5e^{-cN}$ ,*

$$\frac{1}{T} \left| \sum_{t=1}^T \langle \varepsilon_t, \Delta_{(1)} \mathbf{x}_t \rangle \right| \lesssim \sqrt{\frac{\kappa_2 \lambda_{\max}(\Sigma_{\varepsilon}) d_{\mathcal{M}}}{T}} \|\Delta\|_{\text{F}}, \quad \forall \Delta \in \Upsilon.$$

**LEMMA A.5** (Effects of initial values). *Under Assumptions 1–3, if  $T \gtrsim (\kappa_2/\kappa_1) d_{\mathcal{M}}$ , then with probability at least  $1 - \{2 + \sqrt{\kappa_2/\lambda_{\max}(\Sigma_{\varepsilon})}\} \sqrt{N/(\mathcal{R}_1 + \mathcal{R}_2)T}$ ,*

$$|S_1(\Delta)| \lesssim \kappa_1 \|\Delta\|_{\text{F}}^2, \quad |S_i(\Delta)| \lesssim \sqrt{\frac{\kappa_2 \lambda_{\max}(\Sigma_{\varepsilon}) d_{\mathcal{M}}}{T}} \|\Delta\|_{\text{F}}, \quad i = 2, 3, \quad \forall \Delta \in \Upsilon.$$

Now we prove Theorem 3.2. Note that  $\sum_{j=1}^{t-1} \mathbf{A}_j \mathbf{y}_{t-j} = \mathcal{A}_{(1)} \tilde{\mathbf{x}}_t$ . Due to the optimality of  $\hat{\mathcal{A}}$ , we have

$$\sum_{t=1}^T \|\mathbf{y}_t - \mathcal{A}_{(1)}^* \tilde{\mathbf{x}}_t - \hat{\Delta}_{(1)} \tilde{\mathbf{x}}_t\|_2^2 \leq \sum_{t=1}^T \|\mathbf{y}_t - \mathcal{A}_{(1)}^* \tilde{\mathbf{x}}_t\|_2^2,$$

Then, since  $\mathbf{y}_t - \mathcal{A}_{(1)}^* \tilde{\mathbf{x}}_t = \varepsilon_t + \sum_{j=t}^{\infty} \mathbf{A}_j^* \mathbf{y}_{t-j}$ , it follows that

$$\begin{aligned}
 \frac{1}{T} \sum_{t=1}^T \|\hat{\Delta}_{(1)} \tilde{\mathbf{x}}_t\|_2^2 &\leq \frac{2}{T} \sum_{t=1}^T \langle \varepsilon_t, \hat{\Delta}_{(1)} \tilde{\mathbf{x}}_t \rangle + \underbrace{\frac{2}{T} \sum_{t=1}^T \left\langle \sum_{j=t}^{\infty} \mathbf{A}_j^* \mathbf{y}_{t-j}, \hat{\Delta}_{(1)} \tilde{\mathbf{x}}_t \right\rangle}_{S_2(\hat{\Delta})} \\
 &= \frac{2}{T} \sum_{t=1}^T \langle \varepsilon_t, \hat{\Delta}_{(1)} \mathbf{x}_t \rangle + S_2(\hat{\Delta}) - S_3(\hat{\Delta}),
 \end{aligned}
 \tag{A.20}$$

where  $S_2(\cdot)$  and  $S_3(\cdot)$  are defined as in (A.19),  $\hat{\Delta}_{(1)}\tilde{\mathbf{x}}_t = \sum_{k=1}^{t-1} \hat{\Delta}_k \mathbf{y}_{t-k}$ , and  $\hat{\Delta}_{(1)}\mathbf{x}_t = \sum_{k=1}^{\infty} \hat{\Delta}_k \mathbf{y}_{t-k}$ .

Moreover, applying the inequality  $\|\mathbf{a} - \mathbf{b}\|_2^2 \geq \|\mathbf{a}\|_2^2 - 2\langle \mathbf{a}, \mathbf{b} \rangle$  with  $\mathbf{a} = \hat{\Delta}_{(1)}\mathbf{x}_t = \sum_{j=1}^{\infty} \hat{\Delta}_j \mathbf{y}_{t-j}$  and  $\mathbf{b} = \sum_{k=t}^{\infty} \hat{\Delta}_k \mathbf{y}_{t-k}$ , we can lower bound the left-hand side of (A.20) to further obtain that

$$(A.21) \quad \frac{1}{T} \sum_{t=1}^T \|\hat{\Delta}_{(1)}\mathbf{x}_t\|_2^2 - S_1(\hat{\Delta}) \leq \frac{2}{T} \sum_{t=1}^T \langle \varepsilon_t, \hat{\Delta}_{(1)}\mathbf{x}_t \rangle + S_2(\hat{\Delta}) - S_3(\hat{\Delta}),$$

where  $S_1(\cdot)$  is defined as in (A.19). It is worth pointing out that  $S_i(\hat{\Delta})$  for  $1 \leq i \leq 3$  capture the initialization effect of  $\mathbf{y}_s = \mathbf{0}$  for  $s \leq 0$  on the estimation.

Note that  $\hat{\Delta} = \hat{\mathcal{A}} - \mathcal{A}^* \in \Upsilon$  and  $\kappa_2 \geq \kappa_1$ . Suppose that the high probability events in Lemmas A.3–A.5 hold. Then we can derive the estimation error bound from (A.21):

$$\kappa_1 \|\hat{\Delta}\|_F^2 \lesssim \sqrt{\frac{\kappa_2 \lambda_{\max}(\Sigma_{\varepsilon}) d_{\mathcal{M}}}{T}} \|\hat{\Delta}\|_F, \quad \text{or} \quad \|\hat{\Delta}\|_F \lesssim \sqrt{\frac{\kappa_2 \lambda_{\max}(\Sigma_{\varepsilon}) d_{\mathcal{M}}}{\kappa_1^2 T}}.$$

Furthermore, applying Lemma A.5 again, we can derive the prediction error bound from (A.20) and the above result as follows:

$$\frac{1}{T} \sum_{t=1}^T \|\hat{\Delta}_{(1)}\tilde{\mathbf{x}}_t\|_2^2 \lesssim \frac{\kappa_2 \lambda_{\max}(\Sigma_{\varepsilon}) d_{\mathcal{M}}}{\kappa_1 T}.$$

The proof of this theorem is complete.

## SUPPLEMENTARY MATERIAL

**Supplement for “SARMA: A Computationally Scalable High-Dimensional Time Series Model”** provides all remaining technical proofs and further details for Sections 5 and 6.

## REFERENCES

- [1] ANDERSEN, T. G., BOLLERSLEV, T., CHRISTOFFERSEN, P. F. and DIEBOLD, F. X. (2006). *Volatility and correlation forecasting* **1**. Elsevier.
- [2] ANDERSEN, T. G., BOLLERSLEV, T., DIEBOLD, F. X. and LABYS, P. (2003). Modeling and forecasting realized volatility. *Econometrica* **71** 579–625.
- [3] ATHANASOPOULOS, G. and VAHID, F. (2008). VARMA versus VAR for macroeconomic forecasting. *Journal of Business & Economic Statistics* **26** 237–252.
- [4] BAI, J. and NG, S. (2008). Large dimensional factor analysis. *Foundations and Trends® in Econometrics* **3** 89–163.

- [5] BAI, J. and WANG, P. (2016). Econometric analysis of large factor models. *Annual Review of Economics* **8** 53–80.
- [6] BASU, S., LI, X. and MICHAILIDIS, G. (2019). Low rank and structured modeling of high-dimensional vector autoregressions. *IEEE Transactions on Signal Processing* **67** 1207–1222.
- [7] BASU, S. and MATTESON, D. S. (2021). A survey of estimation methods for sparse high-dimensional time series models. *ArXiv preprint arXiv:2107.14754*.
- [8] BASU, S. and MICHAILIDIS, G. (2015). Regularized estimation in sparse high-dimensional time series models. *The Annals of Statistics* **43** 1535–1567.
- [9] BULUT, U. (2016). Do financial conditions have a predictive power on inflation in Turkey? *International Journal of Economics and Financial Issues* **6** 621–628.
- [10] CHAN, J. C. C., EISENSTAT, E. and KOOP, G. (2016). Large Bayesian VARMA. *Journal of Econometrics* **192** 374–390.
- [11] CHEN, K., CHAN, K.-S. and STENSETH, N. C. (2012). Reduced rank stochastic regression with a sparse singular value decomposition. *Journal of the Royal Statistical Society: Series B* **74** 203–221.
- [12] DE LATHAUWER, L., DE MOOR, B. and VANDEWALLE, J. (2000). A multilinear singular value decomposition. *SIAM Journal on Matrix Analysis and Applications* **21** 1253–1278.
- [13] DIAS, G. F. and KAPETANIOS, G. (2018). Estimation and forecasting in vector autoregressive moving average models for rich datasets. *Journal of Econometrics* **202** 75–91.
- [14] DOWELL, J. and PINSON, P. (2016). Very-short-term probabilistic wind power forecasts by sparse vector autoregression. *IEEE Transactions on Smart Grid* **7** 763–770.
- [15] GANDY, S., RECHT, B. and YAMADA, I. (2011). Tensor completion and low-n-rank tensor recovery via convex optimization. *Inverse Problems* **27** 025010.
- [16] GOODHART, C., HOFMANN, B. et al. (2001). Asset prices, financial conditions, and the transmission of monetary policy. In *Conference on Asset Prices, Exchange Rates, and Monetary Policy, Stanford University* 2–3. Citeseer.
- [17] GORROSTIETA, C., OMBAO, H., BÉDARD, P. and SANES, J. N. (2012). Investigating brain connectivity using mixed effects vector autoregressive models. *NeuroImage* **59** 3347–3355.
- [18] HALLIN, M. and LIPPI, M. (2013). Factor models in high-dimensional time series: A time-domain approach. *Stochastic Processes and their Applications* **123** 2678–2695.
- [19] HAN, R., WILLETT, R. and ZHANG, A. (2021). An optimal statistical and computational framework for generalized tensor estimation. *The Annals of Statistics*. To appear.
- [20] HARTFIEL, D. J. (1995). Dense sets of diagonalizable matrices. *Proceedings of the American Mathematical Society* **123** 1669–1672.
- [21] HATZIUS, J., HOOPER, P., MISHKIN, F. S., SCHOENHOLTZ, K. L. and WATSON, M. W. (2010). Financial conditions indexes: A fresh look after the financial crisis. Working Paper No. 16150, National Bureau of Economic Research.



- [22] HORN, R. A. and JOHNSON, C. R. (2012). *Matrix Analysis*, 2nd ed. Cambridge University Press, New York.
- [23] KOLDA, T. G. and BADER, B. W. (2009). Tensor decompositions and applications. *SIAM Review* **51** 455–500.
- [24] LAM, C. and YAO, Q. (2012). Factor modeling for high-dimensional time series: inference for the number of factors. *The Annals of Statistics* **40** 694–726.
- [25] LOZANO, A. C., ABE, N., LIU, Y. and ROSSET, S. (2009). Grouped graphical Granger modeling for gene expression regulatory networks discovery. *Bioinformatics* **25** i110–i118.
- [26] LÜTKEPOHL, H. (2005). *New Introduction to Multiple Time Series Analysis*. Springer Science & Business Media.
- [27] MCCracken, M. W. and NG, S. (2016). FRED-MD: A monthly database for macroeconomic research. *Journal of Business & Economic Statistics* **34** 574–589.
- [28] METAXOGLU, K. and SMITH, A. (2007). Maximum likelihood estimation of VARMA models using a state-space EM algorithm. *Journal of Time Series Analysis* **28** 666–685.
- [29] NEGAHBAN, S. and WAINWRIGHT, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics* **39** 1069–1097.
- [30] PEÑA, D. and TSAY, R. S. (2021). *Statistical Learning for Big Dependent Data*. John Wiley & Sons, Hoboken, New Jersey.
- [31] RASKUTTI, G., YUAN, M. and CHEN, H. (2019). Convex regularization for high-dimensional multi-response tensor regression. *The Annals of Statistics* **47** 1554–1584.
- [32] STOCK, J. H. and WATSON, M. W. (2005). Implications of dynamic factor models for VAR analysis. *National Bureau of Economic Research Working Paper No. 11467*.
- [33] STOCK, J. H. and WATSON, M. W. (2011). Dynamic factor models. In *Oxford Handbook of Economic Forecasting* (M. P. Clements and D. F. Hendry, eds.) Oxford University Press.
- [34] TSAY, R. S. (2014). *Multivariate Time Series Analysis: With R and Financial Applications*. John Wiley & Sons, Hoboken, New Jersey.
- [35] TUCKER, L. R. (1966). Some mathematical notes on three-mode factor analysis. *Psychometrika* **31** 279–311.
- [36] VELU, R. P., REINSEL, G. C. and WICHERN, D. W. (1986). Reduced rank models for multiple time series. *Biometrika* **73** 105–118.
- [37] WANG, D., ZHENG, Y. and LI, G. (2021). High-dimensional low-rank tensor autoregressive time series modelling. *arXiv preprint arXiv:2101.04276*.
- [38] WANG, D., ZHENG, Y., LIAN, H. and LI, G. (2021). High-dimensional vector autoregressive time series modeling via tensor decomposition. *Journal of the American Statistical Association* 1–19.
- [39] WILMS, I., BASU, S., BIEN, J. and MATTESON, D. (2021). Sparse identification and estimation of large-scale vector autoregressive moving averages. *Journal of the American Statistical Association*. To appear.

- [40] XIA, Q., XU, W. and ZHU, L. (2015). Consistently determining the number of factors in multivariate volatility modelling. *Statistica Sinica* **25** 1025–1044.
- [41] ZHENG, Y. and CHENG, G. (2021). Finite time analysis of vector autoregressive models under linear restrictions. *Biometrika* **108** 469–489.