

## ORIGINAL ARTICLE

## REGULARIZED ESTIMATION OF HIGH-DIMENSIONAL VECTOR AUTOREGRESSIONS WITH WEAKLY DEPENDENT INNOVATIONS

RICARDO P. MASINI,<sup>a</sup> MARCELO C. MEDEIROS<sup>b</sup> AND EDUARDO F. MENDES<sup>c\*</sup> 

<sup>a</sup>Center for Statistics and Machine Learning, Princeton University, Princeton, NJ, USA

<sup>b</sup>Department of Economics, Pontifical Catholic University of Rio de Janeiro, Rio de Janeiro, RJ, Brazil

<sup>c</sup>School of Applied Mathematics (EMAp), Getulio Vargas Foundation, Rio de Janeiro (FGV), Rio de Janeiro, RJ, Brazil

There has been considerable advance in understanding the properties of sparse regularization procedures in high-dimensional models. In time series context, it is mostly restricted to Gaussian autoregressions or mixing sequences. We study oracle properties of LASSO estimation of weakly sparse vector-autoregressive models with heavy tailed, weakly dependent innovations. In contrast to current literature, our innovation process satisfy an  $L^1$  mixingale type condition on the centered conditional covariance matrices. This condition covers  $L^1$ -NED sequences and strong ( $\alpha$ -) mixing sequences as particular examples.

Received 30 October 2020; Accepted 01 October 2021

Keywords: high-dimensional time series, LASSO, VAR, mixing

**JEL.** C32; C55; C58.

MOS subject classification: 62M10.

### 1. INTRODUCTION

Modeling multivariate time series data is an important and vibrant area of research. Applications range from economics and finance, as in Sims (1980), Bauer and Vornik (2011), Chiriac and Voev (2011), or Ramey (2016), to air pollution and ecological studies (Hoek *et al.*, 2013; Ensor *et al.*, 2013; Schweinberger *et al.*, 2017). Among alternatives, the vector autoregressive (VAR) model is certainly one of the most successful in modeling temporal evolution of vectors, networks, and matrices. See Lütkepohl (1991) or Wilson *et al.* (2015) for comprehensive textbook introductions.

The advances in data collection and storage have created data sets with large numbers of time series (*Big Data*), where the number of model parameters to be estimated may exceed the number of available data observations. A common approach to dealing with high-dimensional data is to impose additional structure in the form of (approximate) sparsity and estimate the parameters by some shrinkage method. Examples of estimation techniques range from Bayesian estimation with ‘spike-and-slab’ priors to sparsity-inducing shrinkage, such as the least absolute and shrinkage estimator (LASSO) and its many extensions. See Miranda-Agrippino and Ricco (2019) for a nice survey on Bayesian VARs or Kock *et al.* (2020) for a review on penalized regressions applied to time-series models.

#### 1.1. Our Contributions

In this article, we study non-asymptotic properties of high-dimensional VAR models and their parameter estimates using equation-wise (row-wise or node-wise) LASSO. We show that, with high probability, estimated and

\*Correspondence to: Eduardo F. Mendes, School of Applied Mathematics (EMAp), Getulio Vargas Foundation, Praia de Botafogo, 190, Rio de Janeiro (FGV), RJ, Brazil. E-mail: eduardo.mendes@fgv.br

population parameter vectors are close to each other in the Euclidean norm and discuss restrictions on the rate which the number of parameters can increase as the sample size diverges.

The importance of our results relies on the fact that our non-asymptotic guarantees serve as a fundamental ingredient for the derivation of asymptotic properties of penalized estimators in high-dimensional VAR models, as in Adamek *et al.* (2020). In particular, our results apply with minimal restrictions on the conditional variance model, allowing, for instance, large-dimensional multivariate linear processes in the variance. Moreover, auxiliary results proved in this article are of independent interest and can, for instance, be used to derive finite bounds for other type of penalization such as group/structured lasso, elastic-net, SCAD or non-convex penalties.

The data are assumed to be generated from a covariance-stationary and weakly sparse VAR model, where the innovations are martingale difference with sub-Weibull tails and conditional covariance matrix satisfying a  $L^1$  mixingale assumption. An important feature is that the resulting process  $\{\mathbf{y}_t\}$  is not necessarily mixing. Mixing assumptions can be notoriously difficult to show and we avoid it in this article. Nevertheless, it follows that our conditions cover strong mixing innovations as a particular case.

## 1.2. Literature Review

Some consistency results on model estimation and selection of high-dimensional VAR processes were obtained by Song and Bickel (2011), though under much stronger assumptions, such as Gaussianity. Loh and Wainwright (2012), Basu and Michailidis (2015) and Kock and Callot (2015) developed powerful concentration inequalities that enabled them to establish consistency under weaker conditions and prove that these conditions hold with high probability. In particular, Basu and Michailidis (2015) established consistency of  $\ell_1$ -penalized least squares and maximum likelihood estimators of the coefficients of high-dimensional Gaussian VAR processes and related the estimation and prediction error to the complex dependence structure of VAR processes. Other estimation approaches, including Bayesian approaches, are discussed by Davis *et al.* (2016). Miao *et al.* (2020) proposed a factor-augmented large dimensional VAR and studied finite sample properties and provide estimation results. However, they assume independent and identically distributed errors. More recently, Wong *et al.* (2020) derived finite-sample guarantees for the LASSO in a misspecified VAR model. Authors assume the series is either  $\beta$ -mixing process with sub-Weibull marginal distributions or  $\alpha$ -mixing Gaussian processes. Finally, Adamek *et al.* (2020) develop theoretical results for point estimation and inference in near epoch dependent time series using desparsified lasso under high level conditions on the generating process.

## 1.3. Organization of the Article

The article is organized as follows. In Section 2, we define the model and the main assumptions in the article. In Section 3, we discuss examples of applications of our results. The theoretical results are presented in Section 4, while in Section 5, we provide a discussion of our findings and conclude the article. All technical proofs are relegated to the Appendix.

## 1.4. Notation

Throughout the article we use the following notation. For a vector  $\mathbf{b} = (b_1, \dots, b_k)' \in \mathbb{R}^k$  and  $p \in [1, \infty]$ ,  $\|\mathbf{b}\|_p$  denotes its  $\ell_p$  norm, that is,  $\|\mathbf{b}\|_p = \left(\sum_{i=1}^k |b_i|^p\right)^{1/p}$  for  $p \in [1, \infty)$  and  $\|\mathbf{b}\|_\infty = \max_{1 \leq i \leq k} |b_i|$ . We also define  $\|\mathbf{b}\|_0 = \sum_{i=1}^k I(b_i \neq 0)$ . For a random variable  $X$ ,  $\|X\|_p = (\mathbb{E}|X|^p)^{1/p}$  for  $p \in [1, \infty)$  and  $\|X\|_\infty = \inf\{a \in \mathbb{R} : \Pr(|X| \geq a) = 0\}$ . For a  $m \times n$  matrix  $\mathbf{A}$  with elements  $a_{ij}$ , we denote  $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$ ,  $\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$ , the induced  $\ell_\infty$  and  $\ell_1$  norms respectively, and the maximum elementwise norm  $\|\mathbf{A}\|_{\max} = \max_{i,j} |a_{ij}|$ . Also  $\Lambda_{\min}(\mathbf{A})$  and  $\Lambda_{\max}(\mathbf{A})$  denotes the minimum and maximum eigenvalues of the square matrix  $\mathbf{A}$  respectively.

## 2. MODEL SETUP AND ASSUMPTIONS

Let  $\{\mathbf{y}_t = (y_{t,1}, \dots, y_{t,n})' : t \in \mathbb{Z}\}$  be a vector stochastic process defined in some fixed probability space taking values on  $\mathbb{R}^n$  given by

$$\mathbf{y}_t = \mathbf{A}_1 \mathbf{y}_{t-1} + \dots + \mathbf{A}_p \mathbf{y}_{t-p} + \mathbf{u}_t, \quad (1)$$

where  $\mathbf{u}_t = (u_{t,1}, \dots, u_{t,n})'$  is a zero-mean vector of innovations and  $\mathbf{A}_1, \dots, \mathbf{A}_p$ , are  $n \times n$  parameter matrices. The dimension  $n := n_T$  and order  $p := p_T$  of the process are allowed to increase with the number of observations  $T$ . Write the vector-autoregressive (VAR) process (1) using its first-order representation:

$$\tilde{\mathbf{y}}_t = \mathbf{F}_T \tilde{\mathbf{y}}_{t-1} + \tilde{\mathbf{u}}_t, \quad (2)$$

where  $\tilde{\mathbf{y}}_t = (\mathbf{y}'_t, \dots, \mathbf{y}'_{t-p+1})'$ ,  $\tilde{\mathbf{u}}_t = (\mathbf{u}'_t, \mathbf{0}', \dots, \mathbf{0}')'$ , and

$$\mathbf{F}_T = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \dots & \mathbf{A}_{p-1} & \mathbf{A}_p \\ \mathbf{I}_n & \mathbf{0}_n & \dots & \mathbf{0}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{I}_n & & \mathbf{0}_n & \mathbf{0}_n \\ \vdots & & \ddots & \vdots & \vdots \\ \mathbf{0}_n & \mathbf{0}_n & & \mathbf{I}_n & \mathbf{0} \end{bmatrix}.$$

In high-dimensional statistics the dimension of the model is relatively larger the number of observations  $T$ . In practical terms, it is mathematically modeled as the rate in which the dimension may increase as a function of the number of observations available, hence the dependence of  $n$  and  $p$  on  $T$ . In the following assumptions we need the data generating process to satisfy bound conditions that must hold for *any sample size*  $T$ . Although particular constants can vary with  $T$ , they have to be uniformly bounded. To simplify exposition, we simply state that constants involved do not depend on  $T$ . As our results hold for finite trajectories  $X_1, \dots, X_T$  we formulate Assumptions (A1)–(A3) in a way that they must hold ‘...uniformly in  $1 \leq t \leq T$  and  $T \in \mathbb{N}$ ’. Assumptions (A4) and (A5) involve quantities that are likely to change with the sample size and may affect regularization rate. These sequences appear explicitly in the probability bounds and rates.

**Assumption (A1).** All roots of the reverse characteristic polynomial  $\mathcal{A}(z) = \mathbf{I}_n - \sum_{j=1}^p \mathbf{A}_j z^j$  lie outside the unit disk for each  $p, n \in \mathbb{N}$ , and there exist  $\bar{c}_\Phi > 0$ ,  $c_\Phi > 0$  and  $0 < \gamma_1 \leq 1$  such that for all  $m \in \mathbb{N}$

$$\sum_{k=m}^{\infty} |\phi_{k,i}|_1 \leq \bar{c}_\Phi e^{-c_\Phi m^{\gamma_1}}, \quad (3)$$

uniformly in  $1 \leq i \leq n$ , where  $\Phi_k := \mathbf{J}' \mathbf{F}_T^k \mathbf{J} = (\phi_{k,1}, \dots, \phi_{k,n})'$ ,  $\mathbf{F}_T$  denote the companion matrix and  $\mathbf{J} = (\mathbf{I}_n, \mathbf{0}_n, \dots, \mathbf{0}_n)'$ .

**Assumption (A2).** The sequence  $\{(\mathbf{u}_t, \mathcal{F}_t)\}_t$  is a covariance stationary (for each  $T \in \mathbb{N}$ ) martingale difference process where the filtration  $\{\mathcal{F}_t\}_t$  includes the natural filtration of  $\{\mathbf{u}_t\}$ . The smallest and largest eigenvalues of  $\Sigma := \mathbb{E}(\mathbf{u}_t \mathbf{u}_t')$  are bounded away from 0 and  $\infty$  respectively, uniformly in  $T \in \mathbb{N}$ . Furthermore, for all  $\mathbf{b}_1, \mathbf{b}_2 \in \{\mathbf{v} \in \mathbb{R}^n : |\mathbf{v}|_1 \leq 1\}$  and  $m \in \mathbb{N}$ ,

$$\mathbb{E} \left| \mathbb{E} [\mathbf{b}'_1 (\mathbf{u}_t \mathbf{u}_t' - \Sigma) \mathbf{b}_2 | \mathcal{F}_{t-m}] \right| \leq a_1 e^{-a_2 m^{\gamma_2}},$$

for some  $a_1, a_2 > 0$  and  $0 < \gamma_2 \leq 1$ , uniformly in  $1 \leq t \leq T$  and  $T \in \mathbb{N}$ .

**Assumption (A3).** For all  $\mathbf{b} \in \{\mathbf{v} \in \mathbb{R}^n : |\mathbf{v}|_1 \leq 1\}$  and all  $0 < x < \infty$ ,  $\Pr(|\mathbf{b}'\mathbf{u}_t| > x) \leq 2e^{-x/c_\alpha} \alpha$  for some  $\alpha > 0$ ,  $0 < c_\alpha < \infty$ , uniformly in  $1 \leq t \leq T$  and  $T \in \mathbb{N}$ .

Assumption (A1) requires that the VAR process is stable and admits an infinite-order vector moving average,  $\text{VMA}(\infty)$ , representation for all  $n$  and  $p$  as

$$\mathbf{y}_t = \sum_{i=0}^{\infty} \mathbf{J}' \mathbf{F}_T^i \mathbf{J} \mathbf{u}_{t-i} = \sum_{i=0}^{\infty} \mathbf{\Phi}_i \mathbf{u}_{t-i}. \quad (4)$$

Furthermore, the coefficients of the  $\text{MA}(\infty)$  representations of each  $\{y_{i,t}\}$ ,  $i = 1, \dots, n$ , are absolutely summable with exponentially decaying rate. This condition is satisfied in standard  $\text{VAR}(p)$  models, where  $n$  and  $p$  are fixed. In models that  $n$  is large, Lemma 4 in Appendix B.1 shows that condition (3) is satisfied if  $\sum_{k=1}^p \|\mathbf{A}_k\|_\infty < 1$  and further regularity conditions on the size of the coefficients. Finally, notice that under (A1) it is also true that  $\max_{k,i} |\phi_{k,i}|_\infty \leq \bar{c}_\Phi$ , which means that the coefficients  $\{\mathbf{\Phi}_k\}$  are uniformly upper bounded under the maximum entry-wise norm.

Assumption (A2) requires the error process to be a martingale difference process and satisfy a very weak dependence condition on its conditional variance. The former restricts the model to be correctly specified in the mean. Nevertheless, this assumption is standard in the literature and we are able to derive results covering a broad range of data generating processes and conditional dependence measures. The latter is the  $L^1$  projective dependence measure appearing in Dedecker *et al.* (2007, section 2.2.4). Note that (1) strong mixing (or  $\alpha$ -mixing) sequences with exponential decay of the mixing coefficient satisfy this condition (Davidson, 1994, Theorem 14.2); and (2) uniform mixing sequences ( $\phi$ -mixing) and  $\beta$ -mixing sequences are also strong mixing, but the converse is not true (Bradley, 2005, Eqs. (1.11)–(1.18)). If we denote the centered outer product series  $\mathbf{v}_t = \text{vech}(\mathbf{u}_t \mathbf{u}_t' - \mathbf{\Sigma})$ , this assumption requires that  $\{\mathbf{v}_t\}$  is  $L^1$  mixingale. It means that stochastic process with  $L^1$  bounded,  $L^1$  near-epoch dependent, centered outer product series  $\mathbf{v}_t$  are also contemplated in this setting Andrews (1988). Finally, Assumptions (A1) and (A2) combined ensure that  $\{\mathbf{y}_t\}$  is second-order stationary for each  $n$  and  $p$  (Lütkepohl, 2006, Ch. 2).

Condition (A3) imposes restrictions on the tail behavior of the innovation process  $\{\mathbf{u}_t\}$  that are shared by  $\{\mathbf{y}_t\}$ . More precisely, we impose moment conditions on all linear combinations  $\mathbf{b}'\mathbf{u}_t$ . Lemma 3, in the appendix, shows that each  $\{y_{i,t}\}$  ( $i = 1, \dots, n$ ) also share the same tail properties of  $\{\mathbf{u}_t\}$ . This condition is essential for defining the rate in which  $n$  and  $p$  increase with  $T$ . We focus on the case the tail decays at rate  $O(e^{-c_\alpha x^\alpha})$  for some  $\alpha > 0$ , that is,  $\{\mathbf{b}'\mathbf{u}_t\}$  is sub-Weibull with parameter  $\alpha$  studied in Wong *et al.* (2020, Section 4.1). Note that when  $\alpha \geq 1$  and  $\alpha \geq 2$  we have the sub-exponential and sub-Gaussian tails respectively. However, when  $\alpha \in (0, 1)$  the moment generating function does not exist at any point and these variables are usually called *heavy tailed*.

It is convenient to write the model in stacked form. Let  $\mathbf{x}_t = (\mathbf{y}'_{t-1}, \dots, \mathbf{y}'_{t-p})'$  be the  $np \times 1$  vector of regressors and  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)'$  the  $T \times np$  matrix of covariates. Let  $\mathbf{Y}_i = (y_{i,1}, \dots, y_{i,T})'$  be the  $T \times 1$  vector of observations for the  $i$ th element of  $\mathbf{y}_t$ , and  $\mathbf{U}_i = (u_{i,1}, \dots, u_{i,T})'$  the corresponding vector of innovations. Denote  $\boldsymbol{\beta}_i$  the  $np \times 1$  vector of coefficients corresponding to equation  $i$ . Then, model (1) is equivalent to

$$\mathbf{Y}_i = \mathbf{X} \boldsymbol{\beta}_i + \mathbf{U}_i, \quad i = 1, \dots, n. \quad (5)$$

We now make additional assumptions concerning model (5).

**Assumption (A4).** The true parameter vectors  $\boldsymbol{\beta}_i$ ,  $i = 1, \dots, n$ , satisfy  $\sum_{j=1}^{np} |\beta_{i,j}|^q \leq R_q$  for some  $0 \leq q < 1$  and  $0 < R_q < \infty$  where  $R_q := R_{q,T}$  is allowed to depend on the sample size  $T$ .

**Assumption (A5).** For each  $T \in \mathbb{N}$ , the smallest eigenvalue of  $\boldsymbol{\Gamma} := T^{-1} \mathbb{E}(\mathbf{X}'\mathbf{X})$  is greater than a positive constant  $\sigma_\Gamma^2$  that might depend on  $T$ .

Assumption (A4) imposes *weak sparsity* of the coefficients, in a sense that most of them are small. This condition is slightly stronger than we need in a sense that we may have distinct  $q_i$  and  $R_{q,i}$  for each equation. In the case  $q = 0$  we have sparsity in the standard sense, meaning that  $R_0 = s$ , the number of non-zero coefficients. In practice, we estimate a sparse model that truncates all coefficients close to zero. This assumption is standard for *weak sparsity*, see [Negahban et al. \(2012, section 4.3\)](#) and [Han and Tsay \(2020, Assumption 1\)](#) for an application in time series setting.

Assumption (A5) is often used in the sparse estimation literature (e.g. [Kock and Callot, 2015](#); [Medeiros and Mendes, 2016](#); [Han and Tsay, 2020](#)). [Basu and Michailidis \(2015, Proposition 2.3\)](#) derived bounds for  $\Lambda_{\min}(\mathbf{\Gamma})$  and  $\Lambda_{\max}(\mathbf{\Gamma})$  using properties of the block Toeplitz matrix  $\mathbf{\Gamma}$  and its generating function, the cross-spectral density of the generating VAR( $p$ ) process:

$$\frac{\Lambda_{\min}(\mathbf{\Sigma})}{\max_{|z|=1} \Lambda_{\max}(\mathcal{A}^*(z)\mathcal{A}(z))} \leq \Lambda_{\min}(\mathbf{\Gamma}) \leq \Lambda_{\max}(\mathbf{\Gamma}) \leq \frac{\Lambda_{\max}(\mathbf{\Sigma})}{\max_{|z|=1} \Lambda_{\min}(\mathcal{A}^*(z)\mathcal{A}(z))}, \tag{6}$$

where  $\mathcal{A}^*$  is the conjugate transpose of  $\mathcal{A}$ , the reverse characteristic polynomial, defined in Assumption (A1). [Basu and Michailidis \(2015\)\[Proposition 2.2\]](#) shows that under (A1),

$$\max_{|z|=1} \Lambda_{\max}(\mathcal{A}^*(z)\mathcal{A}(z)) < \left[ 1 + \frac{\sum_{k=1}^p (\|\mathbf{A}_k\|_1 + \|\mathbf{A}_k\|_\infty)}{2} \right]^2.$$

Hence, (A5) is satisfied if, for instance,  $\Lambda_{\min}(\mathbf{\Sigma}) > 0$ ,  $\sum_{k=1}^p \|\mathbf{A}_k\|_1 < \infty$  and  $\sum_{k=1}^p \|\mathbf{A}_k\|_\infty < \infty$ .

### 3. EXAMPLES

We illustrate processes satisfying Assumptions (A2) and (A3). In the first two examples we discuss sufficient conditions involving mixing and near epoch dependent sequences, traditionally found in the literature. In the final two examples, we discuss variance process admitting an AR( $\infty$ ) representation. In all examples, the innovation vector  $\mathbf{u}_t$  has dimension  $n$ , which is a function of  $T$ , but bounds are independent on the dimension of the problem or sample size. Hence, following examples hold for *each*  $n$ , as in the assumptions.

**Example 1** (Strong mixing sequences). Let  $\{\mathbf{u}_t\}$  denote a martingale difference, strong mixing sequence with coefficients  $\alpha_m < b_1 \exp(-b_2 m^{\gamma_2})$  and common covariance matrix  $\mathbf{\Sigma}$  with eigenvalues bounded away from zero and infinity, uniformly in  $n$ . It follows that  $r_t = \mathbf{b}'_1 \mathbf{u}_t \mathbf{u}'_t \mathbf{b}_2$  is also strong mixing of same size and, from [Davidson \(1994, Theorem 14.2\)](#),  $\mathbb{E}[r_t - \mathbb{E}(r_t) | \mathcal{F}_{t-m}] \leq a_1 \exp(-a_2 m^{\gamma_2})$ , for constants  $a_1$  and  $a_2$ .

**Example 2** ( $L^1$  near-epoch dependent process). Let  $\{\mathbf{u}_t\}$  denote a weakly stationary, martingale difference sequence. Suppose  $\mathbf{b}'\mathbf{v}_t = \mathbf{b}'\text{vech}(\mathbf{u}_t \mathbf{u}'_t - \mathbf{\Sigma})$  is a centered,  $L^1$ -NED sequence on  $\mathcal{F}_t = \sigma(\epsilon_t, \epsilon_{t-1}, \dots)$ , where  $\{\epsilon_t\}$  is  $\alpha$ -mixing with coefficients  $\alpha_m \leq c_1 \exp(-c_2 m^{\gamma_1})$ , for all  $\mathbf{b} \in \{\mathbf{b} \in \mathbb{R}^{n(n+1)/2} : \|\mathbf{b}\|_1 \leq 1\}$ . It means that there are finite constants  $\{d_t\}$  and  $\{\psi_m\}$  such that

$$\mathbb{E} \left| \mathbf{b}'(\mathbf{v}_t - \mathbb{E}[\mathbf{v}_t | \mathcal{F}_{t-m:t}]) \right| \leq d_t \psi_m,$$

where  $\mathcal{F}_{t-m:t} = \sigma(\epsilon_t, \dots, \epsilon_{t-m})$  and  $\psi_m \leq \exp(-c_3 m^{\gamma_2})$ . Under Assumption (A3), it follows from [Wong et al. \(2020, Lemma 5\)](#) and Hölder inequality that for any  $r < \infty$

$$\|\mathbf{b}'\mathbf{v}_t\|_r \leq \|\mathbf{b}\|_1^r \max_{1 \leq i \leq j \leq n} \|u_{it} u_{jt}\|_r \leq \max_{1 \leq i \leq n} \|u_{it}\|_{2r} \leq c_4 r^{1/\alpha}.$$

Finally, it follows from [Andrews \(1988, Example 6\)](#) that Assumption (A2) holds with  $a_1 \geq (2 \max_t d_t + c_4 r^{1/\alpha}) (e^{c_3/2\gamma_2} + 6c_1 e^{c_2(r-1)/r2\gamma_2})$  and  $a_2 \leq (c_3 \wedge c_2(r-1)/r) / 2\gamma_2$ .

**Example 3** (Linear process in the variance). Let  $\{v_t, \epsilon_t\}$  denote a sequence of centered independently and identically distributed, sub-Weibull random variables with parameter (at least)  $2\alpha$ , taking values in  $\mathbb{R}^{2n}$  with identity covariance matrix. Let  $u_t = H_t^{1/2}v_t$  where  $H_t^{1/2}$  is the lower diagonal Cholesky decomposition of  $H_t$  and

$$h_t = \text{vech}(H_t) = c + \sum_{j=1}^{\infty} \Psi_j \eta_{t-j}.$$

Here,  $\text{vech}(M)$  stacks the lower diagonal elements of matrix  $M$ ,  $c$  is a vector of constants and  $\eta_t = \text{vech}(\epsilon_t \epsilon_t')$ . For all  $\tilde{b} \in \{b \in \mathbb{R}^{n(n+1)/2} : |b|_1 \leq 1\}$ ,  $\{\Psi_j\}$  satisfy  $\sum_{j=m}^{\infty} |\tilde{b}' \Psi_j|_1 \lesssim e^{-\alpha_2 m^{1/2}}$ .

We first show  $\{u_t\}$  is weakly stationary martingale difference with respect to  $\mathcal{F}_{t-1} = \sigma(v_{t-j}, \epsilon_{t-j} : j = 1, 2, \dots)$ . The process  $\{u_t\}$  is  $\mathcal{F}_t$  measurable and satisfy  $\mathbb{E}[u_t | \mathcal{F}_{t-1}] = H_t^{1/2} \mathbb{E}[v_t | \mathcal{F}_{t-1}] = 0$ . Its covariance matrix is

$$\mathbb{E}[u_t u_t'] = \mathbb{E}\left[H_t^{1/2} \mathbb{E}(v_t v_t' | \mathcal{F}_{t-1}) (H_t^{1/2})'\right] = \mathbb{E}[H_t].$$

Now,  $\mathbb{E}[h_t] = c + \sum_{j=1}^{\infty} \Psi_j \mathbb{E} \eta_{t-j} = c + \sum_{j=1}^{\infty} \Psi_j \text{vech}(I_n) = \Sigma$ , where  $\mathbb{E}(\eta_t) = \text{vech}(\mathbb{E}(\epsilon_t \epsilon_t')) = \text{vech}(I_n)$  for all  $t$ .

For constant vectors  $b_1, b_2 \in \{b \in \mathbb{R}_n : |b| \leq 1\}$ ,

$$\begin{aligned} \mathbb{E}[b_1' (u_t u_t' - \Sigma) b_2 | \mathcal{F}_{t-1}] &= b_1' (H_t - \mathbb{E}H_t) b_2 \\ &= \tilde{b}' (h_t - \mathbb{E}h_t) \\ &= \sum_{j=1}^{\infty} \tilde{b}' \Psi_j (\eta_{t-j} - \mathbb{E}\eta_{t-j}), \end{aligned}$$

where  $\tilde{b} \in \{b \in \mathbb{R}^{n(n+1)/2} : |b|_1 \leq 1\}$ . It follows that

$$\begin{aligned} \mathbb{E}\left|\mathbb{E}[b_1' (u_t u_t' - \Sigma) b_2 | \mathcal{F}_{t-m}]\right| &= \mathbb{E}\left|\sum_{j=1}^{\infty} \tilde{b}' \Psi_j \mathbb{E}(\eta_{t-j} - \mathbb{E}\eta_{t-j} | \mathcal{F}_{t-m})\right| \\ &= \left\|\sum_{j=m}^{\infty} \tilde{b}' \Psi_j (\eta_{t-j} - \mathbb{E}\eta_{t-j})\right\|_1 \\ &\leq 2 \left(\sum_{j=m}^{\infty} |\tilde{b}' \Psi_j|_1\right) \max_{|b|_1 \leq 1} \|b' \epsilon_t\|_2^2, \end{aligned}$$

where in the last line we use the same arguments of Lemma 3 in the appendix, followed by the triangle inequality. Then, Assumption (A2) is satisfied under the condition that  $\sum_{j=m}^{\infty} |\tilde{b}' \Psi_j|_1 \lesssim e^{-\alpha_2 m^{1/2}}$  and  $\|b' \epsilon_t\|_2 \leq c_2 < \infty$ .

It follows from Wong *et al.* (2020, Lemma 5) that  $\{u_t\}$  is sub-Weibull with parameter  $\alpha$  if  $\sup_{d \geq 1} d^{-1/\alpha} \|b' u_t\|_d \leq c_\alpha < \infty$ . For any  $d \geq 1$ ,

$$\begin{aligned} d^{-1/\alpha} \|b' u_t\|_d &= d^{-1/\alpha} \|b' u_t u_t' b\|_{d/2}^{1/2} \\ &= d^{-1/\alpha} \left\|b' H_t^{1/2} v_t v_t' (H_t^{1/2}) b\right\|_{d/2}^{1/2} \\ &= d^{-1/\alpha} \left\|b' H_t b \times \frac{b' H_t^{1/2} v_t v_t' (H_t^{1/2}) b}{b' H_t b}\right\|_{d/2}^{1/2} \end{aligned}$$

$$\begin{aligned} &\leq d^{-1/\alpha} \left\{ \mathbb{E} \left( \left| \mathbf{b}' \mathbf{H}_t \mathbf{b} \right|^{d/2} \mathbb{E} \left[ \left| \frac{\mathbf{b}' \mathbf{H}_t^{1/2} \mathbf{v}_t \mathbf{v}_t' (\mathbf{H}_t^{1/2}) \mathbf{b}}{\mathbf{b}' \mathbf{H}_t \mathbf{b}} \right|^{d/2} \middle| \mathcal{F}_{t-1} \right] \right) \right\}^{1/d} \\ &\leq d^{-1/\alpha} \left\{ \mathbb{E} \left( \left| \mathbf{b}' \mathbf{H}_t \mathbf{b} \right|^{d/2} \sup_{\delta' \delta=1} \mathbb{E} \left[ (\delta' \mathbf{v}_t \mathbf{v}_t' \delta)^{d/2} \middle| \mathcal{F}_{t-1} \right] \right) \right\}^{1/d} \\ &= d^{-1/2\alpha} \left\| \mathbf{b}' \mathbf{H}_t \mathbf{b} \right\|_{d/2}^{1/2} \sup_{\delta' \delta=1} d^{-1/\alpha} \|\delta' \mathbf{v}_t\|_d, \\ &\leq \left( d^{-1/\alpha} \left\| \mathbf{b}' \mathbf{H}_t \mathbf{b} \right\|_{d/2} \right)^{1/2} c_{2\alpha} \end{aligned}$$

where the two last lines follow because  $\{\mathbf{v}_t\}$  is independent and sub-Weibull process with parameter  $2\alpha$ . There is a  $\tilde{\mathbf{b}} \in \{\mathbf{b} \in \mathbb{R}^{n(n-1)/2} : |\mathbf{b}|_1 \leq 1\}$  such that

$$\begin{aligned} d^{-1/\alpha} \left\| \mathbf{b}' \mathbf{H}_t \mathbf{b} \right\|_{d/2} &= d^{-1/\alpha} \left\| \tilde{\mathbf{b}}' \mathbf{h}_t \right\|_{d/2} \\ &= d^{-1/\alpha} \left\| \tilde{\mathbf{b}}' \mathbf{c} + \sum_{j=1}^{\infty} \tilde{\mathbf{b}}' \Psi_j \eta_{t-j} \right\|_{d/2} \\ &\leq d^{-1/\alpha} \left| \tilde{\mathbf{b}}' \mathbf{c} \right|_{d/2} + d^{-1/\alpha} \left\| \sum_{j=1}^{\infty} \tilde{\mathbf{b}}' \Psi_j \eta_{t-j} \right\|_{d/2} \\ &\leq d^{-1/\alpha} \left| \tilde{\mathbf{b}}' \mathbf{c} \right\|_{d/2} + 2 \left( \sum_{j=1}^{\infty} |\tilde{\mathbf{b}}' \Psi_j|_1 \right) \max_{|\mathbf{b}|_1 \leq 1} (d^{-1/2\alpha} \|\mathbf{b}' \epsilon_t\|_d)^2 \\ &\leq \left| \tilde{\mathbf{b}}' \mathbf{c} \right|_{d/2} + c_{2\alpha}^2. \end{aligned}$$

Combining these bounds, process  $\{\mathbf{u}_t\}$  is sub-Weibull with parameter  $\alpha$ , satisfying Condition (A3).

**Example 4** (Stochastic covariance). Let

$$\mathbf{y}_t = \sum_{i=1}^p A_i \mathbf{y}_{t-i} + \mathbf{H}_t^{1/2} \mathbf{v}_t, \quad \mathbf{H}_{t+1} = \mathbf{C}_0 + \Psi \mathbf{H}_t \Psi' + \epsilon_t \epsilon_t',$$

where  $\mathbf{v}_t \stackrel{\text{i.i.d.}}{\sim} \mathbf{N}(\mathbf{0}, \mathbf{I}_n)$ ,  $\epsilon_t \stackrel{\text{i.i.d.}}{\sim} (\mathbf{0}, \mathbf{I}_n)$  is sub-Weibull with parameter  $2\alpha$ , and processes  $\{\mathbf{v}_t\}$  and  $\{\epsilon_t\}$  are independent. Let  $\mathcal{F}_t = \sigma(\mathbf{v}_{t-j}, \epsilon_{t-j} : j = 0, 1, 2, \dots)$ .

Process  $\{\mathbf{H}_t\}$  is a matrix process characterizing the stochastic covariance of  $\mathbf{u}_t = \mathbf{H}_t^{1/2} \mathbf{v}_t$  evolves according to a matrix autoregressive process. The *intercept*  $\mathbf{C}_0$  is symmetric and positive definite matrix and the eigenvalues of  $\Psi$  are inside the unity circle. We also assume  $\|\Psi\|_1 \|\Psi\|_{\infty} < 1$ , which implies that the largest singular value of  $\Psi$  is smaller than one.

Under these conditions, the process  $\{\mathbf{H}_t\}$  is ensured to be positive definite and stationary. To see the latter, vectorize the process to obtain  $\mathbf{h}_t = \text{vec}(\mathbf{H}_t)$ ,  $\mathbf{c}_0 = \text{vec}(\mathbf{C}_0)$ ,  $\boldsymbol{\eta}_t = \text{vec}(\epsilon_t \epsilon_t')$ ,  $\tilde{\Psi} = \Psi \otimes \Psi'$  and

$$(\mathbf{I} - \tilde{\Psi} L) \mathbf{h}_{t+1} = \mathbf{c}_0 + \boldsymbol{\eta}_t \Leftrightarrow \mathbf{h}_{t+1} = \sum_{j=0}^{\infty} \tilde{\Psi}^j \mathbf{c}_0 + \sum_{j=0}^{\infty} \tilde{\Psi}^j \boldsymbol{\eta}_{t-j}.$$

The process is stationary because eigenvalues of  $\bar{\Psi}$  are products of eigenvalues of  $\Psi$ , which is also inside the unity circle. We are exactly in the setting of previous example.

We have to show  $\sum_{j=m+1}^{\infty} |\mathbf{b}'\bar{\Psi}^j|_1 \lesssim e^{-a_2 m^{\gamma_2}}$  for all  $\mathbf{b} \in R^{n^2} : |\mathbf{b}|_1 = 1$ . The term on the left-hand side is bounded by  $\sum_{j>m} \|\bar{\Psi}^j\|_1$  and each  $\|\bar{\Psi}^j\|_1 \leq \|\Psi^j\|_1 \|\Psi^j\|_{\infty} \leq (\|\Psi\|_1 \|\Psi\|_{\infty})^j$ . Under the assumption that  $\|\Psi\|_1 \|\Psi\|_{\infty} = c_{\max} < 1$ ,  $\sum_{j=m+1}^{\infty} |\mathbf{b}'\bar{\Psi}^j|_1 \leq \frac{c_{\max}}{1-c_{\max}} e^{-m \log(1/c_{\max})}$  meaning that condition is satisfied with  $\gamma_2 = 1$  and  $a_2 = \log(1/c_{\max})$ .

Finally, the unconditional covariance of  $\mathbf{u}_t$  is

$$\Sigma := \mathbb{E}[H_t] = \sum_{j=0}^{\infty} \Psi^j (C_0 + I_n) \Psi^{j'}$$

The smallest eigenvalue  $\min_{\delta' \delta = 1} \delta' \Sigma \delta \geq \Lambda_{\min}(C_0) + 1$  and largest eigenvalue of  $\max_{\delta' \delta = 1} \delta' \Sigma \delta \leq (\rho(C_0) + 1) (1 - \rho(\Psi)^2)^{-1}$ , where  $\Lambda_{\min}(A)$  is the smallest eigenvalue of  $A$  and  $\rho(A)$  is the spectral radius of  $A$ .

#### 4. LASSO ESTIMATION BOUNDS

Let  $\mathcal{L}_T(\beta_i) = \frac{1}{T} |\mathbf{Y}_i - \mathbf{X}\beta_i|_2^2$  denote the empirical squared risk, for each  $i = 1, \dots, n$ . We estimate  $\beta_i, i = 1, \dots, n$ , equation-wise using the LASSO procedure

$$\hat{\beta}_i \in \arg \min_{\beta_i \in \mathbb{R}^{np}} \{ \mathcal{L}_T(\beta_i) + \lambda_i |\beta_i|_1 \}, \quad i = 1, \dots, n, \tag{7}$$

where  $\lambda_i$  are positive regularization parameters. For ease of exposition we assume  $\lambda_1 = \dots = \lambda_n = \lambda$ . It is well known that  $\beta_i^* = \arg \min_{\beta_i} \mathbb{E} \{ \mathcal{L}_T(\beta_i) \}$  are the population parameters in (5), under stated conditions.

We follow the steps in Negahban *et al.* (2012) to derive error bounds for the equation-wise LASSO estimator. First define the pair of subspaces  $\mathcal{M}(S) = \{ \mathbf{u} \in \mathbb{R}^{np} | u_i = 0, i \in S^c \}$  and its orthogonal complement  $\mathcal{M}^{\perp}(S) = \{ \mathbf{u} \in \mathbb{R}^{np} | u_i = 0, i \in S \}$ , where  $S \subseteq \{1, \dots, np\}$ . Set  $\mathbf{u}_{\mathcal{M}}$  and  $\mathbf{u}_{\mathcal{M}^{\perp}}$  the projection of  $\mathbf{u}$  on  $\mathcal{M}(S)$  and  $\mathcal{M}^{\perp}(S)$  respectively. Clearly, for any  $\mathbf{u} \in \mathbb{R}^{np}$ ,  $|\mathbf{u}|_1 = |\mathbf{u}_{\mathcal{M}}|_1 + |\mathbf{u}_{\mathcal{M}^{\perp}}|_1$ . We say  $|\cdot|_1$  is decomposable with respect to the pair  $(\mathcal{M}(S), \mathcal{M}^{\perp}(S))$  for any set  $S \subset \{1, \dots, np\}$ .

We have to show two conditions to obtain a finite sample estimation error bound for the parameter vectors. The first condition is known as *restricted strong convexity* (RSC) and restricts the geometry of the loss function around the optimum  $\beta^*$  and is related to the Restricted Eigenvalue (Van De Geer and Bühlmann, 2009). The second condition is known as *deviation bound* and restricts the size of the sup-norm of the gradient  $\nabla \mathcal{L}_T(\beta^*)$ . These conditions are shown to be satisfied in a set of large probability defined in Propositions 1 and 2.

**Definition 1** (Deviation bound (DB)). The *deviation bound* condition holds when the regularization parameter  $\lambda$  satisfies  $\{ \lambda \geq 2|\mathbf{X}'\mathbf{U}_i/T|_{\infty} \}$  for all  $i = 1, \dots, n$ .

Note that one may adopt individual  $\lambda_i$ s for each equation, in which above definition should be modified adequately.

**Definition 2** (Restricted Strong Convexity (RSC)). Define  $\mathbb{C}(\beta^*, \mathcal{M}, \mathcal{M}^{\perp}) = \{ \Delta \in \mathbb{R}^{np} | |\Delta_{\mathcal{M}^{\perp}}|_1 \leq 3|\Delta_{\mathcal{M}}|_1 + 4|\beta_{\mathcal{M}^{\perp}}^*|_1 \}$ . The *restricted strong convexity* holds for parameters  $\kappa_{\mathcal{L}}$  and  $\tau_{\mathcal{L}}$  if for any  $\Delta \in \mathbb{C}$ ,

$$\frac{\Delta' \mathbf{X}' \mathbf{X} \Delta}{T} \geq \kappa_{\mathcal{L}} |\Delta|_2^2 - \tau_{\mathcal{L}}^2 (\beta^*).$$



Negahban *et al.* (2012, Section 4) show these conditions are satisfied by many loss functions and penalties. Basu and Michailidis (2015) show that both DB and RSC are satisfied by Gaussian VAR( $p$ ) models in high dimensions.

If both DB and RSC hold with large probability, Negahban *et al.* (2012, Theorem 1) provides an  $\ell_2$  estimation bound for  $\hat{\beta}_i$ . Our goal is to show that the error bounds are valid for each  $\Delta_i = \hat{\beta}_i - \beta_i^*$ ,  $i = 1, \dots, n$  at the same time.

Lemma 1 characterizes the solutions of the optimization program in (7). We require further notation. Define  $\mathbb{C}_i := \mathbb{C}(\beta_i^*, \mathcal{M}_{i,\eta}, \mathcal{M}_{i,\eta}^\perp)$  for a pair of subsets  $\mathcal{M}_{i,\eta} = \mathcal{M}(S_{i,\eta})$  and  $\mathcal{M}_{i,\eta}^\perp = \mathcal{M}^\perp(S_{i,\eta})$ , where  $S_{i,\eta} = \{j \in \{1, \dots, pn\} \mid |\beta_{i,j}| > \eta\}$  and  $S_{i,\eta}^c = \{j \in \{1, \dots, pn\} \mid |\beta_{i,j}| \leq \eta\}$ . These sets represent the *active parameters* under weak sparsity. In Theorem 1 we set  $\eta = \lambda/\sigma_\Gamma^2$  to derive our results.

**Lemma 1.** Suppose  $\{\mathbf{y}_i\}$  is generated from (1) and Assumptions (A1)–(A3) are satisfied. Set

$$\lambda > \tau^* (\epsilon + \log(Tn^2p))^{2/\alpha} \sqrt{\frac{\epsilon + \log(n^2p)}{T}}, \quad (8)$$

where  $\epsilon > 0$  and  $\tau^* > 0$  depends on  $\tau$ ,  $\alpha$  and  $\bar{c}_\Phi$ . Then, if  $T > \epsilon + \log(n^2p)$ , the event  $\{\forall i = 1, \dots, n : \hat{\beta}_i - \beta_i^* \in \mathbb{C}_i\}$  holds with probability at least  $1 - 10e^{-\epsilon}$ .

Lemma 1 shows that under restrictions on  $\lambda$  the solutions to the optimization program in (7) lie inside the star-shaped sets  $\mathbb{C}_i$  with high probability, as the sample size increases. It restricts the directions in which we should control the variation of our estimators. Next result shows the deviation bound holds with high probability for appropriate choice of  $\lambda$ . To formalize the idea, let

$$\mathcal{D}_i(\lambda) = \left\{ \lambda \geq 2 \left| \frac{1}{T} \mathbf{X}' \mathbf{U}_i \right|_\infty \right\}, \quad i = 1, \dots, n, \quad (9)$$

denote the event ‘DB holds for equation  $i$  with regularization parameter  $\lambda$ ’.

**Proposition 1** (Deviation bound). Suppose that  $\{\mathbf{y}_i\}$  is generated from (1), Assumptions (A1)–(A3) are satisfied and  $T > \epsilon + \log(n^2p)$  for some  $\epsilon > 0$ . Set the penalty parameters  $\lambda$  as in (8). Then,

$$\Pr \left( \bigcup_{i=1}^n \mathcal{D}_i^c(\lambda) \right) \leq \pi_1(\epsilon) := 10e^{-\epsilon}.$$

Suppose  $\epsilon = \log(np)$ ,  $n^2p > T$ . The regularization parameter  $\lambda$  satisfies

$$\lambda \gtrsim [\log(np)]^{2/\alpha} \sqrt{\frac{\log(np)}{T}},$$

and  $\pi_1(\lambda) \propto 1/n^2p$ . This regularization parameter is  $O([\log(np)]^{2/\alpha})$  larger, in rate, than one obtained in Wong *et al.* (2020, Proposition 7). Their results relied heavily in  $\mathbf{y}_i$  being a  $\beta$ -mixing sequence in a sense that the concentration inequality derived in Merlevède *et al.* (2011) depends on it. In our case, the dependence is characterized by the conditional variance of the innovation process and coefficients  $\Phi_1, \Phi_2, \dots$ , and we are not aware of ‘tight’ concentration inequalities that hold under these assumptions. Nevertheless, for fixed  $n$ , it is possible to show that the concentration inequality for sub-Weibull martingales in Lemma 5 is tight (Fan *et al.*, 2012b).

Let  $\Gamma_T = \mathbf{X}'\mathbf{X}/T$  denote the scaled Gram matrix and  $\Gamma$  its expected value. We show that if each element in  $\Gamma_T$  is sufficiently close to its expectation, and Assumptions (A4) and (A5) hold, then RSC is satisfied with high probability.

**Lemma 2** (Restricted strong convexity). Suppose Assumptions (A4) and (A5) hold and that  $\|\|\|\Gamma_T - \Gamma\|\|\|_{\max} \leq \frac{\sigma_\Gamma^2 \eta^q}{64R_q}$ . Then, for any  $\Delta_i \in \mathbb{C}_i$  for  $1 \leq i \leq n$

$$\Delta_i' \Gamma_T \Delta_i \geq \frac{\sigma_\Gamma^2}{2} |\Delta_i|_2^2 - \frac{\sigma_\Gamma^2}{2} R_q \eta^{2-q}. \quad (10)$$

To show RSC holds with high probability for all  $i = 1, \dots, n$  at the same time, we have to bound the event

$$B(a) = \{ \|\|\|\Gamma_T - \Gamma\|\|\|_{\max} \leq a \}. \quad (11)$$

where  $a = \frac{\sigma_\Gamma^{2(1-q)} \lambda^q}{64R_q}$ . If we assume distinct  $R_{q,i}$  and  $q_i$  for each equation, we should work with  $\cap_i B_i$  and  $B_i$  defined accordingly.

**Proposition 2.** Suppose Assumptions (A1)–(A3) hold. If

$$p < \frac{T^{\gamma_1 \wedge \gamma_2}}{\left(\frac{2}{\gamma_1 \wedge \gamma_2 + 1}\right) \left(2 + \frac{1.4}{2\gamma_1 c_\phi \wedge a_2}\right)},$$

and

$$a \geq \sqrt{\frac{2(1 + \xi)^{1+2/\alpha} \tau^2 [\log(npT)]^{1+2/\alpha}}{T}},$$

for some  $\xi > 0$ , then  $\Pr(B^c(a)) \leq \pi_2(a)$ , where

$$\pi_2(a) := \frac{2}{(np)^\xi T^{1+\xi}} + \frac{8}{(np)^\xi T^\xi} + \frac{n^2}{a} \left( b_1 e^{-c_\phi \wedge a_2 (T/2)^{\gamma_2 \wedge \gamma_1}} + b_8 e^{-2\gamma_1 c_\phi (T/2)^{\gamma_1}} \right).$$

This bound controls the proximity between the empirical and population covariance matrices. Similar concentration inequalities were derived by [Kock and Callot \(2015, Lemma 9\)](#), [Loh and Wainwright \(2012, Lemma 14\)](#) and [Medeiros and Mendes \(2017\)](#). Their results, however, cannot be applied in our setting. Explicit expressions for the constants  $b_1$ ,  $b_8$  and  $\tau$  in Proposition 2 are found in Lemma 6. Also, one may replace  $\epsilon$  by its lower bound to remove dependence.

This concentration guided the choice of dependence condition used in this work. Traditionally one uses either a Hanson–Wright inequality or a Bernstein or Hoeffding type inequality to bound the empirical covariance around its mean. We write the centered Gram matrix  $\Gamma_T - \Gamma$  as a sum of martingales and a dependence term. The martingales are handled using a Bernstein type bound and the dependence term is handled using both assumptions (A1) and (A2). Combined, they imply a sub-Weibull type decay on expected value of dependence term.

Finally, we use the bounds  $\pi_1(\cdot)$  and  $\pi_2(\cdot)$  in Proposition 1 and Proposition 2 respectively, to derive an upper bound for the prediction error and for the difference between the lasso parameter estimates and the true parameters in the  $\ell_2$  norm.

**Theorem 1.** Suppose assumptions (A1) – (A3) hold. Under conditions of Proposition 1, there exists  $T_0 > 0$  such that for all  $T \geq T_0$ ,

$$\left| \mathbf{X} \left( \hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i^* \right) / T \right|_2^2 \leq 12 |\boldsymbol{\beta}_i^*|_1 \lambda, \quad i = 1, \dots, n,$$

with probability at least  $1 - \pi_1(\epsilon)$  for  $\epsilon > 0$ . Suppose further that assumptions (A4) and (A5) hold. Set  $\eta = \lambda / \sigma_\Gamma^2$ . Under conditions of Propositions 1 and 2, there exists  $T_0 > 0$  such that for all  $T \geq T_0$ ,

$$|\hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i^*|_2^2 \leq (44 + 2\lambda) R_q \left( \frac{\lambda}{\sigma_\Gamma^2} \right)^{2-q}, \quad i = 1, \dots, n,$$

with probability at least  $1 - \pi_1(\epsilon) - \pi_2 \left( \frac{\sigma_\Gamma^{2(1-q)} \lambda^q}{64 R_q} \right)$ .

Theorem 1 states that, with high probability, estimated and population parameter vectors are close to each other in the Euclidean norm. It requires that Propositions 1 and 2 hold jointly, meaning that  $\lambda$ ,  $R_q$  and  $\sigma_\Gamma^2$  must satisfy rate conditions. We show that if the size of ‘small’ coefficients and smallest eigenvalue  $\sigma_\Gamma^2$  of  $\Gamma$  are restricted, then the rate of  $\lambda$  in Proposition 1 is unaffected. For  $\epsilon = \log(np)$  and  $T < np^2$  Proposition 1 requires, after simplification,

$$\lambda \geq \tau^* \log(np)^{2/\alpha} \sqrt{\frac{4 \log(np)}{T}},$$

for some constant  $\tau^*$ . Replacing  $a$  by  $\frac{\sigma_\Gamma^{2(1-q)} \lambda^q}{64 R_q}$  in Proposition 2 we obtain

$$\lambda^q \gtrsim \left( \log(np)^{2/\alpha} \sqrt{\frac{\log(np)}{T}} \right) \times \left( \frac{R_q}{\sigma_\Gamma^{2(1-q)} \log(np)^{1/\alpha}} \right).$$

However, it is not necessarily a constraint in the rate of  $\lambda$ . Propositions 1 and 2 will hold jointly for  $T$  sufficiently large for  $0 \leq q < 1$  if

$$\frac{R_q}{\sigma_\Gamma^{2(1-q)}} = o \left( \log(np)^{(2q-1)/\alpha} \left( \frac{T}{\log(np)} \right)^{(1-q)/2} \right).$$

In other words, if the *small* parameters are not too large and smallest eigenvalue of  $\boldsymbol{\Sigma}$  is not too small as a function of  $T$ .

#### 4.1. Simulation

To evaluate the finite-sample performance of the LASSO estimator in large VARs with stochastic volatility, we consider the following model:

$$\begin{aligned} \mathbf{y}_t &= \mathbf{A}_0 + \mathbf{A}_1 \mathbf{y}_{t-1} + \mathbf{A}_4 \mathbf{y}_{t-4} + \mathbf{H}_t^{1/2} \mathbf{v}_t, & \mathbf{v}_t &\sim \mathbf{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{H}_{t+1} &= \mathbf{C}_0 + \boldsymbol{\Psi} \mathbf{H}_t \boldsymbol{\Psi}' + \boldsymbol{\epsilon}_t \boldsymbol{\epsilon}_t', & \boldsymbol{\epsilon}_t &\sim \mathbf{N}(\mathbf{0}, \mathbf{I}). \end{aligned} \quad (12)$$

We consider 1000 replications of model (12) with  $T = 100, 300$  observations and the number of series is a function of the sample size, that is,  $n = c \times T$ , where  $c = \{1, 2, 3\}$ .  $\mathbf{A}_1$  and  $\mathbf{A}_4$  are two block-diagonal matrices with

Table I. Simulation results

$T$	$n = T$	$n = 2T$	$n = 3T$
Panel (a): Mean squared estimation error			
100	0.1000	0.1926	0.2871
300	0.0288	0.0576	0.0914
Panel (b): Mean squared forecasting error			
100	1.0900	1.1456	1.2891
300	1.0150	1.0248	1.0345

*Note:* The table reports the simulation results of the LASSO estimation of model (12) for different combinations of sample size and number of variables. The penalty parameter of the LASSO is selected by the BIC. Panel (a) reports the mean squared error of the estimation of the VAR parameters. Panel (b) reports the ratio of the mean squared out-of-sample forecasting error of the LASSO with respect to the Oracle forecast.

blocks of dimension  $5 \times 5$  and entries equal to 0.15 and  $-0.1$  respectively.  $\mathbf{C}_0$  and  $\Psi$  are two diagonal matrices with diagonal elements all equal to  $1e-5$  and 0.8, respectively.

This model satisfies Assumptions (A1)–(A5). First, Assumption (A1) is satisfied because the model is block diagonal, with fixed block size. Assumptions (A2) and (A3) follows as in Example 4 given diagonal specification of  $\mathbf{C}_0$  and  $\Psi$ . Assumption (A4) is satisfied with  $q = 0$  and  $R = 10$ . Finally, assumption (A5) holds because lower bound in (6), is strictly positive ( $\sum_{i=1}^4 (\|A_i\|_1 + \|A_i\|_\infty) < \infty$  and  $\Lambda_{\min}(\Sigma) > 0$ , independently of  $n$  and  $p$ ).

Table I reports the simulation results of the LASSO estimation of model (12) for different combinations of sample size and number of variables. The penalty parameter of the LASSO is selected by the BIC. Panel (a) reports the mean squared error (MSE) of the estimation of the VAR parameters. Panel (b) reports the ratio of the mean squared out-of-sample forecasting error (MSFE) of the LASSO with respect to the Oracle forecast. A couple of findings emerge from the table. First, as expected the performance of the LASSO improves with the sample size. Second, as the number of candidate variables grows, the MSE increases. This is also expected. Finally, the out-of-sample forecasts of the LASSO get quite close to the ones from the Oracle when  $T = 300$ , even when  $n = 3T$ .

## 5. DISCUSSION

This work provides finite sample  $\ell_2$  error bounds for the equation-wise LASSO parameters estimates of a weakly sparse, high-dimensional, VAR( $p$ ) model, with dependent and heavy tailed innovation process. It covers a large collection of specifications as illustrated in Section 3.

A distinctive feature of this work is that the dependence structure of the innovations are characterized by a very weak projective dependence condition that is naturally verifiable in settings where one is interested in the conditional variance of the process. The series of innovations is not necessarily mixing nor the resulting time series  $\{\mathbf{y}_t\}$ .

Our bounds hold under a heavy tailed setting in a sense that we do not require the moment generating function to exist. Despite the tails in  $\{\mathbf{y}_t\}$  being sub-Weibull as in Wong *et al.* (2020), we are not able to recover the same rates and lower bound for the regularization parameter  $\lambda$ . The reason is that Wong *et al.* (2020) bounds rely heavily on the concentration inequality for mixing sequences in Merlevède *et al.* (2011). Given the weak projective dependence adopted, we chose to use a martingale concentration and overcome all together the issue of using the dependence metric for deriving the concentration bound. Nevertheless, we believe the loss in efficiency is minimal. Close inspection of proof of Lemma 5 shows that the loss of efficiency is concentrated in bounding the tail. It amounts to an extra  $\log(T)$  term, which does not change the rates under assumption that  $T < n^2p$ , eventually.

A limitation of this work is the restriction that the model is correctly specified in the mean, in a sense that innovations are martingale differences. Nevertheless, this assumption is standard in the literature and we are able

to derive results covering a broad range of data generating processes and conditional dependence measures. The martingale difference condition cannot be relaxed at this moment as our deviation bound depends on it. Furthermore, we do not require strong sparsity in a sense that near zero coefficients are effectively treated as zero as long as they are concentrated in some slowly increasing  $\ell_q$  ball ( $0 \leq q < 1$ ) around the origin.

Results in this article can be extended to polynomial tails. The strategy is to replace the martingale concentration in Lemma 4, used to prove Propositions 1 and 2 by

$$\Pr \left( \max_{1 \leq i \leq n} \left| \sum_{t=1}^T \xi_{it} \right| > Ta \right) \leq \frac{nK}{(a\sqrt{T})^d},$$

whenever  $\|\xi_{it}\|_d < \infty$ . If available under our dependence conditions, one could employ a Fuk–Nagaev type inequality. Nevertheless, it follows that under appropriate changes to concentration rates, equation-wise LASSO estimators also admit oracle bounds. A direct consequence is that moments conditions on Carrasco and Chen (2002) and Hafner and Preminger (2009a,b) are directly applicable.

Despite working with a relatively simple structure and estimation model, the machinery can be applied to more complex settings. The key points are showing that the empirical covariance concentrates around its mean in terms of its maximum entry-wise norm and the concentration inequality for large dimensional, sub-Weibull martingales. Following development of Negahban *et al.* (2012), the results may be naturally extended to structured regularization with node-wise regression and replacing using the Frobenius norm for system estimation. Finally, the high-dimensional VAR specification encompasses large dimensional vector-panels among other models.

In Adamek *et al.* (2020), authors consider a near epoch dependent time series. This condition covers misspecification and non-Gaussian, conditionally heteroskedastic models, such as GARCH innovations. The main focus of the authors is on inference using the desparsified lasso, but estimation error bounds are also derived. Assumptions are in the same line of Medeiros and Mendes (2016), where concentration bounds are assumed to hold in probability. In contrast, we focus on finite sample error bounds with relatable dependence condition on the conditional covariance. The heavy lifting in our article is to derive the required concentration inequalities. Effectively, one could use our results to provide theoretical justification for the concentration bounds under particular model specifications. On the other hand, Adamek *et al.* (2020) provides theoretical justification for desparsified inference in some large dimensional models discussed in our article.

#### ACKNOWLEDGEMENTS

The authors are grateful to Anders B. Kock, Giuseppe Cavaliere, Etienne JJ Wijler, editor (Robert Taylor) and co-editor of JTSA, and two anonymous referees. Their comments and suggestions led to a much improved version of this manuscript. M. C. Medeiros acknowledges partial support from CNPq/Brazil and CAPES.

#### DATA AVAILABILITY STATEMENT

The code that supports the simulation experiment in this study is available from the corresponding author on request.

#### REFERENCES

- Adamek R., Smeekes S., Wilms I. 2020. *Lasso inference for high-dimensional time series*. arXiv preprint arXiv:2007.10952.
- Andrews DW. 1988. Laws of large numbers for dependent non-identically distributed random variables. *Econometric Theory* **4**(3): 458–467.
- Basu S, Michailidis G. 2015. Regularized estimation in sparse high-dimensional time series models. *The Annals of Statistics* **43**: 1535–1567.
- Bauer G, Vornik K. 2011. Forecasting multivariate realized stock market volatility. *Journal of Econometrics* **160**: 93–101.

- Bradley RC. 2005. Basic properties of strong mixing conditions. A survey and some open questions. *Probability Surveys* **2**(2): 107–144.
- Carrasco M, Chen X. 2002. Mixing and moment properties of various GARCH and stochastic volatility models. *Econometric Theory* **18**: 17–39.
- Chiriac R, Voev V. 2011. Modelling and forecasting multivariate realized volatility. *Journal of Applied Econometrics* **26**: 922–947.
- Davidson J. 1994. *Stochastic Limit Theory*. Oxford: Oxford University Press.
- Davis R, Zang P, Zhang T. 2016. Sparse vector autoregressive modeling. *Journal of Computational and Graphical Statistics* **25**: 1077–1096.
- Dedecker J, Doukhan P, Lang G, León J, Louhichi S, Prieur C. 2007. *Weak Dependence with Examples and Applications*. Berlin: Springer.
- Ensr K, Raun L, Persse D. 2013. A case-crossover analysis of out-of-hospital cardiac arrest and air pollution. *Circulation* **127**: 1192–1199.
- Fan X, Grama I, Liu Q. 2012a. Hoeffding's inequality for supermartingales. *Stochastic Processes and their Applications* **122**(10): 3545–3559.
- Fan X, Grama I, Liu Q. 2012b. Large deviation exponential inequalities for supermartingales. *Electronic Communications in Probability* **17**.
- Hafner CM, Preminger A. 2009a. Asymptotic theory for a factor GARCH model. *Econometric Theory* **25**(02): 336–363.
- Hafner CM, Preminger A. 2009b. On asymptotic theory for multivariate GARCH models. *Journal of Multivariate Analysis* **100**(9): 2044–2054.
- Han Y, Tsay RS. 2020. High-dimensional linear regression for dependent data with applications to nowcasting. *Statistica Sinica* **30**: 1797–1827.
- Hoek G, Krishnan R, Beelen R, Peters A, Ostro B, Brunekreef B, Kaufman J. 2013. Long-term air pollution exposure and cardio-respiratory mortality: a review. *Environmental Health* **12**: 43.
- Kock A, Callot L. 2015. Oracle inequalities for high dimensional vector autoregressions. *Journal of Econometrics* **186**: 325–344.
- Kock A., Medeiros M., Vasconcelos G. 2020. Penalized regressions. In *Macroeconomic Forecasting in the Era of Big Data: Theory and Practice. Advanced Studies in Theoretical and Applied Econometrics*, Vol. 52. Berlin: Springer; 193–228.
- Loh P-L, Wainwright M. 2012. High-dimensional regression with noisy and missing data: provable guarantees with nonconvexity. *The Annals of Statistics* **40**: 1637–1664.
- Lütkepohl H. 1991. *Introduction to Multiple Time Series Analysis*. Berlin: Springer-Verlag.
- Lütkepohl H. 2006. *New Introduction to Multiple Time Series Analysis*. Berlin: Springer-Verlag.
- Medeiros M., Mendes E. 2016.  $\ell_1$  regularization of high-dimensional time-series models with non-Gaussian and heteroskedastic innovations. *Journal of Econometrics* **191**: 255–271.
- Medeiros M, Mendes E. 2017. Adaptive LASSO estimation for ARDL models with GARCH innovations. *Econometric Reviews* **622–637**(6–9).
- Merlevède F, Peligrad M, Rio E. 2011. A Bernstein type inequality and moderate deviations for weakly dependent sequences. *Probability Theory and Related Fields* **151**(3–4): 435–474.
- Miao K, Phillips PC, Su L. 2020. *High-dimensional VARs with common factors*, Cowles Foundation Discussion Papers 2252, Cowles Foundation for Research in Economics, Yale University.
- Miranda-Agrippino S, Ricco G. 2019. *Bayesian vector autoregressions: Estimation*, Oxford Research Encyclopedia of Economics and Finance.
- Negahban SN, Ravikumar P, Wainwright MJ, Yu B. 2012. A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers. *Statistical Science* **27**(4): 538–557.
- Ramey V. 2016. Macroeconomic shocks and their propagation. In *Handbook of Macroeconomics* Amsterdam: Elsevier; 71–162.
- Schweinberger M, Babkin S, Ensor K. 2017. High-dimensional multivariate time series with additional structure. *Journal of Computational and Graphical Statistics* **26**: 610–622.
- Sims C. 1980. Macroeconomics and reality. *Econometrica* **48**: 1–48.
- Song S., Bickel P. J. 2011. *Large vector auto regressions*, arXiv preprint arXiv:1106.3915.
- Van De Geer SA, Bühlmann P. 2009. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics* **3**: 1360–1392.
- Wilson G, Reale M, Haywood J. 2015. *Developments in Multivariate Time Series Modeling*. Boca Raton: CRC Press.
- Wong K, Li Z, Tewari A. 2020. Lasso guarantees for  $\beta$ -mixing heavy tailed time series. *Annals of Statistics* **48**(2): 1124–1142.

APPENDIX A. PROOF OF MAIN RESULTS

**A.1. Proof of Lemma 1**

We apply [Negahban et al. \(2012, Lemma 1\)](#). The empirical loss  $\mathcal{L}_T(\beta_i)$  is convex for each  $i$ . Proposition 1 ensures each (9) hold with desired probabilities. □

**A.2. Proof of Proposition 1**

Write the event  $\mathcal{A}_i = \{\max_j |\mathbf{u}'_i \mathbf{x}_j| < T\lambda_0/2\}$ . We shall derive probability bounds for  $\Pr(\cap_{i=1}^n \mathcal{A}_i) \geq 1 - \Pr(\max_{i,j} |\mathbf{u}'_i \mathbf{x}_j| \geq T\lambda_0/2)$ . We bound the probability using Corollary 1.

Under Assumption (A2),  $\mathbf{u}'_i \mathbf{x}_j = \sum_{t=p}^T u_{it} y_{t-s,j}$  ( $s = 1, \dots, p$  and  $i, j = 1, \dots, n$ ) is a martingale and each  $u_{it} y_{t-s,j}$  is a martingale difference process. Hence, we follow by applying Corollary 1. Conditions on  $T$  and  $\lambda_0/2$  are already satisfied. We need to show that  $u_{it} y_{t-s,j}$  is sub-Weibull. For each  $d \geq 1$ ,  $\|u_{it} y_{t-s,j}\|_d \leq \|u_{it}\|_{2p}^{1/2} \|y_{t-s,j}\|_{2p}^{1/2} \leq \bar{c}_\phi \max_{|b| \leq 1} \|\mathbf{b}' \mathbf{u}\|_{2p}$ , by Lemma 3. Then, it follows from [Wong et al. \(2020, Lemmas 5 and 6\)](#) and Assumption (A3) that  $u_{it} y_{t-s,j}$  is sub-Weibull with parameter  $\alpha/2$ . Hence, there is some constant  $\tau^*$  depending on  $\tau$ ,  $\bar{c}_\phi$  and  $\alpha$  such that  $\Pr(|u_{it} y_{t-s,j}| > x) \leq 2 \exp(-|x/\tau^*|^{\alpha/2})$ . Result follows.

**A.3. Proof of Lemma 2**

For notational simplicity, write  $\|\Gamma_T - \Gamma\|_{\max} \leq \delta \leq \sigma_\Gamma^2/64\psi^2(\mathcal{M}_{i,\eta})$  where  $\psi(\mathcal{M}_{i,\eta}) = \sup_{u \in C_i} |u|_1/|u|_2 = \sqrt{|S_{i,\eta}|}$ . Using the arguments in ([Negahban et al., 2012, section 4.3](#)),  $|\beta_{i,\mathcal{M}_{i,\eta}^\perp}|_1 \leq \sum_{j \in S_{i,\eta}^c} |\beta_{ij}|^q |\beta_{ij}|^{1-q} \leq \eta^{1-q} R_q$ , and  $R_q \geq \sum_{j \in S_{i,\eta}} |\beta_{ij}|^q \geq |S_{i,\eta}| \eta^q$ . Hence  $\frac{\sigma_\Gamma^2 \eta^q}{64R_q} \leq \frac{\sigma_\Gamma^2}{64\psi^2(\mathcal{M}_{i,\eta})}$ . It follows that

$$\begin{aligned} \Delta'_i \Gamma_T \Delta_i &= \Delta'_i \Gamma \Delta_i + \Delta'_i [\Gamma_T - \Gamma] \Delta_i \\ &\geq |\Delta_i|_2^2 \inf_{\mathbf{u} \in C_i \setminus \{0\}} \frac{\mathbf{u}' \Gamma \mathbf{u}}{\mathbf{u}' \mathbf{u}} - |\Delta_i|_1 \|[\Gamma_T - \Gamma] \Delta_i\|_\infty \\ &\geq \sigma_\Gamma^2 |\Delta_i|_2^2 - |\Delta_i|_1^2 \|\Gamma_T - \Gamma\|_{\max} \\ &\geq \sigma_\Gamma^2 |\Delta_i|_2^2 - \delta |\Delta_i|_1^2 \\ &\geq \sigma_\Gamma^2 |\Delta_i|_2^2 - \delta \left(4|\Delta_{i,\mathcal{M}_{i,\eta}}|_1 + 4|\beta_{i,\mathcal{M}_{i,\eta}^\perp}|_1\right)^2 \\ &\geq |\Delta_i|_2^2 \left(\sigma_\Gamma^2 - 32\delta\psi(\mathcal{M}_{i,\eta})^2\right) + 32\delta |\beta_{i,\mathcal{M}_{i,\eta}^\perp}|_1^2 \\ &\geq |\Delta_i|_2^2 \frac{\sigma_\Gamma^2}{2} - \frac{\sigma_\Gamma^2}{2\psi^2(\mathcal{M}_{i,\eta})} |\beta_{i,\mathcal{M}_{i,\eta}^\perp}|_1^2 \\ &\geq |\Delta_i|_2^2 \frac{\sigma_\Gamma^2}{2} - \frac{\sigma_\Gamma^2}{2R_q \eta^{-q}} \eta^{2(1-q)} R_q^2 \\ &= |\Delta_i|_2^2 \frac{\sigma_\Gamma^2}{2} - \frac{\sigma_\Gamma^2}{2} \eta^{2-q} R_q, \end{aligned}$$

proving the result. □

**A.4. Proof of Proposition 2**

The proof consists on a trivial application of Lemmas 6 setting  $\epsilon = \sigma_\Gamma^{2(1-q)} \lambda^q / 64R_q$ . □

**A.5. Proof of Theorem 1**

We apply Negahban *et al.* (2012, Theorem 1). Lemma 1 ensures  $\lambda$  is selected accordingly,  $\mathcal{L}_T(\beta_i)$  is a convex function of  $\beta_i$ , Lemma 2 ensures RSC is satisfied with  $\kappa_{\mathcal{L}} = \sigma_{\Gamma}^2/2$  and  $\tau_{\mathcal{L}}^2(\beta_i) = \frac{\sigma_{\Gamma}^2 \eta^{2-q} R_q}{2}$ . Define  $\psi(\mathcal{M}_{i,\eta})$  as in the proof of Lemma 2 and recall  $|S_{i,\eta}| \leq R_q \eta^{-q}$  and  $|\beta_{i,\mathcal{M}_{i,\eta}^{\perp}}|_1 \leq R_q \eta^{1-q}$ , and that  $\eta = \lambda/\sigma_{\Gamma}^2$ . For each  $i$ ,

$$\begin{aligned} |\hat{\beta}_i - \beta_i^*|_2^2 &\leq 9 \frac{\lambda}{\kappa_{\mathcal{L}}^2} \psi^2(\mathcal{M}_{i,\eta}) + \frac{\lambda}{\kappa_{\mathcal{L}}} \left[ 2\tau_{\mathcal{L}}^2(\beta_i^*) + 4|\beta_{i,\mathcal{M}_{i,\eta}^{\perp}}|_1 \right] \\ &\leq 36 \frac{\lambda}{\sigma_{\Gamma}^4} R_q \eta^{-q} + 2 \frac{\lambda}{\sigma_{\Gamma}^2} \left[ R_q \eta^{2-q} \sigma_{\Gamma}^2 + 4R_q \eta^{1-q} \right] \\ &\leq 36 \frac{\lambda}{\sigma_{\Gamma}^4} R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{-q} + 2 \frac{\lambda}{\sigma_{\Gamma}^2} \left[ R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{2-q} \sigma_{\Gamma}^2 + 4R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{1-q} \right] \\ &\leq 36 R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{2-q} + 2\lambda R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{2-q} + 8R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{2-q} \\ &= (44 + 2\lambda) R_q \left( \frac{\lambda}{\sigma_{\Gamma}^2} \right)^{2-q}. \end{aligned}$$

□

APPENDIX B. AUXILIARY LEMMATA

**B.1. Properties of  $y_t$**

We will derive properties of the process  $\{y_t\}$  described in (1)

**Lemma 3.** Suppose that for some norm  $\|\cdot\|_{\psi}$  we have

$$\max_t \max_{|b| \leq 1} \|b' u_t\|_{\psi} \leq c_{\psi},$$

for some constant  $c_{\psi} < \infty$  that only depends on the norm  $\|\cdot\|_{\psi}$ . Then, under conditions (A1) and (A2), for all  $t$  and  $i \in \{1, \dots, n\}$ ,

$$\|y_{i,t}\|_{\psi} \leq c_{\Phi} \times \sum_{j=0}^{\infty} |e_i' \Phi_j|_1.$$

*Proof.* Under assumption (A1) the VAR model in (1) admits the VMA( $\infty$ ) representation (4) for all  $n$  and  $p$ . Let  $\{e_i = (0, \dots, 0, 1, 0, \dots, 0)'\}$ ,  $i = 1, \dots, n$  the canonical basis vectors. Then, for all  $i$ ,  $y_{i,t} = e_i' y_t$  and

$$\begin{aligned} \|e_i' y_t\|_{\psi} &= \left\| \sum_{j=0}^{\infty} e_i' \Phi_j u_{t-j} \right\|_{\psi} \\ &= \left\| \sum_{j=0}^{\infty} \sum_{k=1}^n e_i' \Phi_j e_k u_{k,t-j} \right\|_{\psi} \end{aligned}$$



$$\begin{aligned} &= \left\| \sum_{j=0}^{\infty} |e'_i \Phi_j|_1 \sum_{k=1}^n \frac{e'_i \Phi_j e_k}{|e'_i \Phi_j|_*} u_{k,t-j} \right\|_{\psi} \\ &\leq \left( \sum_{j=0}^{\infty} |e'_i \Phi_j|_1 \right) \max_t \max_{|b| \leq 1} \|b' u_t\|_{\psi} \\ &\leq \sum_{j=0}^{\infty} |e'_i \Phi_j|_1 \times c_{\psi}, \end{aligned}$$

where  $|\cdot|_* := |\cdot|_1 I(|\cdot| > 0) + I(|\cdot| = 0)$ . □

Due stability condition (A1), for each  $n$  and  $p$ , there exists  $\bar{c}_{\Phi}$  such that  $\sum_{i=0}^{\infty} |\phi_{i,j}|_1 \leq \bar{c}_{\Phi}$  for all  $j = 1, \dots, n$ . Let  $\|\cdot\|_{\psi}$  be the Orlicz norm,

$$\|\cdot\|_{\psi} = \inf\{c > 0 : \psi(|\cdot|/c) \leq 1\},$$

where  $\psi(\cdot) : \mathbb{R}^+ \mapsto \mathbb{R}^+$  is convex, increasing function with  $\psi(0) = 0$  and  $\psi(x) \rightarrow \infty$  as  $x \rightarrow \infty$ . Traditional choices of  $\psi(\cdot)$  are (a)  $\psi(x) = x^p, p \geq 1$ , (b)  $\psi(x) = \exp(x^a) - 1, a > 1$ , and (c)  $\psi(x) = (ae)^{1/a} x I(x \leq a^{-1/a}) + \exp(x^a) I(x > a^{-1/a})$ . These choices contemplate sub-Gaussian and sub-exponential tails, as well as process with heavy-tails, such as sub-Weibull and polynomial tails. Note that by combining this result with Wong *et al.* (2020, Lemmas 5 and 6) if  $\{b' u_t\}$  are sub-Weibull, so are  $\{b' y_t\}$ .

Assumption (A1) is satisfied under restrictions on the parameter space. The stability assumption is standard in the literature whereas the tail sum (3) requires further constraints on the parameter matrices. Lemma 4 presents a sufficient set of restrictions on the sparse parameter matrices  $A_1, \dots, A_p$  so that (3) is satisfied.

**Lemma 4.** Suppose that for all  $n$  and  $p$ , there exists some  $\rho > 0$  such that

$$\sum_{k=1}^p \|A_k\|_{\infty} = \sum_{k=1}^p \max_{j=1, \dots, n} |a_{k,j}|_1 \leq e^{-\rho},$$

where  $A_k = [a_{k,1} : \dots : a_{k,n}]'$ . Then for every  $h = 1, \dots, n$ ,

- (i)  $|\phi_{k,h}|_1 \leq \sum_{j=1}^{p \wedge k} \|A_j\|_{\infty} |\phi_{k-j,h}|_1, k = 1, 2, \dots$
- (ii)  $\sum_{k=m}^{\infty} |\phi_{k,h}|_1 \leq c_0 e^{-mp}, m \geq 1$ , provided that for all  $p$ ,

$$\max_{h=1, \dots, n} \max_{k=1, \dots, p} e^{k\rho} \times \sum_{j=1}^k \tilde{\alpha}_j |\phi_{j,h}|_1 \leq (1 - e^{-\rho}) c_0, \tag{B1}$$

where  $\alpha_i = e^{\rho} \|A_i\|_{\infty}$  and  $\tilde{\alpha}_i = \sum_{i: |i|=k-p+j} \prod_{l=1}^{k-p} \alpha_{i_l}$  where  $i = (i_1, \dots, i_{k-p})$  is a multi-index.

*Proof.* Starting from the recursive definition of  $\Phi_k = \sum_{j=1}^{p \wedge k} \Phi_{k-j} A_j$ ,

$$|\phi_{k,h}|_1 = |e'_h \Phi_k|_1 = \left| \sum_{j=1}^{p \wedge k} e'_h \Phi_{k-j} A_j \right|_1 \leq \sum_{j=1}^{p \wedge k} |\phi_{k-j,h} A_j|_1 \leq \sum_{j=1}^{p \wedge k} |\phi_{k-j,h}|_1 \|A_j\|_{\infty}.$$

Suppose  $k \geq p$ , let  $\alpha_j = e^\rho \|A_j\|_\infty$  and verify that  $0 \leq \sum_{j=1}^p \alpha_j \leq 1$ . Iterating on the previous argument  $s \leq k - p$  times yields

$$\begin{aligned} |\phi_{k,h}|_1 &\leq \sum_{j_1=1}^p \cdots \sum_{j_s=1}^p \left( \prod_{l=1}^s |A_{j_l}|_\infty \right) |\phi_{k-\sum_{l=1}^s j_l,h}|_1 \\ &= e^{-s\rho} \sum_{j_1=1}^p \cdots \sum_{j_s=1}^p \left( \prod_{l=1}^s \alpha_{j_l} \right) |\phi_{k-\sum_{l=1}^s j_l,h}|_1 \\ &= \dots \\ &= e^{-\rho(k-p)} \sum_{j=1}^p \left( \sum_{i: |i|_1=k-p+j} \prod_{l=1}^{k-p} \alpha_{i_l} \right) |\phi_{p-j,h}|_1, \end{aligned}$$

where  $i = (i_1, \dots, i_{k-p})$  is a multi-index and the summation is over all combinations satisfying  $|i|_1 = k - p + j$ . The term inside parentheses is  $\tilde{\alpha}_j$  and under the conditions of the lemma

$$|\phi_{k,h}|_1 \leq e^{-\rho k} \times \left[ e^{\rho p} \sum_{j=1}^p \tilde{\alpha}_j |\phi_{p-j,h}|_1 \right] \leq (1 - e^{-\rho}) c_0 e^{-\rho k}.$$

The same result follows trivially for  $k < p$  under the assumptions of the lemma.

Summing over all values of  $k \geq m$ ,

$$\sum_{k=m}^\infty |\phi_{k,h}|_1 \leq c_0 (1 - e^{-\rho}) \sum_{k=m}^\infty e^{-\rho k} = c_0 e^{-m\rho} \frac{\sum_{k=0}^\infty e^{-\rho k}}{(1 - e^{-\rho})^{-1}} = c_0 e^{-m\rho}.$$

□

### B.2. Concentration Inequality for Martingales

We derive concentration bounds for martingales. In the first theorem we consider martingales with at most  $d$  finite moments, whereas in the second we allow the tails of the marginal distributions to decrease at a sub-Weibull, sub-exponential or, even sub- and super-Gaussian rate.

**Lemma 5** (Concentration bounds for high dimensional martingales). Let  $\{\xi_t\}_{t=1,\dots,T}$  denote a multivariate martingale difference process with respect to the filtration  $\mathcal{F}_t$  taking values on  $\mathbb{R}^n$  and assume  $\mathbb{E}(\xi_{it}^2)$  is finite for all  $1 \leq i \leq n$  and  $1 \leq t \leq T$ . Then,

$$\Pr \left( \left| \sum_{t=1}^T \xi_t \right|_\infty > Tx \right) \leq 2n \exp \left( -\frac{Tx^2}{2M^2 + xM} \right) + 4 \Pr \left( \max_{1 \leq i \leq n} |\xi_i|_\infty > M \right),$$

for all  $M > 0$ .

*Proof.* Write  $\xi_t = (\xi_{1t}, \dots, \xi_{nt})'$ . The proof follows after application of Fan *et al.* (2012a, Corollary 2.3).

Write  $V_k^2(M) = \max_{1 \leq i \leq n} \sum_{t=1}^k \mathbb{E} [\xi_{it}^2 I(\xi_{it} < M) | \mathcal{F}_t]$ ,  $X_{ik} = \sum_{t=1}^k \xi_{it}$  and  $X'_{ik}(M) = \sum_{t=1}^k \xi_{it} I(\xi_{it} \leq M)$ . It follows that for  $v > 0$  and  $x > 0$ ,

$$\begin{aligned} \Pr(|\mathbf{X}_T|_\infty > x) &\leq \Pr(\exists i, k : X_{ik} > x \cap V_k^2(M) \leq v^2) + \Pr(V_T^2(M) > v^2) \\ &\leq \Pr(\exists i, k : X'_{ik}(M) > x \cap V_k^2(M) \leq v^2) + \Pr(V_T^2(M) > v^2) \\ &\quad + \Pr\left(\max_{1 \leq i \leq n} \sum_{t=1}^k \xi_{it} I(\xi_{it} > M) > 0\right) \\ &\stackrel{(1)}{\leq} n \exp\left(-\frac{(Tx/M)^2}{2\left((v/M)^2 + \frac{T}{3}x/M\right)}\right) + \Pr(V_T^2(M) > v^2) \\ &\quad + \Pr\left(\max_{1 \leq t \leq T} |\xi_t|_\infty > M\right) \\ &\stackrel{(2)}{\leq} n \exp\left(-\frac{Tx^2}{2M^2 + Mx}\right) + 2 \Pr\left(\max_{1 \leq t \leq T} |\xi_t|_\infty > M\right). \end{aligned}$$

In (1) we use union bound and Fan *et al.* (2012a, Theorem 2.1) and in (2) we set  $v^2 = T\left(M^2 + \frac{1}{6T}Mx\right)$  and the following:

$$\begin{aligned} \Pr(V_T^2(M) > v^2) &\leq \Pr\left(\max_{1 \leq i \leq n} \sum_{t=1}^T \mathbb{E} [\xi_{it}^2 I(|\xi_{it}| \leq M) | \mathcal{F}_t] \geq v^2\right) \\ &\quad + \Pr\left(\max_{1 \leq i \leq n} \sum_{t=1}^T \mathbb{E} [\xi_{it}^2 I(\xi_{it} < -M) | \mathcal{F}_t] > 0\right) \\ &\leq \Pr\left(\max_{1 \leq i \leq n} \sum_{t=1}^T \mathbb{E} [\xi_{it}^2 I(|\xi_{it}| \leq M) | \mathcal{F}_t] \geq T\left(M^2 + \frac{1}{6T}Mx\right)\right) \\ &\quad + \Pr\left(\max_{1 \leq t \leq T} |\xi_t|_\infty > M\right) \\ &\leq \Pr\left(\max_{1 \leq t \leq T} |\xi_t|_\infty > M\right), \end{aligned}$$

where in the last line we note that  $\sum_{t=1}^T \mathbb{E} [\xi_{it}^2 I(|\xi_{it}| \leq M) | \mathcal{F}_t] \leq TM^2$ .

Finally, write  $\Pr(|\mathbf{X}_T|_\infty \geq Tx) = \Pr(\max_{i \leq n} X_{iT} \geq Tx) + \Pr(\max_{i \leq n} (-X_{iT}) \geq Tx)$  and apply above development in both terms. □

**Corollary 1.** Let  $\left\{ \xi_t = (\xi_{1t}, \dots, \xi_{nt})' \right\}_{t \geq 1}$  denote a multivariate martingale difference process with respect to the filtration  $\mathcal{F}_t$  taking values on  $\mathbb{R}^n$ . Suppose that for each  $\max_{i,t} \Pr(|\xi_{it}| > x) \leq 2e^{-(x/\tau)^\alpha}$ , for all  $x > 0$ , some  $\alpha > 0$  and  $\tau > 0$ . Then,

$$\Pr\left(\left|\sum_{t=1}^T \xi_t\right|_\infty > Tx\right) \leq 2n \exp\left(-\frac{Tx^2}{2M^2 + xM}\right) + 8nT \exp\left(-\frac{M^\alpha}{\tau^\alpha}\right).$$

In particular, if  $x > \tau (\epsilon + \log(nT))^{1/\alpha} \sqrt{\epsilon + \log n} / \sqrt{T}$  and  $T > (\epsilon + \log n)$  for any  $\epsilon > 0$ ,

$$\Pr \left( \left| \sum_{i=1}^T \xi_i \right|_{\infty} > Tx \right) \leq 10e^{-\epsilon}.$$

*Proof.* The first part we combine the union bound with assumption on  $\xi_{it}$ . In the second part, we will need the following bound. Let  $0 < a < b/4 < \infty$ . Then  $\sqrt{a+b} - \sqrt{a} \geq \sqrt{b} \left(1 - 2\sqrt{a/b}\right)^{1/2}$ . To verify that, first note that  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , then  $\left(\sqrt{a+b} - \sqrt{a}\right)^2 = 2a + b - 2\sqrt{a^2 + ab} \geq b - 2\sqrt{ab} = b \left(1 - 2\sqrt{a/b}\right)$ . Now, let  $a = 1/T$  and  $b = 8/(\epsilon + \log n)$  and verify that the choice  $M = x\sqrt{T} / \sqrt{\epsilon + \log n}$  satisfy  $\log n - \frac{Tx^2}{2M^2 + Mx} < -\epsilon$ , then replace  $M$  and  $x$  to obtain the bound.  $\square$

### B.3. Concentration Bound for Empirical Covariance Matrix

We derive concentration bound for  $\|\Gamma_T - \Gamma\|_{\max}$ , where  $\Gamma_T = \mathbf{X}'\mathbf{X}/T$  and  $\Gamma = \mathbb{E}\Gamma_T$ . We first split the problem into a sum of martingales and a tail dependence term. Then, we bound both individually.

**Lemma 6.** Suppose Assumptions (A1)–(A3) hold and

$$p < \frac{T^{\gamma_1 \wedge \gamma_2}}{\left(\frac{2}{\gamma_1 \wedge \gamma_2 + 1}\right) \left(2 + \frac{1.4}{2\gamma_1 c_\phi \wedge a_2}\right)}.$$

If for some  $\xi > 0$

$$\epsilon^2 \geq \frac{2(1 + \xi)^{1+2/\alpha} \tau^2 [\log(npT)]^{1+2/\alpha}}{T},$$

then

$$\Pr \left( \|\Gamma_T - \Gamma\|_{\max} \geq \epsilon \right) \leq \frac{2}{(np)^\xi T^{1+\xi}} + \frac{8}{(np)^\xi T^\xi} + \frac{n^2}{\epsilon} \left( b_1 e^{-c_\phi \wedge a_2 (T/2)^{\gamma_2 \wedge \gamma_1}} + b_8 e^{-2\gamma_1 c_\phi (T/2)^{\gamma_1}} \right) \tag{B2}$$

where  $b_1, b_5$  and  $\tau$  are constants not depending on  $T$ .

*Proof.* Use the union bound to rewrite our probability bound in terms of  $\mathbf{y}_{t-s}$ :

$$\Pr \left( \|\Gamma_T - \Gamma\|_{\max} > \epsilon \right) \leq 2 \sum_{r=0}^p \sum_{s=0}^{p-r} \Pr \left( \left\| \sum_{t=p+1}^T \mathbf{y}_{t-r} \mathbf{y}'_{t-r-s} - \mathbb{E} [\mathbf{y}_{t-r} \mathbf{y}'_{t-r-s}] \right\|_{\max} > T\epsilon \right).$$

Now, use a telescopic expansion of  $\mathbf{y}_t \mathbf{y}'_{t-s}$  to obtain a sum of martingales and a dependence term:

$$\sum_{t=p+1}^T \mathbf{y}_t \mathbf{y}'_{t-s} - \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}] = \underbrace{\sum_{t=p+1}^T \sum_{l=1}^m \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s} | \mathcal{F}_{t-l+1}] - \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s} | \mathcal{F}_{t-l}]}_{I_1}$$

$$\begin{aligned}
 & + \underbrace{\sum_{t=p+1}^T \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s} | \mathcal{F}_{t-m}] - \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}]}_{I_2} \\
 & = I_1 + I_2.
 \end{aligned}$$

Here,

$$I_1 = \sum_{l=1}^m \sum_{t=p+1}^T V_{l,t}^{(s)} \quad \text{and} \quad I_2 = \sum_{t=p+1}^T \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s} | \mathcal{F}_{t-m}] - \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}]$$

where  $\{V_{l,t}^{(s)}\}_t, l = 1, \dots, m$ , are sequences of martingale differences. The same decomposition holds for all terms  $\mathbf{y}_{t-s} \mathbf{y}'_{t-s}$ . Then,

$$\begin{aligned}
 \Pr (\|I_1 + I_2\|_{\max} > T\epsilon) & \leq \sum_{l=1}^m \Pr \left( \left\| \sum_{t=p+1}^T V_{l,t}^{(s)} \right\|_{\max} > \frac{T\epsilon}{2m} \right) \\
 & + \Pr \left( \left\| \sum_{t=p+1}^T \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s} | \mathcal{F}_{t-m}] - \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}] \right\|_{\max} > \frac{T\epsilon}{2} \right) \\
 & \leq 2mn^2 \exp \left( -\frac{T\epsilon^2}{2M^2 + M\epsilon} \right) \\
 & + 4mn^2 T \max_{l,t} \Pr (|V_{l,t}^{(s)}| > M) \tag{B3}
 \end{aligned}$$

$$+ \frac{2}{T\epsilon} \mathbb{E} \left| \max_{1 \leq i,j < n} \mathbb{E} \left[ \sum_{t=p+1}^T \mathbb{E} [e'_i \mathbf{y}_t \mathbf{y}'_{t-s} e_j | \mathcal{F}_{t-m}] - e'_i \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}] e_j \right] \right|, \tag{B4}$$

where  $e_i = (0, \dots, 0, 1, 0, \dots, 0)'$  is the  $i$ th canonical basis vector in  $\mathbb{R}^n$ .

**Bounding the tail (B3):**

The martingale differences  $\{V_{l,t}^{(s)}\} (l = 1, \dots, n \text{ and } s = 0, \dots, p)$  are sub-Weibull with parameter  $\alpha/2$ . For any random variables  $(X, Y)$  and  $\sigma$ -algebras  $\mathcal{F}$  and  $\mathcal{G}$ ,

$$\|\mathbb{E}[XY|\mathcal{F}] - \mathbb{E}[XY|\mathcal{G}]\|_p \leq 2\|XY\|_p \leq 2\|Y^2\|_p^{1/2} \|X^2\|_p^{1/2}.$$

Therefore, it follows from Wong *et al.* (2020, Lemmas 5 and 6) that if both  $X$  and  $Y$  are sub-Weibull with parameter  $\alpha$ , then  $XY$  is sub-Weibull with parameter  $\alpha/2$ . Therefore, there exists some  $\tau^*$  such that  $\Pr (|V_{l,t}^{(s)}| > s) \leq 2 \exp (-|x/\tau^*|^{\alpha/2})$ , bounding (B3).

**Bounding covariances (B4):**

Now we move toward bounding the dependence term (B4). Write

$$\mathbf{y}_t \mathbf{y}'_{t-s} = \sum_{j=0}^{s-1} \Phi_j \mathbf{u}_{t-j} \sum_{j=0}^{\infty} \mathbf{u}'_{t-s-j} \Phi'_j + \sum_{j=0}^{\infty} \Phi_{s+j} \mathbf{u}_{t-s-j} \sum_{j=0}^{\infty} \mathbf{u}'_{t-s-j} \Phi'_j$$

$$\begin{aligned}
 &= \sum_{j=0}^{s-1} \Phi_j \mathbf{u}_{t-j} \sum_{j=0}^{\infty} \mathbf{u}'_{t-s-j} \Phi'_j \\
 &\quad + \sum_{j=0}^{\infty} \Phi_{j+s} \mathbf{u}_{t-s-j} \mathbf{u}_{t-s-j} \Phi_j \\
 &\quad + \sum_{k=1}^{\infty} \sum_{j=0}^{\infty} \Phi_{j+s} \mathbf{u}_{t-s-j} \mathbf{u}_{t-s-j-k} \Phi_{j+k} \\
 &\quad + \sum_{k=1}^{\infty} \sum_{j=0}^{\infty} \Phi_{j+s+k} \mathbf{u}_{t-s-j-k} \mathbf{u}_{t-s-j} \Phi_j.
 \end{aligned}$$

It follows that  $\mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}] = \sum_{j=0}^{\infty} \Phi_{j+s} \Sigma \Phi'_j$ . Recall that  $\mathcal{F}_{t-m} = \sigma(\mathbf{u}_{t-i} : i = m, m + 1, \dots)$ , then, for  $m > s$ ,

$$\begin{aligned}
 \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s} | \mathcal{F}_{t-m}] - \mathbb{E} [\mathbf{y}_t \mathbf{y}'_{t-s}] &= \sum_{j=0}^{m-s-1} \Phi_j \mathbb{E} \left[ \mathbf{u}_{t-s-j} \mathbf{u}'_{t-s-j} - \Sigma | \mathcal{F}_{t-m} \right] \Phi'_{j+s} \\
 &\quad + \sum_{j=0}^{\infty} \Phi_{m+j} (\mathbf{u}_{t-m-j} \mathbf{u}_{t-m-j} - \Sigma) \Phi_{m-s+j} \\
 &\quad + \sum_{k=1}^{\infty} \sum_{j=0}^{\infty} \Phi_{m+j} \mathbf{u}_{t-m-j} \mathbf{u}'_{t-m-j-k} \Phi_{m-s+j+k} \\
 &\quad + \sum_{k=1}^{\infty} \sum_{j=0}^{\infty} \Phi_{m+j+k} \mathbf{u}_{t-m-j-k} \mathbf{u}'_{t-m-j} \Phi_{m-s+j} \\
 &= A_1(t, s, m) + A_2(t, s, m) + A_3(t, s, m) + A_4(t, s, m). \tag{B5}
 \end{aligned}$$

We shall bound  $\mathbb{E} \left| \sum_{r=0}^p \sum_{s=0}^{p-r} \mathbf{e}'_k A_i(t-r, s, m) \mathbf{e}_l \right|$  individually, for all  $\{\mathbf{e}_i, i = 1, \dots, n\}$  the canonical basis vector in  $\mathbb{R}^n$ .

**(a) Bounding  $\mathbb{E} \left| \sum_{r=0}^p \sum_{s=0}^{p-r} A_1(t-r, s, m) \right|$ :**

It follows from Assumption (A2) that for all  $\mathbf{b}_1, \mathbf{b}_2 \in \{\mathbf{b} \in \mathbb{R}^n : |\mathbf{b}|_1 = 1\}$

$$\max_t \mathbb{E} \left| \mathbb{E} \left[ \mathbf{b}'_1 (\mathbf{u}_t \mathbf{u}'_t - \Sigma)' \mathbf{b}_2 | \mathcal{F}_{t-m} \right] \right| \leq a_1 \exp(-a_2 m).$$

Set  $\{\mathbf{e}_i, i = 1, \dots, n\}$  the canonical basis vector in  $\mathbb{R}^n$ . It follows from Assumptions (A1) and (A2) that for  $j \leq m - s - 1$ :

$$\begin{aligned}
 &\mathbb{E} \left| \mathbf{e}'_k \Phi_j \mathbb{E} \left[ \left( \mathbf{u}_{t-s-j} \mathbf{u}'_{t-s-j} - \Sigma \right) | \mathcal{F}_{t-m} \right] \Phi'_{j+s} \mathbf{e}_l \right| \\
 &\leq |\phi_{j,k}| \left| |\phi_{j+s,l}| \max_t \mathbb{E} \left[ \mathbf{b}'_1 (\mathbf{u}_{t-s-j} \mathbf{u}'_{t-s-j} - \Sigma) \mathbf{b}_2 | \mathcal{F}_{t-m} \right] \right| \\
 &\leq \bar{c}_\Phi e^{-c_\Phi(j+s)\gamma} e^{-2c_\Phi j\gamma} \left[ a_1 e^{-a_2(m-j-s)\gamma} \right]
 \end{aligned}$$

Let  $0 < \gamma \leq 1$  and  $\frac{m^\gamma}{p} > \left( \frac{2}{\gamma+1} \right) \left( 2 + \frac{1.4}{c} \right)$  then  $c(m-p)^\gamma - \log(p+1) \geq c(m/2)^\gamma$ . Rewriting the inequality, we have to show that  $(m/p - 1)^\gamma - (m/2p)^\gamma > 2 \log(p+1)/cp^\gamma$ . In the LHS, a second order Taylor series expansion

yields  $(2a - 1)^\gamma - a^\gamma \geq \gamma \frac{a-1}{a} a^\gamma \left(1 - \frac{1-\gamma}{2} \frac{a-1}{a}\right) \geq \frac{\gamma+1}{2} \frac{a-1}{a} a^\gamma$ , for  $a > 1$ . Set  $a = m/2p$ , so that  $\frac{\gamma+1}{2} \left(\frac{m-2p}{m}\right) \left(\frac{m}{2p}\right)^\gamma \geq \log(p+1)/cp^\gamma$ . As for the RHS,  $\log(p+1)/p < 0.7$ . Combining bounds above we show our claim.

Note that for  $a, b \geq 0$  and  $0 < \gamma \leq 1$   $a^\gamma + b^\gamma = (a+b)^\gamma [x^\gamma + (1-x)^\gamma]$  where  $x = a/(a+b)$ , and  $1 \leq [x^\gamma + (1-x)^\gamma] \leq 2$  for  $x \in [0, 1]$ . Now, let  $0 < \gamma \leq 1$  and  $c > 0$ , then  $\sum_{j=0}^n e^{-c_j^\gamma} \leq \int_0^n e^{-cx^\gamma} dx = c^{1/\gamma} / \gamma \Gamma(1/\gamma, n) \uparrow c^{1/\gamma} \Gamma(1/\gamma + 1) < \infty$ , as  $n \rightarrow \infty$ , where  $\Gamma(a, n) = \int_0^n x^{a-1} e^{-x} dx \uparrow \int_0^\infty x^{a-1} e^{-x} dx = \Gamma(a)$  are the incomplete gamma function and gamma function respectively.

Then,

$$\begin{aligned} \mathbb{E} \left| \sum_{r=0}^p \sum_{s=0}^{p-r} e'_k A_1(t-r, s, m) e_l \right| &= \sum_{r=0}^p \sum_{s=0}^{p-r} \sum_{j=0}^{m-r-s-1} \mathbb{E} \left| e'_k \Phi_j \mathbb{E} \left[ \mathbf{u}_{t-r-s-j} \mathbf{u}'_{t-r-s-j} - \Sigma | \mathcal{F}_{t-m} \right] \Phi'_{j+s} e_l \right| \\ &\leq \bar{c}_\Phi a_1 \sum_{r=0}^p \sum_{s=0}^{p-r} \sum_{j=0}^{m-r-s-1} e^{-c_\Phi(j+s)^{\gamma_1}} e^{-2c_\Phi j^{\gamma_1}} e^{-a_2(m-r-j-s)^2} \\ &\leq \bar{c}_\Phi a_1 \sum_{r=0}^p e^{-(c_\Phi \wedge a_2)(m-r)^2 \wedge \gamma_1} \sum_{s=0}^{p-r} \sum_{j=0}^{m-r-s-1} e^{-c_\Phi j^{\gamma_1}} \\ &\leq \bar{c}_\Phi a_1 \frac{c_\Phi^{1/\gamma_1} \Gamma(1/\gamma_1)}{2\gamma_1} (p+1)^2 e^{-(c_\Phi \wedge a_2)(m-p)^2 \wedge \gamma_1} \\ &\leq b_1 e^{-(c_\Phi \wedge a_2)(m/2)^2 \wedge \gamma_1} \end{aligned}$$

**(b) Bounding**  $\mathbb{E} \left| \sum_{r=0}^p \sum_{s=0}^{p-r} e'_k A_2(t-r, s, m) e_l \right|$ :

Let  $\max_{\mathbf{b}_1, \mathbf{b}_2, t} \mathbb{E} \left| \mathbf{b}'_1 (\mathbf{u}_t \mathbf{u}'_t - \Sigma) \mathbf{b}_2 \right| \leq 2\Lambda_{\max}(\Sigma)$  where  $\mathbf{b}_1, \mathbf{b}_2 \in \{\mathbf{b} \in \mathbb{R}^n : |\mathbf{b}|_1 = 1\}$ . It follows from Lemma 8 after rearranging terms:

$$\begin{aligned} \mathbb{E} \left| \sum_{r=0}^p \sum_{s=0}^{p-r} e'_k A_2(t-r, s, m) e_l \right| &\leq \sum_{r=0}^p \sum_{s=0}^{p-r} \sum_{j=0}^\infty \mathbb{E} \left| e'_k \Phi_{m+j} \left( \mathbf{u}_{t-r-m-j} \mathbf{u}'_{t-r-m-j} - \Sigma \right) \Phi'_{m-s+j} e_l \right| \\ &\leq \sum_{r=0}^p \sum_{s=0}^{p-r} \sum_{j=m}^\infty |\phi_{j,k}|_1 |\phi_{j-s,l}|_1 \max_{\mathbf{b}_1, \mathbf{b}_2, t} \mathbb{E} \left| \mathbf{b}'_1 (\mathbf{u}_t \mathbf{u}'_t - \Sigma) \mathbf{b}_2 \right| \\ &\leq 2\Lambda_{\max}(\Sigma) \bar{c}_\Phi^2 \sum_{r=0}^p \sum_{s=0}^{p-r} \sum_{j=m}^\infty e^{-c_\Phi j^{\gamma_1}} e^{-c_\Phi(j-r)^{\gamma_1}} \\ &= 2b_2 \left( \sum_{r=0}^p \sum_{s=0}^{p-r} 1 \right) \sum_{j=m-p}^\infty e^{-2c_\Phi(j-r)^{\gamma_1}} \\ &= 2b_2 \left[ (p+1)(m-p)^{1/2} e^{-c_\Phi(m-p)^{\gamma_1}} \right]^2 \\ &\leq 2b_2 e^{-c_\Phi \left(\frac{m}{2}\right)^{\gamma_1}}, \end{aligned}$$

were  $b_2 = \Lambda_{\max}(\Sigma) \bar{c}_\Phi^2$ . In the last line, recall that  $(c_\Phi/2)(m-p)^{\gamma_1} - \log(p+1) \geq (c_\Phi/2)(m/2)^{\gamma_1}$ , then

$$\begin{aligned} (p+1)(m-p)^{1/2} e^{-c_\Phi(m-p)^{\gamma_1}} &= (p+1) e^{\frac{1}{2} \log(m-p) - c_\Phi(m-p)^{\gamma_1}} \\ &\leq e^{\log(p+1) - \frac{c_\Phi}{2} \left(\frac{m}{2}\right)^{\gamma_1}} \leq e^{-\frac{c_\Phi}{2} \left(\frac{m}{2}\right)^{\gamma_1}}. \end{aligned}$$

(c) **Bounding**  $\sum_{r=0}^p \sum_{s=0}^{p-r} A_j(t-r, s, m)$  ( $j = 3, 4$ ):

Under Assumption (A1)-(A3), for all  $\mathbf{b} \in \mathbb{R}^n$  with  $\|\mathbf{b}\|_1 = 1$ ,

$$\begin{aligned} \mathbb{E}|e'_r \Phi_{m+j} \mathbf{u}_{t-m-j} \mathbf{u}'_{t-m-j-k} \Phi_{m-s+j+k} e_s| &\leq |\phi_{m+j,r}|_1 |\phi_{m-s+j+k,s}|_1 \max_t \|\mathbf{b}' \mathbf{u}_t\|_2^2 \\ &\leq b_2 e^{-c_\phi(m+j)^{\gamma_1} - c_\phi(m-s+j+k)^{\gamma_1}}, \end{aligned}$$

where  $b_2 = \bar{c}_\phi^2 \Lambda_{\max}(\Sigma)$ . As before, if we have  $\mathbf{u}_{t-r-s-j}$  we must replace  $m$  by  $m-r$ . It follows from Lemma 8,  $\frac{m^{\gamma_1}}{p} > \left(\frac{2}{\gamma_1+1}\right) \left(2 + \frac{1.4}{2\gamma_1 c_\phi}\right)$  and  $m$  sufficiently large:

$$\begin{aligned} &\sum_{r=0}^p \sum_{j=0}^{\infty} e^{-c_\phi(m-r+j)^{\gamma_1}} \sum_{s=0}^{p-r} \sum_{k=0}^{\infty} e^{-c_\phi(m-r+j+k-s)^{\gamma_1}} \\ &\leq \left(\sum_{r=0}^p \sum_{s=0}^{p-r} 1\right) \sum_{j=0}^{\infty} e^{-c_\phi(m-p+j)^{\gamma_1}} \sum_{k=0}^{\infty} e^{-c_\phi(m-p+j+k)^{\gamma_1}} \\ &\leq \frac{1}{2} (p+1)^2 \sum_{j=m-p}^{\infty} \sum_{k=0}^{\infty} e^{-c_\phi j^{\gamma_1} - c_\phi(j+k)^{\gamma_1}} \\ &\leq b_3 [(p+1)(m-p)e^{-c_\phi(m-p)^{\gamma_1}}]^2 \\ &\leq b_3 \left[e^{\log(p+1) - \frac{c_\phi}{2}(m-p)^{\gamma_1}}\right]^2 \\ &\leq b_3 e^{-c_\phi \left(\frac{m}{2}\right)^{\gamma_1}}, \end{aligned}$$

where  $b_3 = \left(1 + \frac{2}{\gamma_1}\right) / 2$ .

Then, it follows that

$$\sum_{r=0}^p \sum_{s=0}^{p-r} \sum_{k=1}^{\infty} \sum_{j=0}^{\infty} \mathbb{E}|e'_t \Phi_{m+j} \mathbf{u}_{t-m-j} \mathbf{u}'_{t-m-j-k} \Phi_{m-s+j+k} e_t| \leq b_3 e^{-c_\phi \left(\frac{m}{2}\right)^{\gamma_1}}, \tag{B6}$$

where  $b_4 = b_2 b_3 = \bar{c}_\phi^2 \Lambda_{\max}(\Sigma) \left(1 + \frac{2}{\gamma_1}\right) / 2$ .

(d) **Combining bounds:**

Finally combining the three bounds above and setting  $m$  satisfying  $\frac{m^{\gamma_1 \wedge \gamma_2}}{p} > \left(\frac{2}{\gamma_1 \wedge \gamma_2 + 1}\right) \left(2 + \frac{1.4}{2\gamma_1 c_\phi \wedge a_2}\right)$ :

$$\begin{aligned} &\sum_{r=0}^p \sum_{s=0}^{p-r} \mathbb{E} \left| \max_{1 \leq i, j < n} \frac{1}{T} \mathbb{E} \left| \sum_{t=p+1}^T \mathbb{E} [e'_i \mathbf{y}_{t-r} \mathbf{y}'_{t-r-s} e_j | \mathcal{F}_{t-m}] - \mathbb{E} [e'_i \mathbf{y}_{t-r} \mathbf{y}'_{t-r-s} e_j] \right| \right| \\ &\leq n^2 \left( b_1 e^{-(c_\phi \wedge a_2)(m/2)^{\gamma_2 \wedge \gamma_1}} + b_5 e^{-c_\phi(m/2)^{\gamma_1}} \right), \end{aligned} \tag{B7}$$

where  $b_5 = 2b_2 + b_4 = \bar{c}_\phi^2 \Lambda_{\max}(\Sigma) \left[2 + \left(1 + \frac{2}{\gamma_1}\right) / 2\right]$



**Combining tail and covariance:**

Set  $m = T$ , then we require that

$$p < \frac{T^{\gamma_1 \wedge \gamma_2}}{\left(\frac{2}{\gamma_1 \wedge \gamma_2 + 1}\right) \left(2 + \frac{1.4}{2\gamma_1 c_\phi \wedge a_2}\right)}.$$

Then, combining bounds:

$$\begin{aligned} \Pr(\|\Gamma_T - \Gamma\|_{\max} \geq \epsilon) &\leq 2\frac{1}{T} \exp\left(2 \log(npT) - \frac{T\epsilon^2}{2M^2 + M\epsilon}\right) \\ &\quad + 8e^{2\log(npT) - M^\alpha/\tau^\alpha} \\ &\quad + \frac{n^2}{\epsilon} \left(b_1 e^{-a_2 \wedge c_\phi (T/2)^{\gamma_2 \wedge \gamma_1}} + b_5 e^{-c_\phi (T/2)^{\gamma_1}}\right), \end{aligned}$$

where  $\tau$ ,  $b_1$  and  $b_5$  are as above. Let  $\xi > 0$  and  $M^\alpha = (2 + \xi)\tau^\alpha \log(npT)$ , by assumption  $T\epsilon^2 \geq 2(1 + \xi)^{1+2/\alpha} \tau^2 [\log(npT)]^{1+2/\alpha}$  which implies that  $2 \log(npT) - \frac{T\epsilon^2}{2M^2 + M\epsilon} \leq -\xi \log(npT)$ . Finally,

$$\Pr(\|\Gamma_T - \Gamma\|_{\max} \geq \epsilon) \leq \frac{2}{(np)^\xi T^{1+\xi}} + \frac{8}{(np)^\xi T^\xi} + \frac{n^2}{\epsilon} \left(b_1 e^{-a_2 \wedge c_\phi (T/2)^{\gamma_2}} + b_5 e^{-2\gamma_1 c_\phi (T/2)^{\gamma_1}}\right).$$

□

**Lemma 7.** Let  $X$  denote a positive random variable with  $P(X \geq u) \leq e^{-cu^\alpha}$  for positive constants  $c > 0$  and  $\alpha > 0$ . Then, for  $M > 0$  and  $p \geq 1$ ,

$$M^p P(X \geq M) \leq \mathbb{E}[X^p I(X \geq M)] \leq \left(1 + \frac{p}{\alpha}\right) M^p e^{-cM^\alpha}. \tag{B8}$$

In particular, let  $P(X \geq u) = e^{-cu^\alpha}$ . Then,

$$M^p e^{-cM^\alpha} \leq \mathbb{E}[X^p I(X \geq M)] \leq \left(1 + \frac{p}{\alpha}\right) M^p e^{-cM^\alpha}.$$

*Proof.* The lower bound follows trivially:  $\mathbb{E}[X^p I(X \geq M)] \geq M^p P(X \geq M)$ . Now, verify that  $P(XI(X > M) \geq u) = I(u < M^p)P(X \geq M) + I(u \geq M^p)P(X^p \geq u)$ . We have

$$\begin{aligned} \mathbb{E}[X^p I(X \geq M)] &= \int_0^\infty P(XI(X > M) \geq u) du \\ &= \int_0^{M^p} du P(X \geq M) + \int_{M^p}^\infty P(X^p \geq u) du \\ &\leq M^p e^{-cM^\alpha} + \int_{M^p}^\infty e^{-cu^\alpha/p} du \\ &= M^p e^{-cM^\alpha} + \frac{p}{\alpha} c^{-p/\alpha} \Gamma(p/\alpha, cM^\alpha) \\ &\leq \left(1 + \frac{p}{\alpha}\right) M^p e^{-cM^\alpha}, \end{aligned}$$

where  $\Gamma(s, x) = \int_x^\infty t^{s-1} e^{-t} dt$  is the upper incomplete gamma function. □

**Lemma 8.** Let  $0 < a \leq 1$ ,  $b > 0$ ,  $n \geq 1$ . Then

$$\sum_{i=n}^{\infty} \sum_{j=0}^{\infty} e^{-bi^a - b(i+j)^a} \leq \left(2 + \frac{1}{a}\right) n^2 e^{-2bn^a}, \quad (\text{B9})$$

and

$$\sum_{i=n}^{\infty} e^{-bi^a} \leq 2ne^{-bn^a}. \quad (\text{B10})$$

*Proof.* Let  $V$  denote a Weibull( $a, 2b$ ) random variable.

$$\begin{aligned} \sum_{i=n}^{\infty} \sum_{j=0}^{\infty} e^{-bi^a - b(i+j)^a} &= \sum_{j=n}^{\infty} (j - n + 1) e^{-2bj^a} \\ &= \sum_{j=n}^{\infty} (j - n + 1) \mathbb{E}[I(V \geq j)] \\ &= \sum_{j=n}^{\infty} (j - n + 1)^2 \mathbb{E}[I(j \leq V < j + 1)] \\ &\leq \mathbb{E}[(V - n + 1)^2 I(V \geq n)]. \end{aligned}$$

Expanding squares and using Lemma 7,

$$\begin{aligned} \mathbb{E}[(V - n + 1)^2 I(V \geq n)] &= \mathbb{E}[V^2 I(V \geq n)] - 2(n - 1) \mathbb{E}[VI(V \geq n)] + (n - 1)^2 P(V \geq n) \\ &\leq (1 + 2/a)n^2 e^{-2bn^a} - 2(n - 1)ne^{-2bn^a} + (n - 1)e^{-2bn^a} \\ &= (2/a)e^{-2bn^a} n^2 + e^{-2bn^a} [n - (n - 1)]^2 \\ &\leq \left(1 + \frac{2}{a}\right) n^2 e^{-2bn^a}. \end{aligned}$$

Similarly, let  $V$  be a Weibull( $a, b$ ) random variable,

$$\begin{aligned} \sum_{i=n}^{\infty} e^{-bi^a} &\leq \mathbb{E}[(V - n + 1)I(V > n)] \\ &\leq ne^{-bn^a} + (n - 1)e^{-bn^a} \\ &< 2ne^{-bn^a}. \end{aligned}$$

□