

Journal of the American Statistical Association



ISSN: (Print) (Online) Journal homepage: www.tandfonline.com/journals/uasa20

An Interpretable and Efficient Infinite-Order Vector Autoregressive Model for High-Dimensional Time Series

Yao Zheng

To cite this article: Yao Zheng (11 Mar 2024): An Interpretable and Efficient Infinite-Order Vector Autoregressive Model for High-Dimensional Time Series, Journal of the American Statistical Association, DOI: 10.1080/01621459.2024.2311365

To link to this article: https://doi.org/10.1080/01621459.2024.2311365







An Interpretable and Efficient Infinite-Order Vector Autoregressive Model for High-Dimensional Time Series

Yao Zheng

Department of Statistics, University of Connecticut, Storrs, CT

ABSTRACT

As a special infinite-order vector autoregressive (VAR) model, the vector autoregressive moving average (VARMA) model can capture much richer temporal patterns than the widely used finite-order VAR model. However, its practicality has long been hindered by its non-identifiability, computational intractability, and difficulty of interpretation, especially for high-dimensional time series. This article proposes a novel sparse infinite-order VAR model for high-dimensional time series, which avoids all above drawbacks while inheriting essential temporal patterns of the VARMA model. As another attractive feature, the temporal and crosssectional structures of the VARMA-type dynamics captured by this model can be interpreted separately, since they are characterized by different sets of parameters. This separation naturally motivates the sparsity assumption on the parameters determining the cross-sectional dependence. As a result, greater statistical efficiency and interpretability can be achieved with little loss of temporal information. We introduce two ℓ_1 -regularized estimation methods for the proposed model, which can be efficiently implemented via block coordinate descent algorithms, and derive the corresponding nonasymptotic error bounds. A consistent model order selection method based on the Bayesian information criteria is also developed. The merit of the proposed approach is supported by simulation studies and a real-world macroeconomic data analysis. Supplementary materials for this article are available online including a standardized description of the materials available for reproducing the work.

ARTICLE HISTORY

Received September 2022 Accepted January 2024

KEYWORDS

Granger causality; High-dimensional time series; Infinite-order vector autoregression; Sparse estimation; VARMA

1. Introduction

Let $y_t \in \mathbb{R}^N$ be the observation of an N-dimensional time series at time t. The need for modeling y_t with a large dimension N is ubiquitous, ranging from economics and finance (Nicholson et al. 2020; Wilms et al. 2023) to biology and neuroscience (Lozano et al. 2009; Gorrostieta et al. 2012), and to environmental and health sciences (Dowell and Pinson 2016; Davis, Zang, and Zheng 2016). For modeling y_t , three issues are of particular importance:

- (II) Flexibility of temporal dynamics: As N increases, it is more likely that y_t contains component series with complex temporal dependence structures. Then information further in the past may be needed to generate more flexible temporal dynamics.
- (I2) Efficiency: It is important that the estimation is efficient both statistically and computationally under large *N*, so that accurate forecasts can be obtained.
- (I3) Interpretability: Ideally, the model should have easy interpretations, such as direct implications of Granger causality (Granger 1969) among the *N* component series.

The finite-order vector autoregressive (VAR) model, coupled with dimension reduction techniques such as sparse (Basu and Matteson 2021) and low-rank (Wang et al. 2022) methods, has been widely studied for high-dimensional time series. This

model is highly popular due to its theoretical and computational tractability, and the coefficient matrices have intuitive interpretations analogous to those in the multivariate linear regression. However, in practice, a large lag order is often required for the VAR model to adequately fit the data (Chan, Eisenstat, and Koop 2016; Nicholson et al. 2020). Thus, it is more realistic to assume that the data follow the more general, infinite-order VAR $(VAR(\infty))$ process:

$$y_t = \sum_{h=1}^{\infty} A_h y_{t-h} + \varepsilon_t, \tag{1.1}$$

where ε_t are the innovations, and $A_h \in \mathbb{R}^{N \times N}$ are the AR coefficient matrices; in particular, it reduces to the VAR(P) model when $A_h = \mathbf{0}$ for h > P. In fact, if a sample $\{y_t\}_{t=1}^T$ is generated from (1.1), we can approximate it by a VAR(P) model provided that $P \to \infty$ at an appropriate rate as the sample size $T \to \infty$ (Lütkepohl 2005), which in turn explains the practical need for a large P. Nonetheless, for y_t in (1.1) to be stationary, A_h must diminish quickly as $h \to \infty$; otherwise, the infinite sum will be ill-defined. The decay property of A_h , coupled with a large P, will not only pose difficulties in high-dimensional estimation, but make the fitted VAR(P) model hard to interpret. Take the Lasso estimator of the VAR(P) model with sparse A_h 's. Since all entries of A_h must be small at even moderately large h, the Lasso may fail to capture the significant

yet small entries. Moreover, the sparsity pattern of A_h for the fitted model generally varies substantially across h, making it even more difficult to interpret A_h 's simultaneously (Shojaie, Basu, and Michailidis 2012; Nicholson et al. 2020).

In the literature on multivariate time series, an alternative approach to infinite-order VAR modeling is to consider the vector autoregressive moving average (VARMA) model. For example, the VARMA(1, 1) model is

$$y_t = \mathbf{\Phi} y_{t-1} + \boldsymbol{\varepsilon}_t - \mathbf{\Theta} \boldsymbol{\varepsilon}_{t-1}, \tag{1.2}$$

where $\Phi, \Theta \in \mathbb{R}^{N \times N}$ are the AR and MA coefficient matrices. Assuming that (1.2) is invertible, that is, all eigenvalues of Θ are less than one in absolute value, (1.2) can be written as the $VAR(\infty)$ process in (1.1) with $A_h = A_h(\Phi, \Theta) = \Theta^{h-1}(\Phi - \Theta)$ for $h \geq 1$. Note that A_h diminishes quickly as $h \rightarrow \infty$ due to the exponential factor $\mathbf{\Theta}^{h-1}$, so the VAR(∞) process is well defined. Hence, the MA part of the model is the key to parsimoniously generating $VAR(\infty)$ -type temporal dynamics. For the general VARMA(p, q) model, $\mathbf{y}_t = \sum_{i=1}^p \mathbf{\Phi}_i \mathbf{y}_{t-i} + \boldsymbol{\varepsilon}_t - \mathbf{v}_t$ $\sum_{i=1}^{q} \Theta_{i} \varepsilon_{t-j}$, the richness of temporal patterns will increase with p and q, but with only small orders p and q, the VARMA model can usually provide more accurate forecasts than largeorder VAR models in practice (Athanasopoulos and Vahid 2008; Chan, Eisenstat, and Koop 2016). Compared with finite-order VAR models, the VARMA model is more favorable in terms of (I1) but suffers from severe drawbacks regarding (I2), as its computation is generally complicated due to the following two problems:

- (P1) Non-identifiability: For example, in the VARMA(1, 1) case, there are multiple pairs of (Θ, Φ) corresponding to the same process. The root cause of this problem is the matrix multiplications in the parametric form of $A_h(\Phi, \Theta)$ = $\mathbf{\Theta}^{h-1}(\mathbf{\Phi} - \mathbf{\Theta}).$
- (P2) High-order matrix polynomials: Consider as an example the ordinary least squares (OLS) estimation of the VARMA(1, 1) model. For a sample $\{y_t\}_{t=1}^T$, since $A_h(\mathbf{\Phi}, \mathbf{\Theta})$ is an hth-order matrix polynomial for $1 \le h \le T$, the loss function will have a computational complexity of $O(T^2N^3)^1$, hence, unscalable under large N.

While recent attempts have been made to improve the feasibility of VARMA models (Metaxoglou and Smith 2007; Chan, Eisenstat, and Koop 2016; Dias and Kapetanios 2018; Wilms et al. 2023), they do not tackle (P1) and (P2) directly, but rather resort to sophisticated identification constraints and optimization methods. Moreover, high-dimensional VARMA models can be difficult to interpret due to their latent MA structures. Particularly, while it may be natural to assume that Θ and Φ in (1.2) are sparse under large N (Wilms et al. 2023), this does not necessarily result in a sparse $VAR(\infty)$ model; that is, $A_h(\Phi, \Theta)$'s may not be sparse. Thus, the sparse VARMA model is not particularly attractive in terms of (I3).

For high-dimensional time series, we aim to develop a sparse $VAR(\infty)$ model that is favorable in all of (I1)–(I3). The proposed approach is motivated by reparameterizing the VAR(∞) form of the VARMA(p, q) model into formulation (1.1) with

$$A_h = \sum_{k=1}^d \ell_{h,k}(\boldsymbol{\omega}) G_k \quad \text{for} \quad h \ge 1, \tag{1.3}$$

where $G_1, \ldots, G_d \in \mathbb{R}^{N \times N}$ are unknown coefficient matrices, $\{\ell_{h,k}(\cdot)\}_{h=1}^{\infty}$ for $1 \leq k \leq d$ are different sequences of realvalued functions characterizing the exponential decay pattern of A_h , with $\ell_{h,k}(\omega) \rightarrow 0$ as $h \rightarrow \infty$ for each k, and ω is an unknown low-dimensional parameter vector; see also Huang, Lu, and Zheng (2023) for a high-dimensional Tuckerlow-rank time series model concurrently developed from (1.3) with different techniques and interpretations. Similar to the orders (p,q) of the VARMA model, d can be viewed as the overall order that controls the complexity of temporal patterns of the VAR(∞) model; see Section 2 for the detailed model formulation. Note that (1.3) preserves the essential temporal patterns of the VARMA process, since it is derived directly from the former with little loss of generality. Thus, it is fundamentally more flexible than finite-order VAR models, that is, more desirable regarding (I1). Moreover, each $A_h = A_h(\omega, G_1, \dots, G_d)$ in (1.3) is a linear combination of matrices. Hence, unlike $A_h(\Phi, \Theta)$ mentioned above, this form of A_h gets rid of all matrix multiplications. As a result, both problems (P1) and (P2) are eliminated, and then (I2) can be achieved. To tackle the high dimensionality, we assume that G_k 's are sparse, leading to the proposed sparse parametric $VAR(\infty)$ (SPVAR(∞)) model. In addition to improving the estimation efficiency as required by (I2), the sparsity assumption enables greater interpretability, that is, (I3), thanks to the novel separation of temporal and cross-sectional dependence in parameterizing the VARMA-type dynamic structure:

- (D1) Temporal dependence: In (1.3), the decay pattern of A_h as $h \to \infty$ is fully characterized by the scalar weights
- (D2) Cross-sectional dependence: The G_k 's, independent of the above decay pattern as $h \to \infty$, fully capture the crosssectional dependence.

As a result of (D2), the Granger causal network of the N component series of y_t is directly linked to the aggregate sparsity pattern of G_k 's. Moreover, as detailed in Section 2.1, $\{\ell_{h,k}(\boldsymbol{\omega})\}_{h=1}^{\infty}$'s in (1.3) are specifically defined such that $A_k = G_k$ for $1 \le k \le p$, whereas A_{p+j} for $j \geq 1$ are expressed as linear combinations of G_{p+1}, \ldots, G_d , where p is the AR order of the VARMA(p, q)model from which (1.3) originates. Consequently, there is an interesting dichotomy in the interpretations of different G_k 's: On the one hand, each G_k with $1 \le k \le p$ has the same interpretation as the lag-k AR coefficient matrix of the VAR(p) model, capturing the short-term cross-sectional dependence. On the other hand, the "MA" coefficient matrices G_{p+1}, \ldots, G_d encapsulate the cross-sectional dependence associated with the VARMA-type temporal structure, that is, the long-term influence among the component series that extends into high lags. It is worth noting that the Granger causal network each G_k

¹The computational complexity in this article is calculated in a model of computation where field operations (addition and multiplication) take constant

individually captures is specific to a particular temporal pattern characterized by $\{\ell_{h,k}(\boldsymbol{\omega})\}_{h=1}^{\infty}$. This granularity provides a more detailed perspective on Granger causality from a temporal standpoint; see Section 2.2 for details. Additionally, in view of (D1), the sparsity of G_k 's incurs little loss of temporal information, so the essential VARMA-type temporal pattern is well preserved. This is a distinct advantage over regularized VARMA models (Chan, Eisenstat, and Koop 2016; Wilms et al. 2023).

In fact, even compared to sparse finite-order VAR models, the proposed model can be more interpretable for the following two reasons. First, while the AR coefficient matrices A_h must diminish quickly as $h \rightarrow \infty$ to ensure stationarity of y_t , G_k 's do not need to decay thanks to the diminishing $\ell_{h,k}(\omega)$'s. Consequently, G_k 's, which have relatively strong signals, can be easier to interpret than the diminishing A_h 's. Second, similar to the orders (p, q) of VARMA models, the required d is generally small in practice. For example, d = 2 works well for the macroeconomic data in Section 6, so we only need to interpret two adjacency matrices G_1 and G_2 . However, if the VAR(P) model were fitted, we would have to interpret P adjacency matrices, where the required *P* would be much larger.

We summarize the main contributions of this article as follows:

- (i) A sparse parametric $VAR(\infty)$ model is introduced for highdimensional time series, which is favorable regarding (I1)-(I3), while avoiding problems (P1) and (P2).
- (ii) We develop two ℓ_1 -regularized estimators, which can be implemented via efficient block coordinate descent algorithms, and derive their nonasymptotic error bounds under weak sparsity; particularly, our theory takes into account the effect of initializing $y_t = 0$ for $t \le 0$, which is needed for feasible estimation of VAR(∞) models.
- (iii) A high-dimensional Bayesian information criterion (BIC) is proposed for model order selection, and its consistency is established.

The remainder of this article is organized as follows. Section 2 introduces the proposed model and its interpretation. Section 3 presents two ℓ_1 -regularized estimators and their nonasymptotic theory. Section 4 introduces the proposed BIC. Sections 5 and 6 provide simulation and empirical studies. Section 7 concludes with a brief discussion. The block coordinate descent algorithms for implementing the estimation, additional simulation and empirical results, and all technical proofs are provided in a separate supplementary file.

Unless otherwise specified, we denote scalars, vectors and matrices by lowercase letters (e.g., x), boldface lowercase letters (e.g., x), and boldface capital letters (e.g., X), respectively. Let $\mathbb{I}_{\{\cdot\}}$ be the indicator function taking value one when the condition is true and zero otherwise. For any $a, b \in \mathbb{R}$, let $a \vee b = \max\{a, b\}$ and $a \wedge b = \min\{a, b\}$. The ℓ_q -norm of any $\mathbf{x} \in \mathbb{R}^p$ is denoted by $\|x\|_q = (\sum_{i=1}^p |x_i|^q)^{1/q}$ for q > 0. For any $X \in \mathbb{R}^{d_1 \times d_2}$, let X^{\top} , $\sigma_{\max}(X)$ (or $\sigma_{\min}(X)$), $\lambda_{\max}(X)$ (or $\lambda_{\min}(X)$), $\operatorname{vec}(X)$, $||X||_{\operatorname{op}}$, and $||X||_F$ be its transpose, largest (or smallest) singular value, largest (or smallest) eigenvalue, vectorization, operator norm $\|X\|_{\text{op}} = \sigma_{\text{max}}(X)$, and Frobenius norm $\|X\|_{\text{F}} = \sqrt{\text{tr}(X^{\top}X)}$, respectively. We use C > 0 (or c > 0) to denote generic large (or

small) absolute constants. For any sequences x_n and y_n , denote $x_n \lesssim y_n$ (or $x_n \gtrsim y_n$) if there is C > 0 such that $x_n \leq Cy_n$ (or $x_n \ge Cy_n$). We write $x_n \asymp y_n$ if $x_n \lesssim y_n$ and $x_n \gtrsim y_n$. In addition, $x_n \gg y_n$ if $y_n/x_n \to 0$ as $n \to \infty$.

2. Proposed Model

2.1. Motivation: Reparameterization of VARMA Models

This section introduces the motivation behind the proposed model. Recall that the shared root cause of problems (P1) and (P2) of the VARMA(1, 1) model, as discussed in Section 1, lies in the matrix multiplications involved in computing the AR coefficient matrices $A_h(\Phi, \Theta) = \Theta^{h-1}(\Phi - \Theta)$ in the VAR(∞) form of the model. Thus, the key to overcoming both problems is to eliminate the matrix multiplications in the parameterization

To this end, we show that a reparameterization of $A_h(\Phi, \Theta)$ free of matrix multiplications can be derived via the following two main steps: (a) Block-diagonalize Θ via the Jordan decomposition, $\Theta = BJB^{-1}$, where $B \in \mathbb{R}^{N \times N}$ is an invertible matrix, and $I \in \mathbb{R}^{N \times N}$ is the real Jordan form containing eigenvalues of Θ ; see (2.1) for details. (b) Then, merge **B** with all remaining components in the expression of $A_h(\Phi, \Theta)$.

Specifically, by Theorem 1 in Hartfiel (1995), for any 0 < $n \leq N$, real matrices with n distinct nonzero eigenvalues are dense in the set of all $N \times N$ real matrices with rank at most n. Thus, with only a little loss of generality, we can assume that Θ is a real matrix with n distinct nonzero eigenvalues, where n = $rank(\Theta)$; a more general result allowing repeated eigenvalues is derived in the technical appendix of Huang, Lu, and Zheng (2023). Then suppose that Θ has r nonzero real eigenvalues, $\lambda_1, \ldots, \lambda_r$, and s conjugate pairs of nonzero complex eigenvalues, $(\lambda_{r+2m-1}, \lambda_{r+2m}) = (\gamma_m e^{i\theta_m}, \gamma_m e^{-i\theta_m})$ for $1 \le m \le s$, where $|\lambda_j| \in (0,1)$ for $1 \le j \le r$, $\gamma_m \in (0,1)$ and $\theta_m \in (0,\pi)$ for $1 \le m \le s$, and *i* represents the imaginary unit. Therefore, n = r + 2s, and the real Jordan form of Θ is a real block diagonal matrix:

$$J = \operatorname{diag} \{\lambda_1, \dots, \lambda_r, C_1, \dots, C_s, \mathbf{0}\},$$

$$C_m = \gamma_m \cdot \begin{pmatrix} \cos \theta_m & \sin \theta_m \\ -\sin \theta_m & \cos \theta_m \end{pmatrix} \in \mathbb{R}^{2 \times 2},$$
(2.1)

where $1 \le m \le s$; see chap. 3 in Horn and Johnson (2012).

Let $A_1 = \Phi - \Theta := G_1$. Substituting the Jordan decomposition $\Theta = BJB^{-1}$ into the expression of A_h , we can show that for all $h \ge 2$, $A_h = BJ^{h-1}B^{-1}(\Phi - \Theta) = \sum_{j=1}^r \lambda_j^{h-1}G_{1+j} +$ $\sum_{m=1}^{s} \gamma_m^{h-1} \left[\cos\{(h-1)\theta_m\} G_{1+r+2m-1} + \sin\{(h-1)\theta_m\} G_{1+r+2m} \right], \text{ where } G_2, \dots, G_{1+r+2s} \in \mathbb{R}^{N \times N} \text{ are determined}$ jointly by **B** and $B^{-1}(\Phi - \Theta)$; see the proof of Proposition 1 in the supplementary file for details. This result is a reparameterization of A_h 's in terms of the scalars λ_i 's, γ_m 's, θ_m 's, and matrices G_1, \ldots, G_{1+r+2s} . As each A_h is a linear combination of G_1, \ldots, G_{1+r+2s} , problems (P1) and (P2) are tackled at their root: It not only ensures the identifiability of the parameters λ_j 's, γ_m 's, θ_m 's, and the **G**-matrices, up to a permutation in the indices j and m, but also leads to a significantly reduced computational complexity, such as $O(TN^2 + T^2N)$ for the squared loss function.

In general, the VARMA(p,q) model is given by $\mathbf{y}_t = \sum_{i=1}^p \mathbf{\Phi}_i \mathbf{y}_{t-i} + \boldsymbol{\varepsilon}_t - \sum_{j=1}^q \mathbf{\Theta}_j \boldsymbol{\varepsilon}_{t-j}$, where $\mathbf{\Phi}_i, \mathbf{\Theta}_j \in \mathbb{R}^{N \times N}$ for $1 \le i \le p$ and $1 \le j \le q$. Assuming invertibility, it has the following VAR(∞) representation:

$$\mathbf{y}_{t} = \sum_{h=1}^{\infty} \left(\sum_{i=0}^{p \wedge h} \mathbf{P} \underline{\mathbf{\Theta}}^{h-i} \mathbf{P}^{\top} \mathbf{\Phi}_{i} \right) \mathbf{y}_{t-h} + \boldsymbol{\varepsilon}_{t}, \tag{2.2}$$

$$\underline{\boldsymbol{\Theta}} = egin{pmatrix} \boldsymbol{\Theta}_1 & \boldsymbol{\Theta}_2 & \cdots & \boldsymbol{\Theta}_{q-1} & \boldsymbol{\Theta}_q \\ I & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I & \cdots & \mathbf{0} & \mathbf{0} \\ dots & dots & \ddots & dots & dots \\ \mathbf{0} & \mathbf{0} & \cdots & I & \mathbf{0} \end{pmatrix},$$

where $\Phi_0 = -I$ and $P = (I_N, \mathbf{0}_{N \times N(q-1)})$ are constant matrices, $\underline{\Theta}$ is called the MA companion matrix, and all eigenvalues of $\underline{\Theta}$ are less than one in absolute value; see Lütkepohl (2005). Similar to the VARMA(1, 1) case, the following reparameterization can be derived.

Proposition 1. Suppose that all nonzero eigenvalues of $\underline{\Theta}$ are distinct, and there are r distinct nonzero real eigenvalues of $\underline{\Theta}$, $\lambda_j \in (-1,0) \cup (0,1)$ for $1 \leq j \leq r$, and s distinct conjugate pairs of nonzero complex eigenvalues of $\underline{\Theta}$, $(\lambda_{r+2m-1}, \lambda_{r+2m}) = (\gamma_m e^{i\theta_m}, \gamma_m e^{-i\theta_m})$ with $\gamma_m \in (0,1)$ and $\theta_m \in (0,\pi)$ for $1 \leq m \leq s$. Then for all $h \geq 1$, we have

$$A_{h} = \sum_{k=1}^{p} \mathbb{I}_{\{h=k\}} G_{k} + \sum_{j=1}^{r} \mathbb{I}_{\{h \geq p+1\}} \lambda_{j}^{h-p} G_{p+j}$$

$$+ \sum_{m=1}^{s} \mathbb{I}_{\{h \geq p+1\}} \gamma_{m}^{h-p} \left[\cos\{(h-p)\theta_{m}\} G_{p+r+2m-1} + \sin\{(h-p)\theta_{m}\} G_{p+r+2m} \right],$$
(2.3)

where $G_k = A_k$ for $1 \le k \le p$, and $\{G_k\}_{k=p+1}^{p+r+2s}$ are determined jointly by \widetilde{B} and \widetilde{B}_- , with $\widetilde{B} = PB$ and $\widetilde{B}_- = B^{-1} \left(\sum_{i=0}^{p} \underline{\Theta}^{p-i} P^{\top} \Phi_i \right)$. In addition, the corresponding term in (2.3) is suppressed if p, r or s is zero.

Throughout this article, we denote d=p+r+2s. Let $\boldsymbol{\omega}=(\lambda_1,\ldots,\lambda_r,\eta_1^\top,\ldots,\eta_s^\top)^\top\in\mathbb{R}^{r+2s}$, where $\eta_m=(\gamma_m,\theta_m)^\top$ for $1\leq m\leq s$, and $\boldsymbol{g}=\operatorname{vec}(\boldsymbol{G})\in\mathbb{R}^{N^2d}$, where $\boldsymbol{G}=(\boldsymbol{G}_1,\ldots,\boldsymbol{G}_d)\in\mathbb{R}^{N\times Nd}$. Then, we can succinctly write (2.3) in the parametric form of $\boldsymbol{A}_h=\boldsymbol{A}_h(\boldsymbol{\omega},\boldsymbol{g})=\sum_{k=1}^d\ell_{h,k}(\boldsymbol{\omega})\boldsymbol{G}_k$ for all $h\geq 1$. Here $\ell_{h,k}(\cdot)$'s are real-valued functions predetermined according to (2.3), which can be defined conveniently through a matrix as follows: for any $h\geq 1$ and $1\leq k\leq d$, $\ell_{h,k}(\boldsymbol{\omega})$ is the (h,k)th entry of the $\infty\times d$ matrix,

$$L(\boldsymbol{\omega}) = (\ell_{h,k}(\boldsymbol{\omega}))_{h \geq 1,1 \leq k \leq d}$$

$$= \begin{pmatrix} \mathbf{I}_p & \mathbf{0}_{p \times 1} & \cdots & \mathbf{0}_{p \times 1} & \mathbf{0}_{p \times 2} & \cdots & \mathbf{0}_{p \times 2} \\ \mathbf{0}_{\infty \times p} & \boldsymbol{\ell}^I(\lambda_1) & \cdots & \boldsymbol{\ell}^I(\lambda_r) & \boldsymbol{\ell}^{II}(\boldsymbol{\eta}_1) & \cdots & \boldsymbol{\ell}^{II}(\boldsymbol{\eta}_s) \end{pmatrix}$$

$$\in \mathbb{R}^{\infty \times d},$$

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

$$\boldsymbol{\ell}^{I}(\lambda) = (\lambda, \lambda^{2}, \lambda^{3}, \dots)^{\top} \in \mathbb{R}^{\infty},$$

$$\boldsymbol{\ell}^{II}(\boldsymbol{\eta}) = \begin{pmatrix} \gamma \cos(\theta) & \gamma^2 \cos(2\theta) & \gamma^3 \cos(3\theta) & \cdots \\ \gamma \sin(\theta) & \gamma^2 \sin(2\theta) & \gamma^3 \sin(3\theta) & \cdots \end{pmatrix}^{\top}$$
$$\in \mathbb{R}^{\infty \times 2}.$$

2.2. Proposed Sparse Parametric $VAR(\infty)$ Model

Motivated by the discussion in Section 2.1, we propose the following $VAR(\infty)$ model for high-dimensional time series:

$$y_{t} = \sum_{h=1}^{\infty} A_{h}(\boldsymbol{\omega}, \boldsymbol{g}) y_{t-h} + \boldsymbol{\varepsilon}_{t}$$

$$= \sum_{k=1}^{d} G_{k} \sum_{h=1}^{\infty} \ell_{h,k}(\boldsymbol{\omega}) y_{t-h} + \boldsymbol{\varepsilon}_{t},$$
(2.4)

where $\omega \in (-1,1)^r \times \Pi^s \subset \mathbb{R}^{r+2s}$ is a parameter vector, with $\Pi = [0,1) \times (0,\pi)$, $\ell_{h,k}(\cdot)$'s are known real-valued functions defined as in Section 2.1, $G_k \in \mathbb{R}^{N \times N}$ for $1 \leq k \leq d$ are parameter matrices with d = p + r + 2s. To handle the high-dimensionality, we assume that G_k 's are sparse matrices. In this section, we will focus on the exact sparsity as it is instrumental for model interpretability. However, it will be relaxed to weak sparsity in our theoretical analysis; see Assumptions 4 and 4' in Section 3. We call model (2.4) with exactly or weakly sparse G_k 's the Sparse Parametric VAR(∞) (SPVAR(∞)) model.

Note that if no sparsity assumption is imposed on G_k 's, then (2.4) provides an alternative low-dimensional time series model comparable to the VARMA model; see Section 2.3 for its stationarity condition. While formulation (2.4) is derived from the VARMA model, it is worth clarifying that it relaxes the restrictions on G_{p+j} for $1 \le j \le r + 2s$. Specifically, by Proposition 1, if $\{y_t\}$ is indeed generated from a VARMA model, then G_{p+j} 's would fulfill certain restrictions as determined by the Jordan decomposition of the MA companion matrix $\underline{\Theta}$. By contrast, (2.4) treats these matrices as free parameters.

The resemblance between (2.4) and the VARMA model is mainly achieved by $\ell_{h,k}(\cdot)$'s, which yield VARMA-type decay patterns of A_h as $h \to \infty$. According to (2.3), $\ell_{h,k}(\cdot)$'s implicitly depend on the orders (p,r,s). Note that p and (r,s) are counterparts of the AR and MA orders of the VARMA model, respectively. In fact, when r=s=0, (2.4) reduces to the VAR(p) model, $y_t=\sum_{h=1}^p G_h y_{t-h} + \varepsilon_t$. For this reason, we call G_1,\ldots,G_p and G_{p+1},\ldots,G_d the AR and MA coefficient matrices of the model, respectively. While larger (p,r,s) allow for more complex temporal patterns, similar to the VARMA model, usually it suffices to use small orders in practice; see Section 6 for empirical evidence.

The proposed model can be directly used to infer the multivariate Granger causality (MGC), which concerns Granger causal (GC) relations (Granger 1969) between any pair of component series in $y_t = (y_{1,t}, \ldots, y_{N,t})^{\top}$; see Shojaie and Fox (2021) for an excellent review. By definition, $\{y_{j,t}\}$ is GC for $\{y_{i,t}\}$ if the past information of $y_{j,t}$ can improve the forecast of $y_{i,t}$, where $1 \leq i \neq j \leq N$. Most existing works study the MGC under the finite-order VAR for its convenience: Under the model $y_t = \sum_{h=1}^{P} A_h y_{t-h} + \varepsilon_t$, $\{y_{j,t}\}$ is GC for $\{y_{i,t}\}$ if $a_{i,j,h} \neq 0$ for some $h \in \{1, \ldots, P\}$, where $a_{i,j,h}$ is the (i,j)th entry of A_h , for $1 \leq i \neq j \leq N$. Notably, while working with A_h 's would be

$\{y_{2,t}\}$ is not Granger Causal for $\{y_{1,t}\}$		$\{y_{2,t}\}$ is Granger Causal for $\{y_{1,t}\}$					
		(1) Influence at lag 1 only		(2) Influence at all lags ≥ 2		(3) Influence across all lags	
G_1	G_2	G_1	G_2	G_1	G_2	G_1	G_2
0	0	X	0	0	X	X	X

Figure 1. Illustration for different scenarios of Granger causality of $\{y_{2,t}\}$ for $\{y_{1,t}\}$ when (p,r,s)=(1,1,0) and N=3, as determined by the (1,2)th entry of G_1 and G_2 . Cell (1,2) of G_k is marked with "0" when $g_{1,2,k}=0$, and "X" when $g_{1,2,k}\neq 0$.

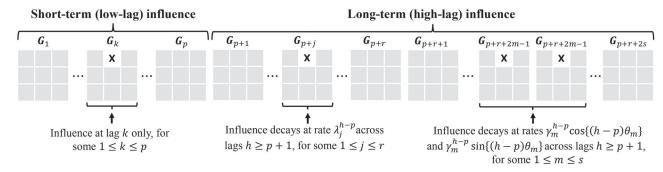


Figure 2. Illustration for different types of lagged influence of $\{y_{2,t}\}$ on $\{y_{1,t}\}$ under general orders (p,r,s) and N=3. Cell (1,2) of G_k is marked with "X" when $g_{1,2,k}\neq 0$.

infeasible when $P = \infty$, we can directly infer the MGC through G_k 's: By (2.4), we have that $\{y_{j,t}\}$ is GC for $\{y_{i,t}\}$ if $g_{i,j,k} \neq 0$ for some $k \in \{1, \ldots, d\}$, where $g_{i,j,k}$ is the (i,j)th entry of G_k , for $1 \leq i \neq j \leq N$; see Figure 1 for an illustration with (i,j) = (1,2), (p,r,s) = (1,1,0), and N=3.

More interestingly, since each G_k captures a piece of cross-sectional information associated with a particular sequence $\{\ell_{h,k}(\boldsymbol{\omega})\}_{h=1}^{\infty}$, we can discern the decay pattern of any GC relations over time, achieving a more granular understanding of the MGC. For simplicity, consider the model for $y_{1,t}$ when (p,r,s)=(1,1,0): $y_{1,t}=\sum_{j=1}^N g_{1,j,1}y_{j,t-1}+\sum_{j=1}^N g_{1,j,2}\sum_{h=2}^\infty \lambda^{h-1}y_{j,t-h}+\varepsilon_{1,t}$, where $g_{i,j,k}$ denotes the (i,j)th entry of G_k . First, it is clear that $\{y_{j,t}\}$ is GC for $\{y_{1,t}\}$ if $g_{1,j,1}$ and $g_{1,j,2}$ are not both zero. Second, if this GC relation exists, the lagged influence of $\{y_{j,t}\}$ on $\{y_{1,t}\}$ can be classified into the following three scenarios: (a) lag-one only, if $g_{1,j,1}\neq 0$ and $g_{1,j,2}=0$; (b) all lags beyond lag one, if $g_{1,j,1}=0$ and $g_{1,j,2}\neq 0$. In scenarios (b) and (c), the exponential decay of the influence over time is determined by λ ; see Figure 1 for an illustration for j=2.

In general, with orders (p,r,s), the model equation for $y_{1,t}$ will consist of two conditional mean terms: The first term involves the sum of $g_{1,j,k}y_{j,t-k}$ for lags $1 \le k \le p$, whereas the second term captures the influence beyond lag p. The latter involves a weighted mixture of r distinct exponential decay rates and s distinct pairs of damped cosine and sine waves. Then the lagged influence of $\{y_{j,t}\}$ on $\{y_{1,t}\}$ can be generalized to the following three scenarios, if the GC relation exists: (1) short-term only, if $g_{1,j,k} \ne 0$ for some $1 \le k \le p$, while $g_{1,j,p+1} = \cdots = g_{1,j,d} = 0$; (2) long-term only, if $g_{1,j,1} = \cdots = g_{1,j,p} = 0$, while $g_{1,j,k} \ne 0$ for some $p+1 \le k \le d$; and (3) both short-term and long-term influences, if $g_{1,j,k} \ne 0$ for some $1 \le k \le p$ and some $p+1 \le k \le d$. A more detailed illustration is given in Figure 2.

Remark 1. In many applications, the cross-sectional dependence may not be time-invariant; for example, Barigozzi and Brownlees (2017) found that the estimated Granger causal network in a sparse VAR system for stock volatilities may be time-varying. Time-varying cross-sectional dependence is also common in behavioral and neural studies: for example, different segments of video time series of freely moving animals may correspond to distinct behaviors (Costacurta et al. 2022), and discrete shifts in the dynamics of neural activity may reflect changes in underlying brain state (Fiecas et al. 2023). To accommodate such applications, the proposed model can be extended to allow G_k 's to be time varying; for example, a Markov-switching SPVAR(∞) model may be developed along the lines of Li, Safikhani, and Shojaie (2022).

Remark 2. In VAR models, the GC relations as captured by the coefficient matrices A_h 's correspond to lagged cross-sectional dependence, whereas the instantaneous cross-sectional dependence is captured by the variance-covariance matrix Σ_{ε} of ε_t . While this section focuses on the former, Σ_{ε} can also be estimated based on residuals from the fitted SPVAR(∞) model; see Remark 5 in Section 3.1.

Remark 3. We can also conduct impulse response analysis based on the VMA(∞) form of the proposed model; see Theorem 1 in Section 2.3 for the VMA(∞) representation. For example, when (p,r,s)=(1,1,0), the corresponding MA coefficient matrices are $\Psi_1=G_1, \Psi_2=G_1^2+\lambda G_2, \Psi_3=G_1^3+\lambda G_1G_2+\lambda G_2G_1+\lambda^2G_2$, etc. When G_1 and G_2 are both sparse with their nonzero entries in sufficiently different positions, all Ψ_j 's will also tend to be sparse; this is indeed the case for the empirical example in Section 6. Thus, we can alternatively interpret the high-dimensional time series via the impulse response analysis.

2.3. Stationarity Condition

We provide a sufficient condition on ω and G_k 's for the existence of a unique strictly stationary solution for (2.4) in the following theorem, which is valid whether G_k 's are sparse or not. Similar to the AR companion matrix of a VARMA(p, q) model, denote

$$\underline{G}_1 = \begin{pmatrix} G_1 & G_2 & \cdots & G_{p-1} & G_p \\ I & 0 & \cdots & 0 & 0 \\ 0 & I & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & I & 0 \end{pmatrix}.$$

Theorem 1. Suppose that there exists $0 < \bar{\rho} < 1$ such that

$$\max\{|\lambda_1|,\ldots,|\lambda_r|,\gamma_1,\ldots,\gamma_s\} \leq \bar{\rho}$$
 and

$$\rho(\underline{G}_1) + \frac{\bar{\rho}}{1 - \bar{\rho}} \sum_{k=1}^{r+2s} \rho(G_{p+k}) < 1,$$

where $\rho(\cdot)$ denotes the spectral radius of a matrix, and $\rho(\underline{G}_1)$ disappears when p=0. Moreover, $\{\boldsymbol{\varepsilon}_t\}$ is a strictly stationary sequence. Then there exists a unique strictly stationary solution to the model equation in (2.4), given by $\boldsymbol{y}_t = \boldsymbol{\varepsilon}_t + \sum_{j=1}^{\infty} \boldsymbol{\Psi}_j \boldsymbol{\varepsilon}_{t-j}$, where $\boldsymbol{\Psi}_j = \sum_{k=1}^{\infty} \sum_{j_1+\dots+j_k=j} \boldsymbol{A}_{j_1} \cdots \boldsymbol{A}_{j_k}$ for $j \geq 1$, with $\boldsymbol{A}_h = \sum_{k=1}^d \ell_{h,k}(\boldsymbol{\omega}) \boldsymbol{G}_k$ for $h \geq 1$.

When r=s=0, the condition in Theorem 1 reduces to $\rho(\underline{G}_1)<1$, which coincides with the necessary and sufficient condition for the strict stationarity of the VAR(p) model. When r and s are not both zero, the stationarity region for G_k 's in Theorem 1 will be larger if $\bar{\rho}$ becomes smaller, that is, if A_h diminishes more quickly as $h\to\infty$.

Remark 4. If $\{y_t\}$ is a VARMA(p,q) process fulfilling the representation in (2.4), it is known that the necessary and sufficient condition for its strict stationarity is simply $\rho(\underline{G}_1) < 1$; see Lütkepohl (2005). This suggests that the sufficient condition in Theorem 1 could sometimes be restrictive. Indeed, the condition on ω and G_k 's in Theorem 1 is derived from the necessary and sufficient condition: $\sum_{j=1}^{\infty} \|\Psi_j\| < \infty$, where Ψ_j 's are functions of A_h 's as defined in the VMA (∞) form of $\{y_t\}$ in Theorem 1, and $\|\cdot\|$ is any submultiplicative matrix norm. This motivates us to recommend a more general numerical method to check stationarity for practical use: first compute the sequence $\{\Psi_j\}$ using the parameters ω and G_k 's, and then numerically check whether the partial sum $\sum_{j=1}^J \|\Psi_j\|$ converges as $J \to \infty$. This method is applied in Section 6 to check the stationarity of the fitted model.

3. High-Dimensional Estimation

3.1. ℓ_1 -Regularized Joint Estimator

We first propose an ℓ_1 -regularized estimator for the SPVAR(∞) model via jointly fitting all component series of y_t . An alternative estimator will be introduced in the next section.

For $\{y_t\}_{t=1}^T$ generated from (2.4) with orders (p, r, s), the squared loss is $\mathbb{L}_T(\omega, \mathbf{g}) = T^{-1} \sum_{t=1}^T \| \mathbf{y}_t - \sum_{h=1}^\infty \mathbf{A}_h(\omega, \mathbf{g}) \mathbf{y}_{t-h} \|_2^2 = T^{-1} \sum_{t=1}^T \| \mathbf{y}_t - \sum_{k=1}^d \mathbf{G}_k \sum_{h=1}^\infty$

 $\ell_{h,k}(\boldsymbol{\omega})\boldsymbol{y}_{t-h}\|_2^2$. Here $\boldsymbol{g}=\operatorname{vec}(\boldsymbol{G})$, where $\boldsymbol{G}=(\boldsymbol{G}_1,\ldots,\boldsymbol{G}_d)\in\mathbb{R}^{N\times Nd}$. Since the loss function depends on observations in the infinite past, initial values for $\{\boldsymbol{y}_t,t\leq 0\}$ will be needed in practice. We set them to zero as $\mathbb{E}(\boldsymbol{y}_t)=\boldsymbol{0}$, and then the corresponding loss becomes

$$\widetilde{\mathbb{L}}_{T}(\boldsymbol{\omega}, \mathbf{g}) = \frac{1}{T} \sum_{t=1}^{T} \left\| \mathbf{y}_{t} - \sum_{h=1}^{t-1} A_{h}(\boldsymbol{\omega}, \mathbf{g}) \mathbf{y}_{t-h} \right\|_{2}^{2}$$

$$= \frac{1}{T} \sum_{t=1}^{T} \left\| \mathbf{y}_{t} - \sum_{k=1}^{d} G_{k} \sum_{h=1}^{t-1} \ell_{h,k}(\boldsymbol{\omega}) \mathbf{y}_{t-h} \right\|_{2}^{2}.$$
(3.1)

The initialization effect will be taken into account in our theoretical analysis, and its negligibility is confirmed by our simulation study; see Lemmas S6–S8 and Section S2 in the supplementary file. We propose the ℓ_1 -regularized joint estimator (JE) as follows:

$$(\widehat{\boldsymbol{\omega}}, \widehat{\boldsymbol{g}}) = \arg\min_{\boldsymbol{\omega} \in \mathbf{\Omega}, \boldsymbol{g} \in \mathbb{R}^{N^2 d}} \left\{ \widetilde{\mathbb{L}}_T(\boldsymbol{\omega}, \boldsymbol{g}) + \lambda_{\boldsymbol{g}} \| \boldsymbol{g} \|_1 \right\}, \quad (3.2)$$

where $\lambda_g > 0$ is the regularization parameter, and $\Omega \subset (-1,1)^r \times \Pi^s$ denotes the parameter space of ω . Let a = vec(A), where $A = (A_1,A_2,\ldots)$ is the horizontal concatenation of $\{A_h\}_{h=1}^{\infty}$. Note that $a = (L(\omega) \otimes I_{N^2})g$. Based on (3.2), the estimator of A_h is $\widehat{A}_h = \sum_{k=1}^d \ell_{h,k}(\widehat{\omega})\widehat{G}_k$ for $h \geq 1$. Then, $\widehat{a} = \text{vec}(\widehat{A}) = (L(\widehat{\omega}) \otimes I_{N^2})\widehat{g}$, where $\widehat{A} = (\widehat{A}_1, \widehat{A}_2, \ldots)$.

Denote the true value of any parameter with the superscript "*", for example, \mathbf{g}^* , $\boldsymbol{\omega}^*$, and \mathbf{a}^* . For $\boldsymbol{\omega}^* \in \Omega$, let $\nu_{\text{lower}}^* = (\min_{1 \leq j \leq r} |\lambda_j^*|) \wedge (\min_{1 \leq m \leq s} |\gamma_m^*|)$ and $\nu_{\text{gap}}^* = \min_{1 \leq j \neq k \leq r+2s} |x_j^* - x_k^*|$, where $x_j^* = \lambda_j^*$ for $1 \leq j \leq r$ and $(x_{r+2m-1}^*, x_{r+2m}^*) = (\gamma_m^* e^{i\theta_m^*}, \gamma_m^* e^{-i\theta_m^*})$ for $1 \leq m \leq s$. The assumptions for our theoretical analysis are presented as follows.

Assumption 1 (Parameter space and stationarity). (i) There exists an absolute constant $0 < \bar{\rho} < 1$ such that $|\lambda_1|, \ldots, |\lambda_r|, \gamma_1, \ldots, \gamma_s \leq \bar{\rho}$ for all $\omega \in \Omega$; and (ii) the time series $\{y_t\}$ is stationary.

Assumption 2 (Separability). (i) There exists an absolute constant $c_{\nu} > 0$ such that $\nu_{\text{lower}}^* \ge c_{\nu}$ and $\nu_{\text{gap}}^* \ge c_{\nu}$; and (ii) r and s are fixed.

Assumption 3 (Sub-Gaussian errors). Let $\boldsymbol{\varepsilon}_t = \boldsymbol{\Sigma}_{\varepsilon}^{1/2} \boldsymbol{\xi}_t$, where $\boldsymbol{\xi}_t$ is a sequence of i.i.d. random vectors with zero mean and $\operatorname{var}(\boldsymbol{\xi}_t) = \boldsymbol{I}_N$, and $\boldsymbol{\Sigma}_{\varepsilon}$ is a positive definite covariance matrix. In addition, the coordinates $(\boldsymbol{\xi}_{it})_{1 \leq i \leq N}$ within $\boldsymbol{\xi}_t$ are mutually independent and σ^2 -sub-Gaussian.

Assumption 1(i) ensures that $|\lambda_j|$'s and γ_m 's are bounded away from one. A sufficient condition for Assumption 1(ii) is given in Theorem 1. Under stationarity, $\{y_t\}$ has the VMA(∞) form $y_t = \Psi_*(B)\varepsilon_t$, where $\Psi_*(B) = I_N + \sum_{j=1}^{\infty} \Psi_j^* B^j$, and B is the backshift operator; see Theorem 1. Let $\mu_{\min}(\Psi_*) = \min_{|z|=1} \lambda_{\min}(\Psi_*(z)\Psi_*^{\mathsf{H}}(z))$ and $\mu_{\max}(\Psi_*) = \max_{|z|=1} \lambda_{\max}(\Psi_*(z)\Psi_*^{\mathsf{H}}(z))$, where $\Psi_*^{\mathsf{H}}(z)$ is the conjugate transpose of $\Psi_*(z)$ for $z \in \mathbb{C}$. It can be verified that $\mu_{\min}(\Psi_*) > 0$; see also Basu and Michailidis (2015). Then we define the positive constants $\kappa_1 = \lambda_{\min}(\Sigma_\varepsilon)\mu_{\min}(\Psi_*)$ and $\kappa_2 = \lambda_{\max}(\Sigma_\varepsilon)\mu_{\max}(\Psi_*)$. Assumption 2(i) requires that

different λ_j^* 's or η_m^* 's are bounded away from zero and from each other. Since these parameters lie in bounded parameter spaces, this also entails that r and s must be fixed; see Assumption 2(ii). Assumption 3 relaxes the Gaussian assumption commonly used in the literature on high-dimensional time series models (e.g., Basu and Michailidis 2015) to sub-Gaussianity.

Let $\mathbf{g}_{AR} = \text{vec}(\mathbf{G}_{AR})$ and $\mathbf{g}_{MA} = \text{vec}(\mathbf{G}_{MA})$, where $\mathbf{G}_{AR} = (\mathbf{G}_1, \dots, \mathbf{G}_p) \in \mathbb{R}^{N \times Np}$ and $\mathbf{G}_{MA} = (\mathbf{G}_{p+1}, \dots, \mathbf{G}_d) \in \mathbb{R}^{N \times N(r+2s)}$. Let $g_{i,j,k}$ be the (i,j)th entry of \mathbf{G}_k . Then, we define the weak sparsity of \mathbf{g}_{AR}^* and \mathbf{g}_{MA}^* by restricting them into the ℓ_q -"balls", $\mathbb{B}_q(R_q^{AR}) := \{\mathbf{g}_{AR} \in \mathbb{R}^{N^2p} \mid \sum_{k=1}^p \sum_{i=1}^N \sum_{j=1}^N |g_{i,j,k}|^q \leq R_q^{AR}\}$ and $\mathbb{B}_q(R_q^{MA}) := \{\mathbf{g}_{MA} \in \mathbb{R}^{N^2(r+2s)} \mid \sum_{k=p+1}^d \sum_{i=1}^N \sum_{j=1}^N |g_{i,j,k}|^q \leq R_q^{MA}\}$, respectively, which is a more general assumption than exact sparsity.

Assumption 4 (Weak sparsity). There exists $q \in [0,1]$ such that $\mathbf{g}_{\mathrm{AR}}^* \in \mathbb{B}_q(R_q^{\mathrm{AR}})$ and $\mathbf{g}_{\mathrm{MA}}^* \in \mathbb{B}_q(R_q^{\mathrm{MA}})$ for some radii $R_q^{\mathrm{AR}}, R_q^{\mathrm{MA}} > 0$.

Assumption 4 implies that $\mathbf{g}^* \in \mathbb{B}_q(R_q)$, where $R_q := R_q^{\mathrm{AR}} + R_q^{\mathrm{MA}}$ and $\mathbb{B}_q(R_q) := \{\mathbf{g} \in \mathbb{R}^{N^2d} \mid \sum_{k=1}^d \sum_{i=1}^N \sum_{j=1}^N |g_{i,j,k}|^q \leq R_q \}$. If q = 0, Assumption 4 becomes the exact sparsity constraints— $\mathbf{g}_{\mathrm{AR}}^*$ and $\mathbf{g}_{\mathrm{MA}}^*$ have at most R_q^{AR} and R_q^{MA} nonzero entries, respectively. If $q \in (0,1]$, the ℓ_q -"balls" enforce a certain decay rate on the absolute values of the entries in \mathbf{g}^* as the dimension N grows. Note that we do not require R_q^{AR} and R_q^{MA} to be fixed.

A main theoretical challenge is that the loss function $\mathbb{L}_T(\boldsymbol{\omega}, \boldsymbol{g})$ is highly nonconvex with respect to $\boldsymbol{\omega}$. Consequently, the global statistical consistency commonly established for highdimensional convex M-estimators is not available. However, if the nonconvex loss function exhibits a benign convex curvature over local regions, then a form of local statistical consistency can be established; see, for example, Loh (2017). For many nonconvex M-estimators, certain convexity holds within a constantradius neighborhood of the true parameter value; for the highdimensional setup, this is termed as local restricted strong convexity in Loh (2017). Then it can be shown that all local optima within this region can enjoy the same convergence rate as the ℓ_1 -regularized least squared estimator for linear regression; see also Janková and van de Geer (2021) and Wang and He (2022) for other works on local statistical guarantees for estimators with nonconvex losses or regularizers. Our method is reminiscent of that for high-dimensional nonconvex M-estimators in the literature. However, our setting is special in that $\mathbb{L}_T(\boldsymbol{\omega}, \boldsymbol{g})$ is only partially nonconvex, as it is convex with respect to g, for any fixed ω . Thus, unlike Loh (2017), we only need to restrict ω within a local region of restricted curvature around ω^* , while **g** can be free.

Let $\underline{\alpha}_{\mathrm{MA}} = \min_{1 \leq j \leq r+2s} \| \mathbf{G}_{p+j}^* \|_{\mathrm{F}}$ and $\overline{\alpha}_{\mathrm{MA}} = \max_{1 \leq j \leq r+2s} \| \mathbf{G}_{p+j}^* \|_{\mathrm{F}}$, which are both allowed to grow with N. Then let $\alpha = \overline{\alpha}_{\mathrm{MA}}/\underline{\alpha}_{\mathrm{MA}}$. The local convexity of our loss function around $\boldsymbol{\omega}^*$ is an immediate consequence of the following proposition.

Proposition 2. Suppose that $\underline{\alpha}_{MA} > 0$. Then under Assumptions 1(i) and 2, there exists a constant $c_{\omega} = \min(2, c/\alpha) > 0$

such that for any $\boldsymbol{\omega} \in \boldsymbol{\Omega}$ with $\|\boldsymbol{\omega} - \boldsymbol{\omega}^*\|_2 \leq c_{\boldsymbol{\omega}}$, it holds $\|\boldsymbol{g} - \boldsymbol{g}^*\|_2 + \underline{\alpha}_{\mathrm{MA}} \|\boldsymbol{\omega} - \boldsymbol{\omega}^*\|_2 \lesssim \|\boldsymbol{a} - \boldsymbol{a}^*\|_2^2 \lesssim \|\boldsymbol{g} - \boldsymbol{g}^*\|_2 + \overline{\alpha}_{\mathrm{MA}} \|\boldsymbol{\omega} - \boldsymbol{\omega}^*\|_2$, where $\boldsymbol{a} = (\boldsymbol{L}(\boldsymbol{\omega}) \otimes \boldsymbol{I}_{N^2}) \boldsymbol{g}$.

Proposition 2 shows that the mapping $(\omega, g) \to a$ is linear within a constant-radius neighborhood of ω^* . Then, since the squared loss of our model is convex with respect to a, it is also convex with respect to (ω, g) jointly within the local region of ω^* . Note that the radius c_ω is a constant independent of N and T under the mild condition that $\underline{\alpha}_{\mathrm{MA}} \asymp \overline{\alpha}_{\mathrm{MA}}$, in which case $\{\|G_{p+j}^*\|_{\mathrm{F}}\}_{j=1}^{r+2s}$ are of the same order of magnitude.

Since Proposition 2 relies on confining ω to a local neighborhood of ω^* , the theoretical guarantees derived in this article are applicable to local estimators. That is, to derive nonasymptotic error bounds, we need to assume that the estimator $\widehat{\omega}$ obtained from (3.2) lies within the local region of ω^* defined in Proposition 2. We will discuss the practical aspect of this assumption after stating the main result. For simplicity, denote

$$\eta_T = \sqrt{\frac{\kappa_2 \lambda_{\max}(\mathbf{\Sigma}_{\varepsilon}) \log\{N(p \vee 1)\}}{\kappa_1^2 T}}$$
 and $\varpi = \frac{\lambda_{\max}(\mathbf{\Sigma}_{\varepsilon})}{\kappa_2(p \vee 1)}$.

Theorem 2. Suppose that Assumptions 1–4 hold with $\sum_{j=0}^{\infty}\|\boldsymbol{\Psi}_{j}^{*}\|_{\mathrm{op}}^{2}<\infty,\ R_{q}\lesssim\varpi/\eta_{T}^{2-q},\ \alpha^{2}\lesssim R_{q}/R_{q}^{\mathrm{MA}},$ $\varpi\lesssim\overline{\alpha}_{\mathrm{MA}}^{2}R_{q}/R_{q}^{\mathrm{MA}},$ and $\underline{\alpha}_{\mathrm{MA}}>0$. In addition, assume that $\log N\gtrsim(\kappa_{2}/\kappa_{1})^{2},\ T\gtrsim\max\{\kappa_{2}(p\vee1)^{4},(\kappa_{2}/\kappa_{1})^{2}(p\vee1)\log\{(\kappa_{2}/\kappa_{1})\alpha N(p\vee1)\}\},$ and we solve (3.2) with $\lambda_{g}\asymp\sqrt{\kappa_{2}\lambda_{\mathrm{max}}(\boldsymbol{\Sigma}_{\varepsilon})\log\{N(p\vee1)\}/T}.$ If $\|\widehat{\boldsymbol{\omega}}-\boldsymbol{\omega}^{*}\|_{2}\leq c_{\boldsymbol{\omega}},$ then with probability at least $1-C(p\vee1)e^{-c(\kappa_{1}/\kappa_{2})^{2}\log N},$

$$\begin{split} \|\widehat{\pmb{a}} - \pmb{a}^*\|_2 &\lesssim \eta_T^{1-q/2} \sqrt{R_q} \quad \text{and} \\ \frac{1}{T} \sum_{t=1}^T \left\| \sum_{h=1}^{t-1} (\widehat{\pmb{A}}_h - \pmb{A}_h^*) \pmb{y}_{t-h} \right\|_2^2 &\lesssim \frac{\eta_T^{2-q} R_q}{\kappa_1^{1-q}}. \end{split}$$

Combining Theorem 2 with Proposition 2, we immediately have the estimation error bounds $\|\widehat{\boldsymbol{g}} - \boldsymbol{g}^*\|_2 \lesssim \eta_T^{1-q/2} \sqrt{R_q}$ and $\|\widehat{\boldsymbol{\omega}} - \boldsymbol{\omega}^*\|_2 \lesssim \underline{\alpha}_{\mathrm{MA}}^{-1} \eta_T^{1-q/2} \sqrt{R_q}$. In particular, under exact sparsity, when r = s = 0, the bound for $\|\widehat{\boldsymbol{a}} - \boldsymbol{a}^*\|_2$ in Theorem 2 matches that for the Lasso estimator of VAR(p) models in Basu and Michailidis (2015), while the Gaussian assumption is relaxed. Also note that we do not require the uniqueness of the optimal solution to (3.2), that is, Theorem 2 is valid for all local optima within the constant-radius neighborhood of $\boldsymbol{\omega}^*$.

The JE can be efficiently implemented via the block coordinate descent algorithm; see Section S1.1 of the supplementary file for details. While the value of c_{ω} is unknown in practice, it is known to be independent of N and T under the mild condition that $\underline{\alpha}_{\mathrm{MA}} \asymp \overline{\alpha}_{\mathrm{MA}}$. The practical implication of the condition $\|\widehat{\boldsymbol{\omega}} - \boldsymbol{\omega}^*\|_2 \le c_{\omega}$ is that a reasonably good initialization for ω will be needed for the optimization algorithm of (3.2). For nonconvex estimators, to meet such requirements, commonly a convex preliminary estimator is used to initialize the algorithm (e.g., Janková and van de Geer 2021). However, for our model, the initialization task can be simplified, because the r values

for details.

 $\lambda_1, \ldots, \lambda_r \in (-1, 1)$ and the *s* values $\eta_1, \ldots, \eta_s \in [0, 1) \times (0, \pi)$ are restricted to bounded spaces and must be well separated from one another; see Assumptions 1(i) and 2(i). In fact, when *r* and *s* are larger, the initialization of ω will be even easier, as the selected *r* and *s* values will be denser on the bounded space and hence naturally tend to be closer to the true values. In practice, we recommend considering several different initial values for ω and selecting the solution of the optimization with minimum insample squared loss; see Section S1.2 of the supplementary file

Remark 5. Following the method for sparse VAR(P) models in Krampe and Paparoditis (2021), under a weak sparsity assumption on Σ_{ε} , we can construct a high-dimensional estimator of Σ_{ε} as $\widehat{\Sigma}_{\varepsilon} = \mathrm{THR}_{\lambda_{\varepsilon}}(T^{-1}\sum_{t=1}^{T}\widehat{\boldsymbol{\varepsilon}_{t}}\widehat{\boldsymbol{\varepsilon}}_{t}^{\top})$, where the residuals $\widehat{\boldsymbol{\varepsilon}}_{t}$ are obtained based on $\widehat{\boldsymbol{A}}_{h}$'s, and $\mathrm{THR}_{\lambda_{\varepsilon}}(\cdot)$ is the entrywise thresholding function with a chosen threshold parameter $\lambda_{\varepsilon} > 0$; see Krampe and Paparoditis (2021) for details. Then, based on $\widehat{\Sigma}_{\varepsilon}$ and $\widehat{\boldsymbol{A}}_{h}$'s, we can estimate var(y_{t}), so the instantaneous cross-sectional dependence can be interpreted. We leave a rigorous theoretical study of this estimation for future research.

Remark 6. While Theorem 2 establishes statistical error bounds, an interesting avenue for future research is to develop a more comprehensive estimation theory that integrates both statistical and algorithmic convergence analyses; see similar works such as Agarwal, Negahban, and Wainwright (2012) and Loh (2017). To tackle the theoretical challenges arising from the nonconvexity of the loss function, Proposition 2 may be leveraged to transform the problem into a convex one within a local region around ω^* .

3.2. ℓ_1 -Regularized Rowwise Estimator

While Theorem 2 allows R_q to grow with N, it requires $R_q \lesssim \varpi/\eta_T^{2-q}$; for example, if q=0, then this essentially will become $R_0 \lesssim T/\log\{N(p\vee 1)\}$. However, this requirement could be stringent when T is relatively small. To relax the sparsity requirement, we further introduce a rowwise estimator (RE) based on separately fitting each row of the proposed model.

For $1 \leq i \leq N$, the *i*th row of model (2.4) is $y_{i,t} = \sum_{h=1}^{\infty} \boldsymbol{a}_{i,h}^{\top} y_{t-h} + \varepsilon_{i,t}$, where $\boldsymbol{a}_{i,h} = \sum_{k=1}^{d} \ell_{h,k}(\boldsymbol{\omega}) \boldsymbol{g}_{i,k} \in \mathbb{R}^N$ is the *i*th row of \boldsymbol{A}_h , and $\boldsymbol{g}_{i,k} \in \mathbb{R}^N$ is the *i*th row of \boldsymbol{G}_k . Then, the squared loss for the *i*th row is $\mathbb{L}_{i,T}(\boldsymbol{\omega},\boldsymbol{g}_i) = T^{-1} \sum_{t=1}^{T} (y_{i,t} - \sum_{h=1}^{\infty} \boldsymbol{a}_{i,h}^{\top} y_{t-h})^2 = T^{-1} \sum_{t=1}^{T} \{y_{i,t} - \sum_{k=1}^{d} \boldsymbol{g}_{i,k}^{\top} \sum_{h=1}^{\infty} \ell_{h,k}(\boldsymbol{\omega}) \boldsymbol{y}_{t-h} \}^2$, where $\boldsymbol{g}_i = (\boldsymbol{g}_{i,1}^{\top}, \dots, \boldsymbol{g}_{i,d}^{\top})^{\top} \in \mathbb{R}^{Nd}$ is the *i*th row of $\boldsymbol{G} = (\boldsymbol{G}_1, \dots, \boldsymbol{G}_d)$. Note that joint loss function as defined in the previous section can be decomposed as $\mathbb{L}_T(\boldsymbol{\omega},\boldsymbol{g}) = \sum_{i=1}^{N} \mathbb{L}_{i,T}(\boldsymbol{\omega},\boldsymbol{g}_i)$. Thus, the rowwise losses $\mathbb{L}_{i,T}(\cdot)$'s can be minimized separately with respect to \boldsymbol{g}_i for $1 \leq i \leq N$. Meanwhile, since $\boldsymbol{\omega}$ is shared by all $\mathbb{L}_{i,T}(\cdot)$'s, each rowwise minimization can yield a consistent estimator of $\boldsymbol{\omega}$. This motivates us to consider the following ℓ_1 -regularized RE for $1 \leq i \leq N$:

$$(\widehat{\boldsymbol{\omega}}_i, \widehat{\boldsymbol{g}}_i) = \arg\min_{\boldsymbol{\omega} \in \Omega, \, \boldsymbol{g}_i \in \mathbb{R}^{Nd}} \left\{ \widetilde{\mathbb{L}}_{i,T}(\boldsymbol{\omega}, \boldsymbol{g}_i) + \lambda_g \|\boldsymbol{g}_i\|_1 \right\}, \quad (3.3)$$

where $\lambda_g > 0$ is the regularization parameter, and $\widetilde{\mathbb{L}}_{i,T}(\boldsymbol{\omega}, \boldsymbol{g}_i)$ is defined by setting the initial values $\{y_{i,s}, s \leq 0\}$ to zero,

that is, $\widetilde{\mathbb{L}}_{i,T}(\boldsymbol{\omega}, \boldsymbol{g}_i) = T^{-1} \sum_{t=1}^T (y_{i,t} - \sum_{h=1}^{t-1} \boldsymbol{a}_{i,h}^\top \boldsymbol{y}_{t-h})^2 = T^{-1} \sum_{t=1}^T \{y_{i,t} - \sum_{k=1}^d \boldsymbol{g}_{i,k}^\top \sum_{h=1}^{t-1} \ell_{h,k}(\boldsymbol{\omega}) \boldsymbol{y}_{t-h} \}^2$. Let $\boldsymbol{a}_i = (\boldsymbol{a}_{i,1}^\top, \boldsymbol{a}_{i,2}^\top, \dots)^\top \in \mathbb{R}^{\infty}$ be the ith row of $\boldsymbol{A} = (\boldsymbol{A}_1, \boldsymbol{A}_2, \dots)$ for $1 \leq i \leq N$. Note that $\boldsymbol{a}_i = (\boldsymbol{L}(\boldsymbol{\omega}) \otimes \boldsymbol{I}_N) \boldsymbol{g}_i$. Based on (3.3), we have $\widehat{\boldsymbol{a}}_i = (\widehat{\boldsymbol{a}}_{i,1}^\top, \widehat{\boldsymbol{a}}_{i,2}^\top, \dots)^\top = (\boldsymbol{L}(\widehat{\boldsymbol{\omega}}) \otimes \boldsymbol{I}_N) \widehat{\boldsymbol{g}}_i$, where $\widehat{\boldsymbol{g}}_i = (\widehat{\boldsymbol{g}}_{i,1}^\top, \dots, \widehat{\boldsymbol{g}}_{i,d}^\top)^\top$, and $\widehat{\boldsymbol{a}}_{i,h} = \sum_{k=1}^d \ell_{h,k}(\widehat{\boldsymbol{\omega}}_i) \widehat{\boldsymbol{g}}_{i,k}$. The algorithm for the RE is provided in Section S1.1 of the supplementary file.

Similar to the previous section, we can derive the nonasymptotic error bounds for the RE. For $1 \leq i \leq N$, let $\mathbf{g}_{i,\mathrm{AR}} = (\mathbf{g}_{i,1}^\top, \ldots, \mathbf{g}_{i,p}^\top)^\top \in \mathbb{R}^{Np}$ and $\mathbf{g}_{i,\mathrm{MA}} = (\mathbf{g}_{i,p+1}^\top, \ldots, \mathbf{g}_{i,d}^\top)^\top \in \mathbb{R}^{N(r+2s)}$. To define the weak sparsity of $\mathbf{g}_{i,\mathrm{AR}}^*$ and $\mathbf{g}_{i,\mathrm{MA}}^*$, we consider the ℓ_q -"balls", $\mathbb{B}_q(R_{i,q}^\mathrm{AR}) := \{\mathbf{g}_{i,\mathrm{AR}} \in \mathbb{R}^{Np} \mid \sum_{k=1}^p \sum_{j=1}^N |g_{i,j,k}|^q \leq R_{i,q}^\mathrm{AR} \}$ and $\mathbb{B}_q(R_{i,q}^\mathrm{MA}) := \{\mathbf{g}_{i,\mathrm{MA}} \in \mathbb{R}^{N(r+2s)} \mid \sum_{k=p+1}^d \sum_{j=1}^N |g_{i,j,k}|^q \leq R_{i,q}^\mathrm{MA} \}$. The following is the row-wise counterpart of Assumption 4.

Assumption 4' (Rowwise weak sparsity). For $1 \leq i \leq N$, there exists $q \in [0,1]$ such that $\mathbf{g}_{i,AR}^* \in \mathbb{B}_q(R_{i,q}^{AR})$ and $\mathbf{g}_{i,MA}^* \in \mathbb{B}_q(R_{i,q}^{MA})$ for some radii $R_{i,q}^{AR}$, $R_{i,q}^{MA} > 0$.

Let $R_{i,q} = R_{i,q}^{\mathrm{AR}} + R_{i,q}^{\mathrm{MA}}$, and then by Assumption 4', $\mathbf{g}_i^* \in \mathbb{B}_q(R_{i,q}) := \{\mathbf{g}_i \in \mathbb{R}^{Nd} \mid \sum_{k=1}^d \sum_{j=1}^N |g_{i,j,k}|^q \leq R_{i,q} \}$. Moreover, Assumption 4' implies the overall sparsity level in Assumption 4, since it leads to $\mathbf{g}_{\mathrm{AR}}^* \in \mathbb{B}_q(R_q^{\mathrm{AR}})$, $\mathbf{g}_{\mathrm{MA}}^* \in \mathbb{B}_q(R_q^{\mathrm{MA}})$, and consequently $\mathbf{g}^* \in \mathbb{B}_q(R_q)$, where $R_q^{\mathrm{AR}} = \sum_{i=1}^N R_{i,q}^{\mathrm{AR}}$, $R_q^{\mathrm{MA}} = \sum_{i=1}^N R_{i,q}^{\mathrm{MA}}$, and $R_q = R_q^{\mathrm{MA}} + R_q^{\mathrm{AR}} = \sum_{i=1}^N R_{i,q}$. For $1 \leq i \leq N$, let $\underline{\alpha}_{i,\mathrm{MA}} = \min_{1 \leq j \leq r+2s} \|\mathbf{g}_{i,p+j}^*\|_2$ and

For $1 \leq i \leq N$, let $\underline{\alpha}_{i,\text{MA}} = \min_{1 \leq j \leq r+2s} \| \mathbf{g}^*_{i,p+j} \|_2$ and $\overline{\alpha}_{i,\text{MA}} = \max_{1 \leq j \leq r+2s} \| \mathbf{g}^*_{i,p+j} \|_2$, which are both allowed to grow with N. Denote $\alpha_i = \overline{\alpha}_{i,\text{MA}} / \underline{\alpha}_{i,\text{MA}}$. The rowwise counterparts of Proposition 2 and Theorem 2 are established as follows.

Proposition 3. Fix $1 \le i \le N$. Suppose that $\underline{\alpha}_{i,\text{MA}} > 0$. Then under Assumptions 1(i) and 2, there exists a constant $c_{i,\omega} = \min(2, c/\alpha_i) > 0$ such that for any $\omega \in \Omega$ with $\|\omega - \omega^*\|_2 \le c_{i,\omega}$, it holds $\|\mathbf{g}_i - \mathbf{g}_i^*\|_2 + \underline{\alpha}_{i,\text{MA}} \|\omega - \omega^*\|_2 \lesssim \|\mathbf{a}_i - \mathbf{a}_i^*\|_2^2 \lesssim \|\mathbf{g}_i - \mathbf{g}_i^*\|_2 + \overline{\alpha}_{i,\text{MA}} \|\omega - \omega^*\|_2$, where $\mathbf{a}_i = (L(\omega) \otimes I_N)\mathbf{g}_i$.

Theorem 3. Suppose that Assumptions 1–3 and 4' hold with $\sum_{j=0}^{\infty}\|\boldsymbol{\Psi}_{j}^{*}\|_{\mathrm{op}}^{2}<\infty,\ R_{i,q}\lesssim\varpi/\eta_{T}^{2-q},\ \alpha_{i}^{2}\lesssim R_{i,q}/R_{i,q}^{\mathrm{MA}},$ $\varpi\lesssim\overline{\alpha}_{i,\mathrm{MA}}^{2}R_{i,q}/R_{i,q}^{\mathrm{MA}},\ \mathrm{and}\ \underline{\alpha}_{i,\mathrm{MA}}>0,\ \mathrm{for}\ 1\leq i\leq N.$ In addition, assume that $\log N\gtrsim(\kappa_{2}/\kappa_{1})^{2},\ T\gtrsim\max\{\kappa_{2}(p\vee1)^{4},(\kappa_{2}/\kappa_{1})^{2}(p\vee1)\log\{(\kappa_{2}/\kappa_{1})\alpha_{\mathrm{max}}N(p\vee1)\}\},$ with $\alpha_{\mathrm{max}}=\max_{1\leq i\leq N}\alpha_{i},\ \mathrm{and}\ \mathrm{we}\ \mathrm{solve}\ (3.3)\ \mathrm{with}\ \lambda_{g}\asymp\sqrt{\kappa_{2}\lambda_{\mathrm{max}}(\boldsymbol{\Sigma}_{\varepsilon})\log\{N(p\vee1)\}/T}.$ For $1\leq i\leq N,$ if $\|\widehat{\boldsymbol{\omega}}_{i}-\boldsymbol{\omega}^{*}\|_{2}\leq c_{i,\boldsymbol{\omega}},$ then with probability at least $1-C(p\vee1)e^{-c(\kappa_{1}/\kappa_{2})^{2}\log N},$

$$\begin{split} \|\widehat{\pmb{a}}_i - \pmb{a}_i^*\|_2 &\lesssim \eta_T^{1-q/2} \sqrt{R_{i,q}} \quad \text{and} \\ \frac{1}{T} \sum_{t=1}^T \left\| \sum_{h=1}^{t-1} (\widehat{\pmb{a}}_{i,h} - \pmb{a}_{i,h}^*)^\top \pmb{y}_{t-h} \right\|_2^2 &\lesssim \frac{\eta_T^{2-q} R_{i,q}}{\kappa_1^{1-q}}. \end{split}$$

Compared to Theorem 3, the sparsity condition in Theorem 3 is much weaker, that is, $R_{i,q} \lesssim \varpi/\eta_T^{2-q}$ for $1 \leq i \leq N$; or

essentially, $R_{i,0} \lesssim T/\log\{N(p \vee 1)\}$ when q = 0. Thus, the RE may be preferred in practice when T is relatively small.

Moreover, by Theorem 3 and Proposition 3, we have $\|\widehat{\mathbf{g}}_i - \mathbf{g}_i^*\|_2 \lesssim \eta_T^{1-q/2} \sqrt{R_{i,q}}$ and $\|\widehat{\boldsymbol{\omega}}_i - \boldsymbol{\omega}^*\|_2 \lesssim \underline{\alpha}_{i,\mathrm{MA}}^{-1} \eta_T^{1-q/2} \sqrt{R_{i,q}}$ for $1 \leq i \leq N$. Note that each RE $\widehat{\boldsymbol{\omega}}_i$ is a consistent estimator of $\boldsymbol{\omega}^*$, and the estimation error is proportional to $\underline{\alpha}_{i,\mathrm{MA}}^{-1} \sqrt{R_{i,q}}$. On the other hand, as implied by Theorem 2, the estimation error of the JE for $\boldsymbol{\omega}^*$ is proportional to $\underline{\alpha}_{\mathrm{MA}}^{-1} \sqrt{R_q}$. For example, if $R_{i,q} \asymp R_q/N$ and $\underline{\alpha}_{i,\mathrm{MA}}^2 \asymp \underline{\alpha}_{\mathrm{MA}}^2/N$, then the two bounds will be comparable. However, intuitively, allowing different estimators $\widehat{\boldsymbol{\omega}}_i$ for different rows may enhance the flexibility in practice, although it may also increase the risk of overfitting. In addition, combining the results for $\widehat{\boldsymbol{a}}_i$, $\widehat{\boldsymbol{g}}_i$ and the prediction error across all rows, we have $\|\widehat{\boldsymbol{a}} - \boldsymbol{a}^*\|_2 \lesssim \eta_T^{1-q/2} \sqrt{R_q}$, $\|\widehat{\boldsymbol{g}} - \boldsymbol{g}^*\|_2 \lesssim \eta_T^{1-q/2} \sqrt{R_q}$, and $T^{-1} \sum_{t=1}^T \|\sum_{h=1}^{t-1} (\widehat{\boldsymbol{A}}_h - \boldsymbol{A}_h^*) \boldsymbol{y}_{t-h} \|_2^2 \lesssim \eta_T^{2-q} R_q/\kappa_1^{1-q}$. Here, with a slight abuse of notation, $\widehat{\boldsymbol{a}}_i$, $\widehat{\boldsymbol{g}}_i$, and $\widehat{\boldsymbol{A}}_h$'s represent the estimates obtained based on merging the RE $\widehat{\boldsymbol{a}}_i$ or $\widehat{\boldsymbol{g}}_i$ for $1 \leq i \leq N$. Note that these bounds match exactly those of the JE in the previous section.

In addition to the above upper bounds analysis, we numerically assess the actual comparative performance of RE and JE via simulations in Section S2.2 of the supplementary file. It is shown that they can perform very similarly for the estimation of g^* , while RE may outperform JE for the estimation of ω^* , resulting in an overall advantage for the estimation of a^* . However, as long as T is not too small compared to R_q , JE and RE tend to have similar out-of-sample forecast accuracy; see the empirical analysis in Section 6 and the simulation study in Section S2.4 of the supplementary file for details. Furthermore, as commented by one referee, the competitive numerical performance of the JE might hint that its more stringent sparsity condition could be an artifact of the proof technique.

4. Model Order Selection

In this section, we introduce a Bayesian information criterion (BIC) based approach to selecting the model orders for the proposed high-dimensional SPVAR(∞) model.

Let $\mathcal{M}^*=(p^*,r^*,s^*)$ denote the true orders. For the feasibility of order selection, it is crucial to ensure that \mathcal{M}^* is irreducible; that is, if $\{y_t\}$ is generated with orders \mathcal{M}^* , there is no alternative parameterization with reduced orders. As established in Lemma S14 in the supplementary file, the irreducibility of r^* and s^* is guaranteed if λ_j^* 's, γ_m^* 's, and $\underline{\alpha}_{\mathrm{MA}}$ are nonzero. On the other hand, p^* is irreducible under the following assumption.

Assumption 5 (Irreducibility).
$$G_{p^*} \neq \sum_{j=1}^{r^*} G_{p^*+j} + \sum_{m=1}^{s^*} G_{p^*+j}$$

To select the model orders, for any $\mathcal{M} = (p, r, s)$, we define the high-dimensional BIC,

$$BIC(\mathcal{M}) = \log \widetilde{\mathbb{L}}_{T}(\widehat{\boldsymbol{\omega}}_{\mathcal{M}}, \widehat{\boldsymbol{g}}_{\mathcal{M}})$$

$$+ \tau_{N} d \left[\frac{\log\{N(p \vee 1)\}}{T} \right]^{1-q/2} \log T,$$
(4.1)

where $\widehat{\boldsymbol{\omega}}_{\mathcal{M}}$ and $\widehat{\boldsymbol{g}}_{\mathcal{M}}$ denote estimates obtained by fitting the model with orders \mathcal{M} using either the JE in (3.2) or the RE in

(3.3). In particular, if the RE is employed, then $\widetilde{\mathbb{L}}_T(\widehat{\boldsymbol{\omega}}_{\mathcal{M}}, \widehat{\boldsymbol{g}}_{\mathcal{M}}) = \sum_{i=1}^N \mathbb{L}_{i,T}(\widehat{\boldsymbol{\omega}}_{i,\mathcal{M}}, \widehat{\boldsymbol{g}}_{i,\mathcal{M}})$, where $\widehat{\boldsymbol{\omega}}_{\mathcal{M}}$ and $\widehat{\boldsymbol{g}}_{\mathcal{M}}$ denote collections of $\widehat{\boldsymbol{\omega}}_{i,\mathcal{M}}$'s and $\widehat{\boldsymbol{g}}_{i,\mathcal{M}}$'s, respectively. Note that for notational simplicity, we suppress the dependence of $\widetilde{\mathbb{L}}_T(\cdot)$ and $\mathbb{L}_T(\cdot)$ on \mathcal{M} in this section. Additionally, $\tau_N > 0$ is a sequence possibly dependent on N satisfying the following condition.

Assumption 6 (Penalty parameter). $\tau_N \gtrsim N^{-1} R_q \{\kappa_2 \lambda_{\max}(\Sigma_{\varepsilon})\}^{1-q/2} / \kappa_1^{3-2q}.$

Assumption 6 ensures that the proposed BIC can rule out any overspecified model, $\mathcal{M} \in \mathcal{M}_{\text{over}} = \{\mathcal{M} \in \mathcal{M} \mid p \geq p^*, r \geq r^* \text{ and } s \geq s^*\} \setminus \mathcal{M}^*$. When the constants κ_1, κ_2 and $\lambda_{\max}(\Sigma_\varepsilon)$ are fixed, Assumption 6 can be simplified to $\tau_N \gtrsim N^{-1}R_q$. While R_q is unknown in practice, to set a reasonable τ_N , we may assume that $R_q \lesssim N$; for example, this will hold if G_k^* 's are (weakly) rowsparse. Then it would suffice to fix $\tau_N \equiv \tau > 0$. In practice, we may simply set q = 0. We recommend $\tau = 0.05$, which performs well in our simulations.

Based on (4.1), we estimate the model orders by

$$\widehat{\mathcal{M}} = (\widehat{p}, \widehat{r}, \widehat{s}) = \arg\min_{\mathcal{M} \in \mathscr{M}} \mathrm{BIC}(\mathcal{M}),$$

where $\mathcal{M} = \{(p,r,s) \mid 0 \le p \le \overline{p}, 0 \le r \le \overline{r}, 0 \le s \le \overline{s}\}$, with $\overline{\mathcal{M}} := (\overline{p},\overline{r},\overline{s})$ being predetermined maximum orders. Since the true orders are usually small in practice, $\overline{\mathcal{M}}$ need not be large; for example $\overline{p} = \overline{r} = \overline{s} = 6$ may be sufficient for most applications. Our simulations show that $\widehat{\mathcal{M}}$ is insensitive to the choice of $\overline{\mathcal{M}}$ as long as it is large enough compared to \mathcal{M}^* .

Let $\mathcal{M}_{\text{mis}} = \{\mathcal{M} \in \mathcal{M} \mid p < p^*, r < r^* \text{ or } s < s^*\}$. To establish the conditions that prevent the proposed BIC from selecting any misspecified model, we need to accurately quantify the minimum difference between any $\mathcal{M} \in \mathcal{M}_{\text{mis}}$ and \mathcal{M}^* . This analysis is challenging since there is no monotonic nested ordering over \mathcal{M} due to the involvement of three different orders, p, r and s. Particularly, $\mathcal{M} \in \mathcal{M}_{\text{mis}}$ may not be nested within \mathcal{M}^* regarding all three orders. For instance, if $\mathcal{M}^* = (1, 1, 0)$, then a misspecified model may be $\mathcal{M}_1 = (\bar{p}, 0, 0)$ or $\mathcal{M}_2 = (0, \bar{r}, \bar{s})$, where, for example, $\bar{p} = \bar{r} = \bar{s} = 6$. Clearly, we cannot simply treat \mathcal{M}_1 or \mathcal{M}_2 as a smaller model than \mathcal{M}^* , as they possess orders as large as \bar{p}, \bar{r} , or \bar{s} .

To uniformly accommodate the possibly nonnested relationship between $\mathcal{M} \in \mathcal{M}_{\text{mis}}$ and \mathcal{M}^* , we leverage their connections with a common model, $\overline{\mathcal{M}} = (\overline{p}, \overline{r}, \overline{s})$. Specifically, we can show that model (2.4) with any orders $\mathcal{M} = (p, r, s) \in \mathcal{M}$ can be reparameterized as the model with $\overline{\mathcal{M}} = (\overline{p}, \overline{r}, \overline{s})$. In addition, the corresponding parameter vectors, denoted $\overline{\boldsymbol{\omega}} \in (-1, 1)^{\overline{r}} \times \overline{\boldsymbol{\Pi}}^{\overline{s}}$ and $\overline{\boldsymbol{g}} \in \mathbb{R}^{N \times N\overline{d}}$, satisfy the following equality constraints:

$$\overline{C}_1^{\mathcal{M}} \overline{\omega} = \mathbf{0}$$
 and $(\overline{C}_2^{\mathcal{M}} (\overline{\omega}) \otimes I_{N^2}) \overline{g} = \mathbf{0}$, (4.2)

where $\overline{C}_1^{\mathcal{M}} \in \mathbb{R}^{(\delta_r+2\delta_s)\times(\overline{r}+2\overline{s})}$ is a constant matrix encoding $(\delta_r+2\delta_s)$ constraints on $\overline{\boldsymbol{\omega}}$, specifying which elements are restricted to zero, and the matrix function $\overline{\boldsymbol{C}}_2^{\mathcal{M}}(\overline{\boldsymbol{\omega}}) \in \mathbb{R}^{\delta_d \times \overline{d}}$ encodes δ_d equality constraints on $\overline{\boldsymbol{g}}$ for any given $\overline{\boldsymbol{\omega}}$, with $\delta_r=\overline{r}-r$, $\delta_s=\overline{s}-s$, and $\delta_d=\overline{d}-d$; see Section S7.3 in the supplementary file for detailed definitions of $\overline{\boldsymbol{C}}_1^{\mathcal{M}}$ and $\overline{\boldsymbol{C}}_2^{\mathcal{M}}(\cdot)$. In particular, increasing p by one amounts to deleting a particular

row from the constraint matrix $\overline{C}_2^{\mathcal{M}}(\cdot)$. On the other hand, increasing r (or s) by one is equivalent to deleting a particular row (or a pair of rows) from both $\overline{C}_1^{\mathcal{M}}$ and $\overline{C}_2^{\mathcal{M}}(\cdot)$.

Note that $\overline{C}_2^{\mathcal{M}}(\cdot)$ cannot reduce to a constant matrix independent of $\overline{\omega}$ except in the special cases where $p = \overline{p} - 1$ or r = s =0. In particular, when $p = \overline{p} - 1$, the second equation in (4.2) is essentially the reducibility condition of \overline{p} , which resembles that for p^* in Assumption 5(i). However, in general, this equation represents much more intricate constraints, since $\overline{C}_2^{\mathcal{M}}(\cdot)$ is a nonlinear function. The complexity of this form can be understood from two perspectives. First, due to the nonlinearity of model (2.4) in ω , the effect of any underspecification in r or s will be highly nonlinear. Second, the order p plays a special role in the definition of $\ell_{h,k}(\cdot)$'s as it is involved in $\mathbb{I}_{\{h \geq p+1\}} \lambda_i^{h-p}$ and $\mathbb{I}_{\{h>p+1\}}\gamma_m^{h-p}$; see (2.3). Then, whenever $p\neq p^*$, the exponent h - p will differ from that under \mathcal{M}^* for all lags $h \geq p + 1$, thereby affecting all $\ell_{h,k}(\cdot)$'s. Consequently, due to the interplay between p and $\ell_{h,k}(\cdot)$'s, an underspecification in p generally will also have a nonlinear effect.

Let $\Gamma_{\mathcal{M}} = \{\overline{\boldsymbol{\omega}} \in (-1,1)^{\overline{r}} \times \boldsymbol{\Pi}^{\overline{s}}, \ \overline{\boldsymbol{g}} \in \mathbb{R}^{N^2 \overline{d}} : \overline{\boldsymbol{C}}_1^{\mathcal{M}} \overline{\boldsymbol{\omega}} = \mathbf{0} \text{ and } (\overline{\boldsymbol{C}}_2^{\mathcal{M}}(\overline{\boldsymbol{\omega}}) \otimes \boldsymbol{I}_{N^2})\overline{\boldsymbol{g}} = \mathbf{0} \}$ denote the restricted parameter space for any candidate model \mathcal{M} . By leveraging (4.2), we can characterize the minimum difference between the true model and the approximated model of orders $\mathcal{M} \in \mathcal{M}_{\text{mis}}$ via the quantity $\delta_{\mathcal{M}} := \kappa_1 \inf_{(\boldsymbol{\omega}, \boldsymbol{g}) \in \Gamma_{\mathcal{M}}} \|(\boldsymbol{L}(\boldsymbol{\omega}) \otimes \boldsymbol{I}_{N^2})\boldsymbol{g} - \boldsymbol{a}^*\|_2^2$; see Proposition S1 and the proof of Theorem 4 in Section S7 of the supplementary file for details. We may regard $\delta_{\mathcal{M}}$ as the signal strength of the misspecification. The following assumption guarantees that $\delta_{\mathcal{M}}$ is large enough for the BIC to detect the misspecification.

Assumption 7 (Minimum signal strength).

(i) $\min_{\mathcal{M} \in \mathcal{M}_{mis}} \delta_{\mathcal{M}}/N \gg (T^{-1} \log N)^{1-q/2} \tau_N \log T$; and (ii) $\max_{\mathcal{M} \in \mathcal{M}_{mis}} \delta_{\mathcal{M}}^{-1} |\widetilde{\mathbb{L}}_T(\widehat{\boldsymbol{\omega}}_{\mathcal{M}}, \widehat{\boldsymbol{g}}_{\mathcal{M}}) - \mathbb{E}\{\mathbb{L}_T(\boldsymbol{\omega}_{\mathcal{M}}^\circ, \boldsymbol{g}_{\mathcal{M}}^\circ)\}| = o_p(1),$ where $(\boldsymbol{\omega}_{\mathcal{M}}^\circ, \boldsymbol{g}_{\mathcal{M}}^\circ)$ is the minima of $\mathbb{E}\{\mathbb{L}_T(\boldsymbol{\omega}_{\mathcal{M}}, \boldsymbol{g}_{\mathcal{M}})\}$ over the parameter space $\boldsymbol{\omega}_{\mathcal{M}} \in (-1, 1)^r \times \boldsymbol{\Pi}^s$ and $\boldsymbol{g}_{\mathcal{M}} \in \mathbb{R}^{N^2d}$.

Note that $\delta_{\mathcal{M}}/N$ can be viewed as the average level of misspecification across N rows of the model equation. As mentioned earlier, we may let $\tau_N \equiv \tau$ under mild condition. Thus, the lower bound in Assumption 7(i) tends to zero as $T \to \infty$. Assumption 7(ii) requires that the empirical loss for any fitted misspecified model converges to some population loss at a rate faster than $\delta_{\mathcal{M}}$ as $T \to \infty$. Here the mispecified model with parameters $(\boldsymbol{\omega}_{\mathcal{M}}^{\circ}, \boldsymbol{g}_{\mathcal{M}}^{\circ})$ can be understood as the best approximation of the process $\{\boldsymbol{y}_t\}$ under the misspecification. Now we are ready to establish the consistency of the estimator $\widehat{\mathcal{M}}$.

Theorem 4. If the JE (or the RE) is used, suppose that for any $\mathcal{M} \in \mathcal{M}_{\text{over}}$, there is a subvector $\widehat{\boldsymbol{\omega}}_{\mathcal{M}^*} \in (-1,1)^{r^*} \times \boldsymbol{\Pi}^{s^*}$ of $\widehat{\boldsymbol{\omega}}_{\mathcal{M}}$ (or $\widehat{\boldsymbol{\omega}}_{i,\mathcal{M}^*} \in (-1,1)^{r^*} \times \boldsymbol{\Pi}^{s^*}$ of $\widehat{\boldsymbol{\omega}}_{i,\mathcal{M}}$ with $1 \leq i \leq N$) such that $\|\widehat{\boldsymbol{\omega}}_{\mathcal{M}^*} - \boldsymbol{\omega}^*\|_2 \leq c_{\boldsymbol{\omega}}$ (or $\|\widehat{\boldsymbol{\omega}}_{i,\mathcal{M}^*} - \boldsymbol{\omega}^*\|_2 \leq c_{i,\boldsymbol{\omega}}$ with $1 \leq i \leq N$), and the conditions in Theorem 2 (or 3) hold with $\mathcal{M} = \mathcal{M}^*$. In addition, suppose that $\overline{\mathcal{M}}$ is fixed, with $\overline{p} \geq p^*, \overline{r} \geq r^*$ and $\overline{s} \geq s^*$. Under Assumptions 5–7, $\mathbb{P}(\widehat{\mathcal{M}} = \mathcal{M}^*) \to 1$ as $N, T \to \infty$.

5. Simulation Experiments

In this section, we present two simulation experiments to verify the estimation error rates of the JE and the consistency of the BIC. Four additional experiments on the estimation error of the RE, its comparison with the JE, sensitivity analysis of the initialization for $\{y_t, t \leq 0\}$, and comparison of the proposed estimators with competing approaches are provided in Section S2 of the supplementary file.

Throughout this section, we generate $\{y_t\}$ from model (2.4), where $\{\varepsilon_t\}$ are generated independently from $N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$ with $\sigma=0.2$, and each \mathbf{G}_k is exactly sparse with cN nonzero entries for $1 \le k \le d$, so the overall sparsity level is $R_0 = cdN$. We generate $\{G_k\}_{k=1}^d$ by drawing their nonzero entries independently from the uniform distribution on [-0.5, 0.5]. Then, to ensure the stationarity of $\{y_t\}$, after setting $\boldsymbol{\omega}$, we rescale all G_k 's by a common factor such that $\rho(\underline{G}_1) + \bar{\rho} \sum_{k=1}^{r+2s} \rho(G_{p+k})/(1-\bar{\rho}) = 0.8$; see Theorem 1.

In the first experiment, we examine the estimation error rates for the JE. Two data generating processes are considered: (p, r, s) = (1, 1, 0) (DGP1) and (1, 0, 1) (DGP2), where $\lambda_1 =$ -0.6 for DGP1, and $(\gamma_1, \theta_1) = (0.6, \pi/4)$ for DGP2. We let all G_k 's be row-sparse matrices with three nonzero entries in each row, that is, $R_0 = 3dN$, where N = 10, 20, 40, or 80. Note that by Theorem 2, we have $\|\hat{a} - a^*\|_2 / \sqrt{N} \lesssim \eta_T \sqrt{R_0/N}$, $\|\widehat{\boldsymbol{g}} - {\boldsymbol{g}}^*\|_2 / \sqrt{N} \lesssim \eta_T \sqrt{R_0/N}$, and $\underline{\alpha}_{\mathrm{MA}} \|\widehat{\boldsymbol{\omega}} - {\boldsymbol{\omega}}^*\|_2 / \sqrt{N} \lesssim$ $\eta_T \sqrt{R_0/N}$, where $\eta_T = \sqrt{T^{-1} \log N}$. To verify these bounds, we choose a grid of equally spaced values for the theoretical rate $\eta_T \sqrt{R_0/N} = \sqrt{3T^{-1}d\log N}$ within the range of $\mathcal{I}_1 =$ [0.3756, 0.4981] for DGP1 and $\mathcal{I}_2 = [0.46, 0.61]$ for DGP2. Then we compute T given the theoretical rate, N and d. The selected ranges \mathcal{I}_1 and \mathcal{I}_2 lead to the same range of T for both DGPs under any N; that is, the ranges of the x-axis in Figure 3 are set such that the corresponding points in upper and lower panels share the same T. Across all settings, T falls in the range of [55, 186]. Figure 3 plots the scaled estimation errors $\|\hat{a}\|$ $a^*\parallel_2/\sqrt{N}$, $\|\widehat{g}-g^*\|_2/\sqrt{N}$, and $\underline{\alpha}_{\rm MA}\|\widehat{\omega}-\omega^*\|_2/\sqrt{N}$, averaged over 500 replications, against the theoretical rate $\eta_T \sqrt{R_0/N}$. An approximately linear relationship can be observed across all settings, confirming our theoretical results.

In the second experiment, we verify the consistency of the proposed BIC. Three cases of true model orders are considered: $(p^*, r^*, s^*) = (0, 0, 1), (0, 1, 1),$ and (1, 0, 1), referred to as DGPs 1, 2, and 3, respectively. We set N=40, $\theta_1=\pi/4$, and $\lambda_1=$ $-\gamma_1 = \bar{\rho}$, where three choices of the decay rate are considered: $\bar{\rho} \in \{0.45, 0.5, 0.5\}$. For $1 \le k \le d$, each G_k contains 3N nonzero entries, so $R_0 = 3dN$, but unlike the first experiment, we do not restrict each row of G_k to have exactly three nonzero entries. We set $\tau = 0.05$ and $\bar{p} = \bar{r} = \bar{s} = 9$; the results are found to be unchanged if the maximum orders are 3. Figure 4 displays the proportion of correct order selection based on 500 replications for each setting, with the models fitted by the JE; the results for the RE are very similar and hence omitted. It shows that the BIC generally performs better as T or $\bar{\rho}$ increases, and the proportion of correct order selection eventually becomes close to one with sufficiently large T. Thus, the consistency of the BIC is verified. Additionally, the required sample size for achieving accurate order selection follows this order among the three DGPs: DGP1 < DGP3 < DGP2. To understand this, first note that R_0 =

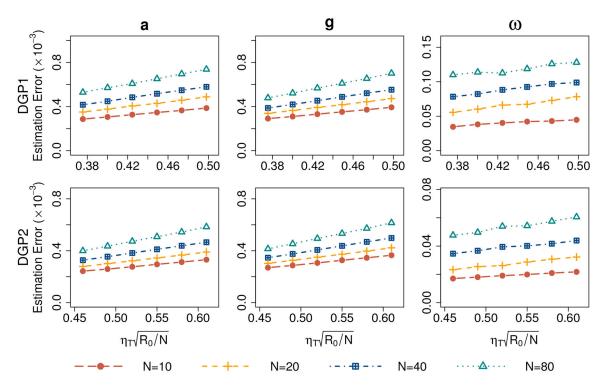


Figure 3. Plots of scaled estimation errors $\|\widehat{a} - a^*\|_2/\sqrt{N}$ (left panel), $\|\widehat{g} - g^*\|_2/\sqrt{N}$ (middle panel), and $\underline{\alpha}_{\text{MA}}\|\widehat{\omega} - \omega^*\|_2/\sqrt{N}$ (right panel) against theoretical rate $\eta_T\sqrt{R_0/N}$ for JE.

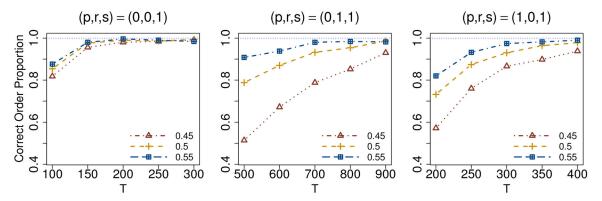


Figure 4. Proportion of correct model order selection for three DGPs and three choices of decay rates, $\bar{\rho} \in \{0.45, 0.5, 0.55\}$.

6N, 9N, and 9N for DGPs 1, 2, and 3, respectively. Thus, the estimation accuracy is highest for DGP1, and so is the order selection accuracy. Moreover, since DGP2 has a more complex temporal structure than DGP3, it leads to greater challenges in estimating ω and, consequently, in order selection.

6. Empirical Analysis

We analyze N=20 quarterly macroeconomic variables of the United States from the first quarter of 1969 to the fourth quarter of 2007. These are key economic and financial indicators collected by Koop (2013), seasonally adjusted as needed. We conduct the transformations following Koop (2013) to make all series stationary, resulting in a sample of length T=194. Then each series is normalized to have zero mean and unit variance; see Table S1 in the supplementary file for detailed descriptions of the 20 variables.

We first fit the proposed model to the entire dataset. Using the JE and the proposed BIC, we select (p,r,s)=(1,1,0), so d=2, and the fitted model is $y_t=\widehat{\mathbf{G}}_1y_{t-1}+\sum_{h=2}^{\infty}(-0.45)^{h-1}\widehat{\mathbf{G}}_2y_{t-h}+$ $\boldsymbol{\varepsilon}_t$, where $\widehat{\mathbf{G}}_1$ and $\widehat{\mathbf{G}}_2$ are displayed in Figure 5; the estimation results based on the RE are roughly similar and provided in the supplementary file. The stationarity of the model is confirmed by the method in Remark 4. As discussed in Section 2.2, $\widehat{\mathbf{G}}_1$ and $\widehat{\mathbf{G}}_2$ captures lag-one (or short-term) and higher-lag (or long-term) dependence, respectively. Note that $\widehat{\mathbf{G}}_1$ is much denser than $\widehat{\mathbf{G}}_2$, suggesting that many dynamic interactions are short-term. However, most of the nonzero entries in $\widehat{\mathbf{G}}_2$ are fairly large in absolute value, supporting the necessity of a VARMA-type model. For the Granger causal (GC) interpretation, take the model equation for real GDP (RGDP) as an example:

$$y_{\text{RGDP},t} = 0.17y_{\text{Cons},t-1} + 0.11y_{\text{IP:total},t-1} + 0.07y_{\text{HStarts:total},t-1} + 0.12y_{\text{S\&P:indust},t-1}$$

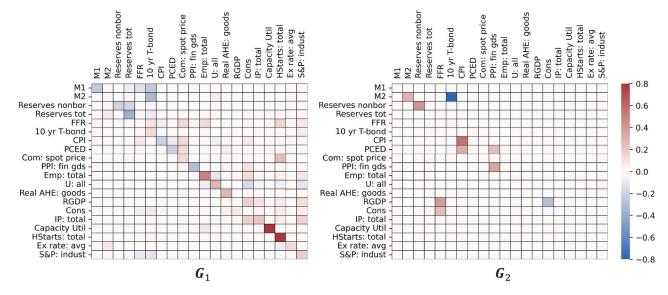


Figure 5. Estimates of G_1 and G_2 for the proposed model based on JE.

+
$$\sum_{h=2}^{\infty} (-0.45)^{h-1} (0.39y_{\text{FFR},t-h} - 0.30y_{\text{Cons},t-h})$$

+ $\varepsilon_{\text{RGDP},t}$,

suppressing other lag-one terms with coefficients less than 0.014 in absolute value for brevity. The above equation indicates that five time series are GC for RGDP and can be categorized as follows: (a) the industrial production index (IP: total), housing starts (HStarts: total), and S&P stock price index (S&P: indust) only have short-term influence on RGDP; (b) the federal funds rate (FFR) only has long-term influence on RGDP; (c) the real personal consumption expenditures (Cons) has both short-term and long-term influence on RGDP. For other insights from the estimation results, see Section S3 in the supplementary file for more discussions.

Next we evaluate the forecasting performance via a rolling procedure: First set the forecast origin to t=166 (Q4-2000). For each $k=1,\ldots,28$, fit the model using the data of $1\leq t\leq T_{\rm train}=165+k$, and then compute the one-step ahead forecast for t=166+k. Thus, rolling forecasts over the period of Q1-2001 to Q4-2007 are obtained. We measure the forecast error by $\|\widehat{\boldsymbol{y}}_t-\boldsymbol{y}_t\|_2$; our findings based on the ℓ_1 -norm are similar and hence are omitted. For the proposed model, we consider both JE and RE, and implement them using a fixed regularization parameter λ_g throughout the forecasting period. Five other competing approaches are considered as follows:

- (i) VAR OLS: As a low-dimensional baseline, we consider the VAR(4) model fitted via the OLS method, where the lag order 4 is employed following Koop (2013).
- (ii) VAR Lasso: Since the VAR(∞) model can be approximated by the VAR(P) with $P \to \infty$ as $T \to \infty$, we fit the sparse VAR(P) model via the Lasso with $P = \lfloor 1.5 \sqrt{T_{\text{train}}} \rfloor$ following the first-stage estimation in Wilms et al. (2023).
- (iii) VAR HLag: Same as (ii) except that the hierarchical lag (HLag) regularization in Nicholson et al. (2020) is used instead of the ℓ_1 -regularization.

- (iv) VARMA ℓ_1 : Sparse VARMA(p,q) (Wilms et al. 2023) with the ℓ_1 -regularization for the second stage and $p=q=\lfloor 0.75\sqrt{T_{\text{train}}} \rfloor$ as in the above article.
- (v) VARMA HLag: Same as (iv) except that the HLag regularization is used at the second stage.

We implement (ii)–(v) by the R package bigtime which offers two regularization parameter selection methods, cross validation (CV) and BIC. We observe that neither one of these two methods uniformly outperforms the other throughout the forecasting period. To better ensure the competitiveness of (ii)–(v), we obtain the forecast errors under both CV and BIC and only report the smaller value for each rolling step.

The average forecast error over the entire forecast period is 5.367, 4.307, 4.069, 4.318, 4.144, 3.971, and 3.968 for VAR OLS, VAR Lasso, VAR HLag, VARMA ℓ_1 , VARMA HLag, SPVAR (∞) JE, and SPVAR(∞) RE, respectively. Among the 28 rolling steps, each of these approaches performs best 4, 4, 0, 2, 2, 10, and 6 times, respectively. Thus, based on these measures, SPVAR(∞) has the highest overall forecast accuracy among all models, and the performance of JE and RE are very similar; see Table S2 in the supplementary file for the forecast errors of all seven methods for each rolling step. Moreover, to check whether the advantage of the SPVAR(∞)-based forecasts is statistically significant, we conduct the model confidence set (MCS) procedure of Hansen, Lunde, and Nason (2011) implemented by the R package MCS. We find that based on either the Tmax or TR statistic, the 97.5% MCS only includes SPVAR(∞) JE and SPVAR(∞) RE, confirming that the proposed model indeed outperforms the competing ones in terms of forecasting for the data.

7. Conclusion and Discussion

This article develops the SPVAR(∞) model as a tractable variant of the VARMA model for high-dimensional time series. It overcomes the drawbacks in identification, computation, and interpretation of the latter, while greater statistical efficiency and Granger causal interpretations are achieved by imposing



sparsity on the parameter matrices capturing the cross-sectional dependence. To the best of our knowledge, it is the first high-dimensional sparse VARMA- or $VAR(\infty)$ -type model with all of the above advantages.

There is a vast literature on nonlinear and nonstationary VAR models (e.g., Kalliovirta, Meitz, and Saikkonen 2016; Zhang and Wu 2021), factor-augmented VAR (Miao, Phillips, and Su 2022), and other extensions. The method in this article can be extended to develop corresponding $VAR(\infty)$ counterparts; for example, (2.4) can be extended to the nonlinear model: $\mathbf{y}_t = f(\mathbf{x}_t^{[1]}, \dots, \mathbf{x}_t^{[d]}) + \boldsymbol{\varepsilon}_t$, where $\mathbf{x}_t^{[k]} = \sum_{h=1}^{\infty} \ell_{h,k}(\boldsymbol{\omega}) \mathbf{y}_{t-h}$ for $1 \le k \le d$ parsimoniously summarize the temporal information over all lags into d predictors. Other interesting extensions include imposing group sparsity on G_k 's to capture groupwise homogeneity (Basu, Shojaie, and Michailidis 2015), extending $\ell_{h,k}(\omega)$'s to polynomial decay functions for long-memory time series (Chung 2002), and incorporating dynamic factor structures (Wang et al. 2022). Lastly, it is important to study the high-dimensional statistical inference under the proposed model, for example, hypothesis testing for Granger causality (Chernozhukov et al. 2021; Babii, Ghysels, and Striaukas 2022).

Supplementary Materials

Online supplementary materials contain the block coordinate descent algorithms for implementing the proposed estimation, additional simulation studies and empirical results, and all technical proofs for the main paper.

Acknowledgments

I am grateful to the Editor, Associate Editor, and two anonymous referees for their valuable comments which led to substantial improvement of this article.

Disclosure Statement

There are no competing interests to declare.

Funding

This research was partially supported by NSF grant DMS-2311178.

References

- Agarwal, A., Negahban, S., and Wainwright, M. J. (2012), "Fast Global Convergence of Gradient Methods for High-Dimensional Statistical Recovery," *The Annals of Statistics*, 40, 2452–2482. [8]
- Athanasopoulos, G., and Vahid, F. (2008), "VARMA versus VAR for Macroeconomic Forecasting," *Journal of Business & Economic Statistics*, 26, 237–252. [2]
- Babii, A., Ghysels, E., and Striaukas, J. (2022), "High-Dimensional Granger Causality Tests with an Application to VIX and News," *Journal of Finan*cial Econometrics, to appear. [13]
- Barigozzi, M., and Brownlees, C. (2017), NETS: Network Estimation for Time Series," *Journal of Applied Econometrics*, 34, 347–364. [5]
- Basu, S., and Matteson, D. S. (2021), "A Survey of Estimation Methods for Sparse High-Dimensional Time Series Models," ArXiv preprint arXiv:2107.14754. [1]
- Basu, S., and Michailidis, G. (2015), "Regularized Estimation in Sparse High-Dimensional Time Series Models," *The Annals of Statistics*, 43, 1535–1567. [6,7]

- Basu, S., Shojaie, A., and Michailidis, G. (2015), "Network Granger Causality with Inherent Grouping Structure," *Journal of Machine Learning Research*, 16, 417–453. [13]
- Chan, J. C., Eisenstat, E., and Koop, G. (2016), "Large Bayesian VARMAs," *Journal of Econometrics*, 192, 374–390. [1,2,3]
- Chernozhukov, V., Härdle, W. K., Huang, C., and Wang, W. (2021), "Lasso-Driven Inference in Time and Space," *The Annals of Statistics*, 49, 1702–1735. [13]
- Chung, C.-F. (2002), "Sample Means, Sample Autocovariances, and Linear Regression of Stationary Multivariate Long Memory Processes," *Econometric Theory*, 18, 51–78. [13]
- Costacurta, J., Duncker, L., Sheffer, B., Gillis, W., Weinreb, C., Markowitz, J., Datta, S. R., Williams, A., and Linderman., S. (2022), "Distinguishing Discrete and Continuousbehavioral Variability Using Warped Autoregressive HMMs," in Advances in Neural Information Processing Systems (Vol. 35), pp. 23838–23850. [5]
- Davis, R. A., Zang, P., and Zheng, T. (2016), "Sparse Vector Autoregressive Modeling," *Journal of Computational and Graphical Statistics*, 25, 1077–1096. [1]
- 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. [2]
- 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. [1]
- Fiecas, M. B., Coffman, C., Xu, M., Hendrickson, T. J., Mueller, B. A., Klimes-Dougan, B., and Cullen, K. R. (2023), "Approximate Hidden Semi-Markov Models for Dynamic Connectivity Analysis in Resting-State FMRi," *Statistics and Its Interface*, 16, 259–277. [5]
- Gorrostieta, C., Ombao, H., Bédard, P., and Sanes, J. N. (2012), "Investigating Brain Connectivity Using Mixed Effects Vector Autoregressive Models," *NeuroImage*, 59, 3347–3355. [1]
- Granger, C. W. (1969), "Investigating Causal Relations by Econometric Models and Cross-Spectral Methods," *Econometrica*, 37, 424–438. [1,4]
- Hansen, P. R., Lunde, A., and Nason, J. M. (2011), "The Model Condence Set," *Econometrica*, 79, 453–497. [12]
- Hartfiel, D. J. (1995), "Dense Sets of Diagonalizable Matrices," *Proceedings of the American Mathematical Society*, 123, 1669–1672. [3]
- Horn, R. A., and Johnson, C. R. (2012), *Matrix Analysis* (2nd ed.), New York: Cambridge University Press. [3]
- Huang, F., Lu, K., and Zheng, Y. (2023). "SARMA: Scalable Low-Rank High-Dimensional Autoregressive Moving Averages via Tensor Decomposition," working paper. [2,3]
- Janková, J., and van de Geer, S. (2021), "De-Biased Sparse PCA: Inference and Testing for Eigenstructures of Large Covariance Matrices," *IEEE Transactions on Information Theory*, 67, 2507–2527. [7]
- Kalliovirta, L., Meitz, M., and Saikkonen, P. (2016), "Gaussian Mixture Vector Autoregression," *Journal of Econometrics*, 192, 485–498.
- Koop, G. M. (2013), "Forecasting with Medium and Large Bayesian VARs," Journal of Applied Econometrics, 28, 177–203. [11,12]
- Krampe, J., and Paparoditis, E. (2021), "Sparsity Concepts and Estimation Procedures for High-Dimensional Vector Autoregressive Models," *Journal Time Series Analysis*, 42, 554–579. [8]
- Li, X., Safikhani, A., and Shojaie, A. (2022), "Estimation of High-Dimensional Markov-Switching VAR Models with an Approximate EM Algorithm," arXiv preprint arXiv:2210.07456. [5]
- Loh, P.-L. (2017), "Statistical Consistency and Asymptotic Normality for High-Dimensional Robust *M*-estimators," *The Annals of Statistics*, 45, 866–896. [7,8]
- 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. [1]
- Lütkepohl, H. (2005), New Introduction to Multiple Time Series Analysis, Berlin: Springer. [1,4,6]
- Metaxoglou, 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. [2]
- Miao, K., Phillips, P. C., and Su, L. (2022), "High-Dimensional VARs with Common Factors," *Journal of Econometrics*. to appear. [13]



- Nicholson, W. B., Wilms, I., Bien, J., and Matteson, D. S. (2020), "High Dimensional Forecasting via Interpretable Vector Autoregression," *Journal of Machine Learning Research*, 21, 1–52. [1,2,12]
- Shojaie, A., Basu, S., and Michailidis, G. (2012), "Adaptive Thresholding for Reconstructing Regulatory Networks from Time-Course Gene Expression Data," *Statistics in Biosciences*, 4, 66–83. [2]
- Shojaie, A., and Fox, E. B. (2021), "Granger Causality: A Review and Recent Advances," arXiv preprint arXiv:2105.02675. [4]
- Wang, D., Zheng, Y., Lian, H., and Li, G. (2022), "High-Dimensional Vector Autoregressive Time Series Modeling via Tensor Decomposition," *Journal of the American Statistical Association*, 117, 1338–1356. [1,13]
- Wang, L., and He, X. (2022), "Analysis of Global and Local Optima of Regularized Quantile Regression in High Dimensions: A Subgradient Approach," *Econometric Theory*, 1–45. [7]
- Wilms, I., Basu, S., Bien, J., and Matteson, D. (2023), "Sparse Identification and Estimation of Large-Scale Vector Autoregressive Moving Averages," *Journal of the American Statistical Association*, 118, 571–582. [1,2,3,12]
- Zhang, D., and Wu, W. B. (2021), "Convergence of Covariance and Spectral Density Estimates for High-Dimensional Locally Stationary Processes," *The Annals of Statistics*, 49, 233–254. [13]