



Published in final edited form as:

J Am Stat Assoc. 2019 ; 114(526): 735–748. doi:10.1080/01621459.2018.1437043.

High-Dimensional Posterior Consistency in Bayesian Vector Autoregressive Models

Satyajit Ghosh, Kshitij Khare, George Michailidis

Department of Statistics and the Informatics Institute, University of Florida

Satyajit Ghosh: satyajitghosh90@ufl.edu; Kshitij Khare: kdkhare@ufl.edu; George Michailidis: gmichail@ufl.edu

Abstract

Vector autoregressive (VAR) models aim to capture linear temporal interdependencies amongst multiple time series. They have been widely used in macroeconomics and financial econometrics and more recently have found novel applications in functional genomics and neuroscience. These applications have also accentuated the need to investigate the behavior of the VAR model in a high-dimensional regime, which provides novel insights into the role of temporal dependence for regularized estimates of the model's parameters. However, hardly anything is known regarding properties of the posterior distribution for Bayesian VAR models in such regimes. In this work, we consider a VAR model with two prior choices for the autoregressive coefficient matrix: a non-hierarchical matrix-normal prior and a hierarchical prior, which corresponds to an *arbitrary* scale mixture of normals. We establish posterior consistency for both these priors under standard regularity assumptions, when the dimension p of the VAR model grows with the sample size n (but still remains smaller than n). A special case corresponds to a shrinkage prior that introduces (group) sparsity in the columns of the model coefficient matrices. The performance of the model estimates are illustrated on synthetic and real macroeconomic data sets.

Keywords

Vector Autoregressive Models; Posterior consistency; Shrinkage prior; Bayesian Lasso

1 Introduction

There has been recent interest in modeling high-dimensional time series data sets. In macroeconomics, Mol et al. (2008) advocated the need to include a large number of variables in econometric models to improve forecastability, while Billio et al. (2012) examined stock returns of many financial institutions to assess systemic risk of the financial system. Similar modeling challenges arise in functional genomics for the reconstruction of regulatory networks as discussed in Basu et al. (2015), while in neuroscience one is interested in understanding functional connectivity between brain regions Seth et al. (2015).

A popular and informative model has been vector autoregressions (VAR), that captures linear temporal dependencies between time series. The VAR model and its properties have been thoroughly explored in low-dimensional settings both from a frequentist (for a comprehensive overview see Lütkepohl (2007)) and a Bayesian perspective Ba bura et al. (2010).

More recently, Basu and Michailidis (2015) provided an in depth analysis of the model for Gaussian data in a *high-dimensional setting under sparsity assumptions*, while Melnyk and Banerjee (2016) extended the results to other regularizers (e.g. group lasso, sparse group lasso, etc.). The results of Basu and Michailidis (2015); Melnyk and Banerjee (2016) and related follow-up work (Raskutti and Yuan (2015); Schweinberger et al. (2017); Lin and Michailidis (2017)) indicate that the resulting estimation error rates are those obtained for independent and identically distributed data times a factor that captures the temporal dependence in the data.

On the Bayesian front, there has been primarily methodological/computational work for low-dimensional VAR models. The so-called Minnesota prior Litterman et al. (1979); Doan et al. (1984) has been a staple of applied econometric work involving VAR models. This is a normal prior distribution on the elements of the transition matrix that puts stronger weights on the “own” lags of each time series, since they are considered more informative for forecasting purposes than lags from “other” time series. For large size VAR models, Ba bura et al. (2010) advocate normal-inverted Wishart distribution that leads to a posterior mean that can be interpreted as a ridge shrinkage estimator, suitable for such models. A first attempt for Bayesian estimation of VAR models combined with variable selection is presented in Korobilis (2013), where an indicator variable is specified for each parameter in the transition matrix that indicates whether the cross-autocorrelation coefficient is included or set to zero. A prior is specified for the indicator variables that in principle can also be combined with the Minnesota prior.

On the other hand, Bayesian investigations into high-dimensional asymptotics of statistical models that *incorporate sparsity with temporally dependent data* are not in general available to the best of our knowledge. Hence, the main objective of this work is to study posterior (estimation) consistency for a VAR model, which asserts that the posterior concentrates around the “true” parameter value (in an appropriate norm) as the sample size increases. There is a rich literature on high-dimensional posterior estimation consistency for linear regression models for independent and identically distributed data. Ghoshal (1999) established posterior consistency and asymptotic normality with a general prior on the p -vector of regression coefficients (with appropriate positivity and Lipschitz assumptions) when $p^3 \log p/n \rightarrow 0$ and $p^4 \log p/n \rightarrow 0$, respectively. Bontemps (2011) extended the work of Ghoshal (1999) by permitting the model to be misspecified and the number of predictor variables to grow proportionally to the sample size. Armagan et al. (2013) focus on shrinkage priors, which are appropriate scale mixtures of normal priors and induce weak sparsity in the vector of regression coefficients (see Carvalho et al. (2010); Griffin and Brown (2010); Armagan et al. (2011, 2013)). They establish posterior consistency under a simple sufficient condition on prior concentration when $p = \alpha(n)$. Lee and Oh (2013) establish posterior consistency under a high dimensional Bayesian PCA regression setup with $p > n$ under appropriate assumptions on the rank of the design matrix. Posterior estimation consistency in linear regression models with g -priors has also been addressed in Sparks et al. (2015).

A crucial difference between the linear regression models considered in the above work and VAR models (expressed as a linear model) is that the design matrix in the latter case is

random, and exhibits dependencies both between its rows and across its columns, and also with the error term in the model (see Section 2). This leads to a significantly more involved and challenging theoretical analysis that we successfully resolve. In this paper, we investigate high-dimensional posterior consistency for Bayesian VAR models in two natural and relevant settings: (a) with a non-hierarchical matrix normal prior on the $dp \times p$ autoregressive parameter matrix and (b) a hierarchical prior which corresponds to a general scale mixture of normals. In particular, this includes spherically symmetric priors such as the multivariate- t and standard shrinkage priors which induce (group) sparsity in the columns of the coefficient matrices, such as the group structure in Basu et al. (2015). Further, we employ a flat (uniform) prior distribution for the error term. Note that the joint maximum likelihood estimation problem for a sparse VAR model, with a sparse error covariance matrix is investigated in Lin and Michailidis (2017). The posterior consistency results are established under mild regularity assumptions on the underlying spectral density and with $p = \alpha(n/\log n)$. The key to handling the dependencies, within the design matrix and also between the design matrix and the error term, is a pair of high-dimensional concentration inequalities established in the supplementary material (Propositions B1 and B3). Note that we are considering the “large p large n ” setting with $p = \alpha(n)$. However, we make no assumption reducing the effective dimension of the “true parameter matrix”. We only assume that the matrix norm of the true parameter matrix is of the order p in the non-hierarchical prior setting and bounded by a constant in the hierarchical prior setting. The large p small n situation, where p is allowed to grow at a much faster rate than n is also of interest, but assumptions such as sparsity/restricted eigenvalue type conditions are required, which in turn reduce the effective dimension of the true parameters. General posterior consistency results for VAR models in the large p small n setting are also not available to the best of our knowledge and are topics of future discussion/research.

The remainder of the paper is organized as follows. In Section 2, we introduce the VAR model and necessary notions of posterior consistency. We consider the non-hierarchical matrix normal prior on the coefficient matrix in Section 3.1 and establish posterior consistency under suitable regularity assumptions. In Section 3.2, we prove posterior consistency considering a hierarchical prior corresponding to a scale-mixture of matrix normals. In Section 4 and Section 5 the methodology/results of this paper are illustrated on simulated and real data sets, respectively. Finally, we conclude with a discussion in Section 6.

1.1 Notation

Throughout this paper, \mathbb{Z} , \mathbb{R} and \mathbb{C} denote the sets of integers, real numbers and complex numbers, respectively. We denote the cardinality of a set J by $|J|$. For a vector $v \in \mathbb{R}^p$,

$\|v\| := \sqrt{\sum v_j^2}$ denotes the ℓ_2 -norm. For a matrix \mathbf{A} , $\|\mathbf{A}\|$ and $\sigma_{\max}(\mathbf{A})$ denote spectral norm i.e.,

$\|\mathbf{A}\| = \sup_{x \neq 0} \frac{\|\mathbf{A}x\|_2}{\|x\|_2}$ and the largest singular value of \mathbf{A} , respectively. For a symmetric or

Hermitian matrix \mathbf{A} , we denote its maximum and minimum eigenvalues by $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$. The vector e_i is used for the i -th unit vector in \mathbb{R}^p . Bold uppercase letters are only used to denote matrices, and vectorized form of such matrices are represented by

corresponding lower cases. For example, if Φ is a $p \times p$ matrix then ϕ is $\text{vec}(\Phi)$. Also, $\mathbf{0}$ represents a zero-matrix of appropriate dimension, and in general vectors are denoted by italicized bold lowercase letters.

2 Model Formulation

For a p -dimensional stationary time series $\{X^t\}$, a vector autoregressive model of lag- d is given by

$$X^t = \mathbf{c} + \sum_{i=1}^d \mathbf{A}_i X^{t-i} + \varepsilon^t. \quad (1)$$

The temporal dependence structure of the VAR model is characterized by the $p \times p$ transition matrices $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d$ and \mathbf{c} is a $p \times 1$ location vector which we choose to be $\mathbf{0}$. In the Gaussian VAR, the errors ε^t are i.i.d $\mathcal{N}_p(\mathbf{0}, \Sigma_\varepsilon)$ where Σ_ε is a $p \times p$ unknown error covariance matrix. The model in (1) can be rewritten in the Yule-Walker representation (Lütkepohl (2007)) as

$$X^t - \mu = \sum_{i=1}^d \mathbf{A}_i (X^{t-i} - \mu) + \varepsilon^t,$$

where $\mu = (\mathbf{I} - \mathbf{A}_1 - \mathbf{A}_2 - \dots - \mathbf{A}_d)^{-1} \mathbf{c}$ is known as the process mean. Usually μ will not be known in advance. In that case, μ may be estimated by the vector of sample means $\bar{X} = \sum_1^n X^t$. An alternative estimator is $\hat{\mu} = (\mathbf{I} - \hat{\mathbf{A}}_1 - \hat{\mathbf{A}}_2 - \dots - \hat{\mathbf{A}}_d)^{-1} \hat{\mathbf{c}}$ in which $\hat{\mathbf{c}}$ and $\hat{\mathbf{A}}_i$'s are the least squares estimator. Henceforth we assume without loss of generality $\mu = \mathbf{0}$. Based on the data $\{X^0, \dots, X^T\}$, we define the response matrix \mathbf{Y} and design matrix \mathbf{X} as follows,

$$\mathbf{Y} = \begin{bmatrix} (X^T)' \\ \vdots \\ (X^d)' \end{bmatrix}_{n \times p} \quad \mathbf{X} = \begin{bmatrix} (X^{T-1})' & \dots & (X^{T-d})' \\ \vdots & \ddots & \vdots \\ (X^{d-1})' & \dots & (X^0)' \end{bmatrix}_{n \times dp}.$$

We can now rewrite the above model in a linear regression setup as

$$\mathbf{Y} = \mathbf{X}\Phi + \mathbf{E} \quad (2)$$

where

$$\Phi = \begin{bmatrix} \mathbf{A}'_1 \\ \mathbf{A}'_2 \\ \vdots \\ \mathbf{A}'_d \end{bmatrix} \quad \mathbf{E} = \begin{bmatrix} (\varepsilon^T)' \\ \vdots \\ (\varepsilon^d)' \end{bmatrix}.$$

In this formulation, the number of samples is $n = T - d + 1$ and the number of unknown parameters is $q = dp^2$, respectively. Vectorizing (column-wise) each matrix, we get

$$\mathbf{y} := \text{vec}(\mathbf{Y}) = \mathbf{Z}\boldsymbol{\phi} + \boldsymbol{\varepsilon},$$

where $\mathbf{Z} := (\mathbf{I}_p \otimes \mathbf{X})$, $\boldsymbol{\phi} = \text{vec}(\boldsymbol{\Phi})$ and $\boldsymbol{\varepsilon} = \text{vec}(\mathbf{E})$. In this paper, we consider a high-dimensional setting where the dimension p of the VAR model increases with the sample size n . However, we assume that the lag d does not vary with n . This basic formulation of regression lends itself easily to a Bayesian analysis in which priors are placed on the unknown parameter matrices $\boldsymbol{\Phi}$ and $\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}}$.

As previously mentioned, we let the dimension $p = p_n$ of the VAR model vary with n , so that our results are relevant to high dimensional settings. We assume that our data come from a true VAR model described as follows: for every, $n \geq 1$, let $\mathcal{Y}_n := (X^{n,0}, \dots, X^{n,n+d-1})$ be the set of observations for sample size n , which satisfy $X^{n,k} = \sum_{i=1}^d \mathbf{A}_{i,0n} X^{n,k-i} + \varepsilon^{n,k}$ for $d \leq k \leq n + d - 1$. The errors $\{\varepsilon^{n,k}\}_{k=d}^{n+d-1}$ are i.i.d. $\mathcal{N}_{p_n}(0, \boldsymbol{\Sigma}_{\boldsymbol{\varepsilon},0n})$. Here $\{\boldsymbol{\Phi}_{0n}\}_{n \geq 1}$ denotes the sequence of the true coefficient matrices given by $\boldsymbol{\Phi}'_{0n} := [\mathbf{A}_{1,0n} \ \mathbf{A}_{2,0n} \ \dots \ \mathbf{A}_{d,0n}]$, and $\{\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon},0n}\}_{n \geq 1}$ denotes the sequence of the true error covariance matrices. Let \mathbb{P}_0 denote the probability measure underlying the true model described above.

Next, consider a Bayesian model which builds on (2) by placing priors on the parameters $(\boldsymbol{\Phi}, \boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}})$. In particular, let $\{\pi_n(\boldsymbol{\Phi}, \boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}})\}_{n \geq 1}$ and $\{\pi_n(\boldsymbol{\Phi}, \boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}} | \mathcal{Y}_n)\}_{n \geq 1}$ denote the sequences of the corresponding (joint) prior and posterior densities. Analogously, $\{\Pi_n(\cdot)\}_{n \geq 1}$ and $\{\Pi_n(\cdot | \mathcal{Y}_n)\}_{n \geq 1}$ denote the corresponding sequences of (joint) prior and posterior distributions. We will also use the notation $\boldsymbol{\pi}_n$ and Π_n to denote the marginal prior and posterior densities/distributions for $\boldsymbol{\Phi}$ and $\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}}$ as appropriate.

Note that our main parameter of interest is $\boldsymbol{\Phi}$, while the error covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}}$ is more of an unknown nuisance parameter that we need to deal with. One would hope that as the sample size n tends to infinity, the posterior probability assigned to any $\boldsymbol{\varepsilon}$ neighborhood of $\boldsymbol{\Phi}_{0n}$ converges to 1 almost surely under \mathbb{P}_0 . We now formally define a notion of posterior consistency that formalizes this.

Definition 1. The sequence of posterior distributions $\Pi_n(\cdot | \mathcal{Y}_n)$ is said to be consistent at $\{\boldsymbol{\Phi}_{0n}\}_{n \geq 1}$, if for every $\boldsymbol{\varepsilon} > 0$, $\Pi_n(\|\boldsymbol{\Phi} - \boldsymbol{\Phi}_{0n}\| > \boldsymbol{\varepsilon} | \mathcal{Y}_n) \rightarrow 0$ as $n \rightarrow \infty$ a.s. \mathbb{P}_0 .

For ease of exposition, we will henceforth denote $\boldsymbol{\Phi}_{0n}$ as $\boldsymbol{\Phi}_0$, and $\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon},0n}$ as $\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon},0}$, and highlight their dependence on n as needed.

2.1 Stability of VAR(d) process

Since VAR models are dynamical systems, the notion of ‘stability’ plays an important role in their analysis and asymptotic properties.

Definition 2. A VAR(d) process defined in (1) is said to be *stable* if the matrix valued polynomials $\mathcal{A}(z) := \mathbf{I}_p - \sum_{i=1}^d \mathbf{A}_i z^i$ satisfies $\det(\mathcal{A}(z)) \neq 0$ on the unit circle of the complex plane $\{z \in \mathbb{C}: |z| = 1\}$.

The autocovariance function of a p -dimensional centered covariance-stationary time series $\{\mathbf{X}^t\}$ is defined as $\Gamma_{\mathbf{X}}(h) = \text{Cov}(\mathbf{X}^t, \mathbf{X}^{t+h})$, $t, h \in \mathbb{Z}$ and the corresponding spectral density is given by $f_{\mathbf{X}}(\theta) := \frac{1}{2\pi} \sum_{h=-\infty}^{\infty} \Gamma_{\mathbf{X}}(h) e^{-ih\theta}$, $\theta \in [-\pi, \pi]$. For a Gaussian stable VAR(d) model the spectral density has a closed form expression,

$$f_{\mathbf{X}}(\theta) = \frac{1}{2\pi} \left(\mathbf{I}_p - \sum_{j=1}^d \mathbf{A}_j e^{-ij\theta} \right)^{-1} \boldsymbol{\Sigma}_\varepsilon \left(\mathbf{I}_p - \sum_{j=1}^d \mathbf{A}_j e^{-ij\theta} \right)^{-1*},$$

where $*$ denotes the Hermitian conjugate of a matrix and $i \equiv \sqrt{-1}$. The autocovariance function which characterizes a centered Gaussian process, can be used to quantify the temporal and cross-sectional dependence for VAR(d) models. The peak of the spectral density, measured by its maximum eigenvalue $\mathcal{M}(f_{\mathbf{X}}) := \max_{\theta \in [-\pi, \pi]} \lambda_{\max}(f_{\mathbf{X}}(\theta))$ can be used as a measure of stability of the process. Also the minimum eigenvalue $\mathfrak{m}(f_{\mathbf{X}}) := \min_{\theta \in [-\pi, \pi]} \lambda_{\min}(f_{\mathbf{X}}(\theta))$ captures cross-dependence among its components. However, as mentioned in Basu and Michailidis (2015) instead of working with $\mathcal{M}(f_{\mathbf{X}})$ and $\mathfrak{m}(f_{\mathbf{X}})$ it is often easier to work with the eigenvalues of $\mathcal{A}^*(z)\mathcal{A}(z)$ over the unit circle $\{z \in \mathbb{C} : |z| = 1\}$. Let

$$\begin{aligned} \mu_{\min}(\mathcal{A}) &:= \min_{|z|=1} \lambda_{\min}(\mathcal{A}^*(z)\mathcal{A}(z)) = \min_{\theta \in [-\pi, \pi]} \lambda_{\min} \left(\left(\mathbf{I}_p - \sum_{i=1}^d \mathbf{A}_i e^{ij\theta} \right) \left(\mathbf{I}_p - \sum_{i=1}^d \mathbf{A}_i e^{-ij\theta} \right) \right) \\ \mu_{\max}(\mathcal{A}) &:= \min_{|z|=1} \lambda_{\max}(\mathcal{A}^*(z)\mathcal{A}(z)) = \max_{\theta \in [-\pi, \pi]} \lambda_{\max} \left(\left(\mathbf{I}_p - \sum_{j=1}^d \mathbf{A}_j e^{ij\theta} \right) \left(\mathbf{I}_p - \sum_{j=1}^d \mathbf{A}_j e^{-ij\theta} \right) \right) \end{aligned}$$

For stable VAR(d) process $0 < \mu_{\min}(\mathcal{A}) \leq \mu_{\max}(\mathcal{A}) < \infty$. Since each \mathcal{A} is i.i.d as $\mathcal{N}_p(0, \boldsymbol{\Sigma}_\varepsilon)$, each row of \mathbf{X} is distributed as $\mathcal{N}_{dp}(0, \mathbf{C}_X)$, where the covariance matrix \mathbf{C}_X has the following structure,

$$\mathbf{C}_X = \begin{bmatrix} \Gamma(0) & \Gamma(1) & \cdots & \Gamma(d-1) \\ \Gamma(1)' & \Gamma(0) & \cdots & \Gamma(d-2) \\ \vdots & \vdots & \ddots & \vdots \\ \Gamma(d-1)' & \Gamma(d-2)' & \cdots & \Gamma(0) \end{bmatrix}_{dp \times dp} \quad (3)$$

The quantities $\mu_{\min}(\mathcal{A})$ and $\mu_{\max}(\mathcal{A})$ provide a useful bound for the eigenvalues of \mathbf{C}_X . As mentioned in Melnyk and Banerjee (2016) and from Proposition 2.3 and equation (2.6) of Basu and Michailidis (2015) we have the following chain of inequalities,

$$\frac{\lambda_{\min}(\boldsymbol{\Sigma}_\varepsilon)}{\mu_{\max}(\mathcal{A})} \leq 2\pi\mathbf{m}(f_X) \leq \lambda_{\min}(\mathbf{C}_X) \leq \lambda_{\max}(\mathbf{C}_X) \leq 2\pi\mathcal{M}(f_X) \leq \frac{\lambda_{\max}(\boldsymbol{\Sigma}_\varepsilon)}{\mu_{\min}(\mathcal{A})}. \quad (4)$$

We finally note that the p dimensional VAR(d) model in (2) can be equivalently written as a dp dimensional VAR(1) process. Let

$$\tilde{\mathbf{X}}^t = \begin{bmatrix} X^t \\ \vdots \\ X^{t-d+1} \end{bmatrix}_{dp \times 1}, \quad \tilde{\mathbf{A}}_1 = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_{d-1} & \mathbf{A}_d \\ \mathbf{I}_p & \mathbf{O} & \cdots & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I}_p & \cdots & \mathbf{O} & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{I}_p & \mathbf{O} \end{bmatrix}_{dp \times dp} \quad \text{and} \quad \omega^t = \begin{bmatrix} \varepsilon^{t+1} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}_{dp \times 1}.$$

Then the new representation becomes

$$\tilde{\mathbf{X}}^t = \tilde{\mathbf{A}}_1 \tilde{\mathbf{X}}^{t-1} + \omega^t \quad t = d, \dots, n+d-1. \quad (5)$$

It follows that $\mathbf{X} = [\tilde{\mathbf{X}}^{n+d-2} \tilde{\mathbf{X}}^{n+d-3} \dots \tilde{\mathbf{X}}^d]'$, i.e., i -th row of \mathbf{X} is denoted as ($dp \times 1$ vector) $\tilde{\mathbf{X}}^{n+d-i-1}$. Note that if the underlying VAR(d) process $\{\mathbf{X}^t\}$ is stable then the process $\tilde{\mathbf{X}}^t$ with the characteristic polynomial, $\tilde{\mathcal{A}}(z) := \mathbf{I}_{dp} - \tilde{\mathbf{A}}_1 z$ is also stable. This is because $\{\tilde{\mathbf{X}}^t\}$ can be viewed as generated according to a stable VAR(1) process with transition matrix $\tilde{\mathbf{A}}_1$ and $\{\tilde{\mathbf{X}}^t\}$ is stable if and only if $\{\mathbf{X}^t\}$ is stable (Lütkepohl (2007)). Based on $\tilde{\mathcal{A}}(z)$ we define

$$\begin{aligned} \mu_{\min}(\tilde{\mathcal{A}}) &:= \min_{\theta \in [-\pi, \pi]} \lambda_{\min}((\mathbf{I}_p - \tilde{\mathbf{A}}_1' e^{i\theta})(\mathbf{I}_p - \tilde{\mathbf{A}}_1 e^{-i\theta})) \\ \mu_{\max}(\tilde{\mathcal{A}}) &:= \max_{\theta \in [-\pi, \pi]} \lambda_{\max}((\mathbf{I}_p - \tilde{\mathbf{A}}_1' e^{i\theta})(\mathbf{I}_p - \tilde{\mathbf{A}}_1 e^{-i\theta})) \end{aligned} \quad (6)$$

While $\mu_{\min}(\tilde{\mathcal{A}})$ and $\mu_{\max}(\tilde{\mathcal{A}})$ are not necessarily the same as $\mu_{\min}(\mathcal{A})$ and $\mu_{\max}(\mathcal{A})$ the inequalities in (4) still hold with $\mu_{\min}(\mathcal{A})$ and $\mu_{\max}(\mathcal{A})$ replaced by $\mu_{\min}(\tilde{\mathcal{A}})$ and $\mu_{\max}(\tilde{\mathcal{A}})$, respectively.

3 Bayesian Estimation and Posterior Consistency

In this section, we first discuss Bayesian estimation of VAR models with non-hierarchical and hierarchical scale mixture matrix normal prior distributions on the parameter matrix Φ (conditioned on $\boldsymbol{\Sigma}_\varepsilon$) and subsequently establish high-dimensional posterior consistency in this setting under mild regularity assumptions. We start by introducing the necessary

notation for the *matrix-variate normal distribution*. Let $\mathbf{M}_{a,b}$ denote the space of $a \times b$ matrices.

Definition 3. An $a \times b$ random matrix \mathbf{X} is defined to follow a matrix-variate normal distribution ($\mathcal{MN}_{a \times b}(\mathbf{M}, \mathbf{B}_1, \mathbf{B}_2)$) if its density function (on the space $\mathbf{M}_{a,b}$) is given by

$$\frac{|\mathbf{B}_1|^{-b/2} |\mathbf{B}_2|^{-a/2} e^{-\frac{1}{2} \text{tr}\{\mathbf{B}_1^{-1}(\mathbf{X} - \mathbf{M})\mathbf{B}_2^{-1}(\mathbf{X} - \mathbf{M})\}}}{(2\pi)^{ab/2}}.$$

Here $\mathbf{M} \in \mathbf{M}_{a,b}$, $\mathbf{B}_1 \in \mathbf{M}_{a,a}$ and $\mathbf{B}_2 \in \mathbf{M}_{b,b}$ which are both positive definite matrices corresponding to the variances among the rows and columns of \mathbf{X} , respectively. Note that the matrix normal distribution is related to the multivariate normal distribution in the following way: $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\mathbf{M}, \mathbf{B}_1, \mathbf{B}_2)$, if and only if $\text{vec}(\mathbf{X}) \sim \mathcal{N}_{np}(\text{vec}(\mathbf{M}), \mathbf{B}_2 \otimes \mathbf{B}_1)$.

3.1 Non-Hierarchical Matrix Normal Prior

We consider a matrix normal prior for Φ conditional on Σ_{ϵ} , and a flat (uniform) prior on Σ_{ϵ} . In particular,

$$\Phi | \Sigma_{\epsilon} \sim \mathcal{MN}_{dp \times p}(\mathbf{O}, \mathbf{U}^{-1}, \Sigma_{\epsilon}) \text{ and } \pi(\Sigma_{\epsilon}) \propto 1, \quad (7)$$

where \mathbf{U} is a $dp \times dp$ known positive definite matrix. Note that under this matrix normal prior \mathbf{U}^{-1} and Σ_{ϵ} are the covariance matrices corresponding to the columns and rows of Φ , respectively. The posterior distribution of Φ (conditional on Σ_{ϵ}) can easily be shown to be $\mathcal{MN}_{dp \times p}(\Phi_{\text{PM}}, (\mathbf{X}'\mathbf{X} + \mathbf{U})^{-1}, \Sigma_{\epsilon})$, where $\Phi_{\text{PM}} := (\mathbf{X}'\mathbf{X} + \mathbf{U})^{-1}\mathbf{X}'\mathbf{Y}$ is the (conditional) posterior mean which does not depend on Σ_{ϵ} . Hence, the unconditional posterior mean of Φ is available in closed form and is given by Φ_{PM} . It follows by standard computations using the multivariate normal density that the marginal posterior density of Σ_{ϵ} is proportional to

$$|\Sigma_{\epsilon}|^{-n/2} e^{-\text{tr}(\Sigma_{\epsilon}^{-1} \widehat{\Sigma}_{\text{res}})},$$

Where $\widehat{\Sigma}_{\text{res}} = \mathbf{Y}'(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X} + \mathbf{U})^{-1}\mathbf{X}')\mathbf{Y}$. This density is proper if and only if $n > 2p$. In this case, the marginal posterior density of Σ_{ϵ} corresponds to the Inverse-Wishart density with scale parameter $\widehat{\Sigma}_{\text{res}}$ and shape parameter $n - p - 1$. We summarize the above observations in the lemma below.

Lemma 3.1. *Under the non-hierarchical prior in (7), the posterior density of $(\Phi, \Sigma_{\epsilon})$ is proper if and only if $n > 2p$. In this case*

$$\begin{aligned}\Phi \mid \Sigma_{\varepsilon}, \mathcal{Y} &\sim \mathcal{M}\mathcal{N}_{d \times p} \left(\Phi_{PM}, (X'X + U)^{-1}, \Sigma_{\varepsilon} \right) \\ \Sigma_{\varepsilon} \mid \mathcal{Y} &\sim \text{Inverse-Wishart}(\tilde{\Sigma}_{res}, n - p - 1).\end{aligned}$$

Assumptions for posterior consistency—We will establish our results under the high-dimensional setting from Section 2. Recall that $\Phi_0 = \Phi_{0n}$ denotes the true underlying parameter matrix, and $\Sigma_{\varepsilon,0} = \Sigma_{\varepsilon,0n}$ denotes the true underlying error covariance matrix in this setting. The quantities $\mu_{\min}(\tilde{\mathcal{A}})$, $\mu_{\max}(\tilde{\mathcal{A}})$ and C_X are as defined in (6) and (3), but with Φ_0 and $\Sigma_{\varepsilon,0}$ as the underlying parameter values. We assume the following:

Assumption A1. The VAR(d) model given in (1) is stable.

Assumption A2. $\frac{1 + \mu_{\max}(\tilde{\mathcal{A}})}{\mu_{\min}(\tilde{\mathcal{A}})}$ is $O\left(\sqrt{\frac{5n}{p}}\right)$ as $n \rightarrow \infty$.

Assumption A3. $0 < \inf_{n \geq 1} \lambda_{\min}(C_X) < \infty$ and $\lambda_{\max}(\Sigma_{\varepsilon,0n}) = O(1)$.

Assumption A4. The true parameter matrix of the VAR model (2), Φ_0 and the hyperparameter U of (7) are such that $\|\Phi_0^T U \Phi_0\| = o(n)$ and $\|U \Phi_0\| = o(n)$.

Assumption A5. $p = o(n)$.

When $d = 1$, we deal with a VAR(1) model and C_X becomes $\Gamma_X(0)$, while $\mu_{\min}(\tilde{\mathcal{A}})$ is the same as $\mu_{\min}(\mathcal{A})$ and $\mu_{\max}(\tilde{\mathcal{A}})$ is also equal to $\mu_{\max}(\mathcal{A})$. Assumption A1 is a standard assumption which ensures that the underlying VAR process is well-behaved. Assumption A2 plays an important role in deriving high-dimensional concentration bounds for $X'X/n$ and $X'E/n$ around C_X and O , respectively (see Propositions B.1 and B.3 in the supplementary material). Assumption A3 is needed to ensure that $\lambda_{\min}(X'X/n)$ is bounded away from 0 with high probability. Further, if we consider that each column of Φ is independently and identically distributed according to a normal prior distribution, that is $U = I_{dp}$, Assumption A4 reduces to $\Phi_0' \Phi_0 = o(n)$ and $\|\Phi_0\| = o(n)$.

We now state the main theoretical result of posterior consistency with a non-hierarchical matrix normal prior distribution on Φ . The proof is given in Appendix C.1 of the supplementary material.

Theorem 3.2 (Posterior consistency for non-hierarchical prior). *For any centered VAR(d) model (2) with non-hierarchical prior (7) on Φ satisfying Assumptions A1-A5, the posterior consistency of the parameter matrix can be achieved i.e. for every fixed every $\varepsilon > 0$*

$$\mathbb{E}_0 \left[\Pi_n \left(\|\Phi - \Phi_0\| > \varepsilon \mid \mathcal{Y} = (X^0, \dots, X^n) \right) \right] \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

where Φ_0 is the true parameter matrix under the model (2).

A natural question to ask is whether the assumption $p = o(n)$ can be relaxed for posterior consistency. In the lemma below, we consider a situation in which Assumptions A1 - A4 are satisfied and p is the same order as n , and prove that the resulting posterior is *not consistent*. The proof is given in Appendix C.2 of the supplementary material.

Lemma 3.3. *Consider a (sequence of) VAR(1) model with $p_n = \gamma^n$, $\Phi_{0n} = \alpha \mathbf{I}_{p_n}$ and $\Sigma_{e,0n} = \mathbf{I}_{p_n}$, where $\gamma \in (0, \frac{1}{2})$, $\alpha \in (0, 1)$ do not depend on n . If we use the non-hierarchical prior (7) on Φ with $\|U\| = o(n)$, then there exists $\varepsilon > 0$ such that*

$$\liminf_{n \rightarrow \infty} \mathbb{E}_0 \left[\mathbb{P}_n \left(\|\Phi - \Phi_0\| > \varepsilon \mid \mathcal{Y} = (X^0, \dots, X^n) \right) \right] > 0.$$

Remark Note that the condition $\|U\| = o(n)$ assumed in Lemma 3.3 corresponds to Assumption A4 in the setting of the Lemma. The reason for making this assumption is that we want to show violating Assumption A5 ($p = \alpha n$) can lead to posterior inconsistency, even if all of Assumptions A1-A4 hold. If we decide to violate Assumption A4 too by assuming $\|U\| = o(n)$ or $\|U\| \gg n$ (goes to ∞ at the same rate or faster than n) then the posterior inconsistency proof becomes comparatively easier. We have provided the corresponding proofs in Supplemental Section C.3 and Supplemental Section C.4, respectively.

3.2 Hierarchical Normal-mixture Prior

Next, we study the posterior consistency of the parameter matrix in model (2) in which Φ has the following hierarchical prior

$$\begin{aligned} \Phi \mid \Sigma_\varepsilon, \mathbf{U} &\sim \mathcal{MN}_{d_p \times p}(\mathbf{O}, \mathbf{U}^{-1}, \Sigma_\varepsilon), \quad (8) \\ \pi(\Sigma_\varepsilon) &\propto 1, \\ \text{and } \mathbf{U} &\sim \pi_{\text{scl}}(\cdot) \end{aligned}$$

where \mathbf{U} is the $d_p \times d_p$ matrix having probability density $\pi_{\text{scl}}(\cdot)$ over the space of $d_p \times d_p$ positive definite matrices, $\mathbb{M}_{d_p}^+$. As shown below, the group lasso and multivariate t distribution prior on Φ can be obtained from (8) using appropriate choices of $\pi_{\text{scl}}(\cdot)$. The lemma below shows that the posterior is proper if $n > (d+1)p$, and provides the form of various conditional and marginal posterior densities. The proof is given in Appendix C.5 of the supplementary material.

Lemma 3.4. *Under the hierarchical normal-mixture prior in (8), the posterior density of $(\Phi, \Sigma_\varepsilon, \mathbf{U})$ is proper if $n > (d+1)p$. In this case*

$$\Phi | \Sigma_\epsilon, U, \mathcal{Y} \sim \mathcal{MN}_{dp \times p}(\Phi_{PM}, (X'X + U)^{-1}, \Sigma_\epsilon)$$

$$\Sigma_\epsilon | U, \mathcal{Y} \sim \text{Inverse-Wishart}(\widehat{\Sigma}_{res}, n - p - 1)$$

$$\pi(U | \mathcal{Y}) \propto \frac{|U|^{dp/2}}{|X'X + U|^{dp/2} |\widehat{\Sigma}_{res}|^{(n-p-1)/2}} \pi_{\text{scl}}(U).$$

The Bayesian group lasso prior was proposed by Kyung et al. (2010) in the context of linear regression. We adapt it to the VAR setting as follows. Suppose the rows of Φ are divided in G groups $\Phi_{[1]}, \dots, \Phi_{[G]}$, where each $\Phi_{[g]}$ is an $m_g \times p$ sub-matrix of Φ (hence $\sum m_g = dp$) and X_g is the submatrix of X of order $n \times m_g$ corresponding to the group $\Phi_{[g]}$. The frequentist group lasso estimator (conditional on Σ_ϵ) is obtained by solving

$$\min_{\Phi_{[1]}, \dots, \Phi_{[G]}} \left\| \Sigma_\epsilon^{-1/2} \left(Y - \sum_{g=1}^G X_g \Phi_{[g]} \right) \right\|_F^2 + \sum_{g=1}^G \lambda_g \left\| \Phi_{[g]} \Sigma_\epsilon^{-1/2} \right\|_F,$$

where λ_g is a tuning parameter corresponding to the group g . The group lasso estimator (conditional on Σ_ϵ) can also be expressed as the maximum a posteriori probability (MAP) estimate under model (2) with the prior

$$\pi(\Phi | \Sigma_\epsilon) \propto \exp \left(- \sum_{g=1}^G \lambda_g \left\| \Phi_{[g]} \Sigma_\epsilon^{-1/2} \right\|_F \right)$$

which is a multivariate generalization of the double exponential prior and can also be expressed as a scale mixture of normals with Gamma hyperpriors (Park and Casella (2008), Kyung et al. (2010)) leading to the group lasso hierarchy,

$$\Phi_{[g]} | \tau_g, \Sigma_\epsilon \stackrel{\text{ind}}{\sim} \mathcal{MN}_{m_g \times p}(\mathbf{0}, \tau_g \mathbf{I}_{m_g}, \Sigma_\epsilon) \quad \text{and} \quad \tau_g \stackrel{\text{ind}}{\sim} \text{Gamma} \left(\frac{m_g + 1}{2}, \frac{\lambda_g^2}{2} \right), \quad g = 1, \dots, G.$$

Here $\text{Gamma}(\alpha, \lambda)$ denotes the Gamma distribution with shape parameter α and rate parameter λ . This can be alternatively presented as

$\Phi | \tau, \Sigma_\epsilon \sim \mathcal{MN}_{dp \times p}(\mathbf{0}, \text{BDiag}(\tau_1, \dots, \tau_G), \Sigma_\epsilon)$ and $\tau_g \stackrel{\text{ind}}{\sim} \text{Gamma} \left(\frac{m_g + 1}{2}, \frac{\lambda_g^2}{2} \right)$ where $\text{BDiag}(\tau_1, \dots, \tau_G)$ denotes a block-diagonal matrix with g -th block to be $\tau_g \mathbf{I}_{m_g}$. Note that under the above hierarchical prior, conditionally on (τ_1, \dots, τ_G) and Σ_ϵ , the columns of Φ are independent. If $m_g = 1 \forall g = 1, \dots, dp$ we get the ordinary Bayesian lasso.

Under the specification given in (8), suppose we assume $\mathbf{U} = \text{Diag}(\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{dp})$ and $1/\tau_i \stackrel{\text{ind}}{\sim} \text{Gamma}(\alpha_i, \lambda_i/2)$ then it can be shown that the prior density for $\boldsymbol{\Phi}$ given only $\boldsymbol{\Sigma}_\epsilon$ is proportional to

$$\prod_{i=1}^{dp} \left(\|\boldsymbol{\Phi}_i \cdot \boldsymbol{\Sigma}_\epsilon^{-1/2}\|_2^2 + \lambda_i \right)^{-(\alpha_i + \frac{1}{2})}$$

which corresponds to the multivariate t -distribution.

Estimation—For the hierarchical model given in (8), the posterior density of $\boldsymbol{\Phi}$ is intractable and quantities such as the posterior mean are not available in closed form. Hence, we develop a Markov Chain Monte Carlo algorithm to generate values from the posterior density. It follows by straightforward calculations that

$$\begin{aligned} \boldsymbol{\Phi} | \boldsymbol{\Sigma}_\epsilon, \mathbf{U}, \mathcal{Y} &\sim \mathcal{M}\mathcal{N}_{dp \times p}(\boldsymbol{\Phi}_{\text{PM}}, (\mathbf{X}'\mathbf{X} + \mathbf{U})^{-1}, \boldsymbol{\Sigma}_\epsilon) \\ \boldsymbol{\Sigma}_\epsilon | \mathbf{U}, \mathcal{Y} &\sim \text{Inverse-Wishart}(\widehat{\boldsymbol{\Sigma}}_{res}, n - p - 1) \\ \pi(\mathbf{U} | \boldsymbol{\Phi}, \boldsymbol{\Sigma}_\epsilon, \mathcal{Y}) &\propto |\mathbf{U}|^{dp/2} \exp\left[-\frac{1}{2}\text{tr}\{\boldsymbol{\Phi}\boldsymbol{\Sigma}_\epsilon^{-1}\boldsymbol{\Phi}'\mathbf{U}\}\right] \pi_{scf}(\mathbf{U}), \end{aligned}$$

where $\widehat{\boldsymbol{\Sigma}}_{res} = \mathbf{Y}'(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X} + \mathbf{U})^{-1}\mathbf{X}')\mathbf{Y}$. While the conditional posterior distribution of $\boldsymbol{\Phi}$ given $\boldsymbol{\Sigma}_\epsilon, \mathbf{U}$ and $\boldsymbol{\Sigma}_\epsilon$ given \mathbf{U} are easy to simulate from (being Matrix-normal and InverseWishart), the tractability of the conditional posterior density of \mathbf{U} given $\boldsymbol{\Phi}, \boldsymbol{\Sigma}_\epsilon$ depends on the form of the prior $\pi_{scf}(\mathbf{U})$. We show below that for three standard choices of $\pi_{scf}(\mathbf{U})$ corresponding to the Wishart prior, the group lasso prior and the multivariate t -prior $\pi(\mathbf{U} | \boldsymbol{\Phi})$ becomes a tractable density and easy to simulate from.

Case 1: Wishart Prior—For a $dp \times dp$ positive definite matrix \mathbf{V} , let $\mathbf{U} \sim \text{Wishart}_{dp}(\mathbf{V}, df = \boldsymbol{\nu} + dp)$ that is $\pi(\mathbf{U}) \propto |\mathbf{U}|^{\frac{\nu-1}{2}} \exp\left[-\frac{1}{2}\text{tr}\{\mathbf{V}^{-1}\mathbf{U}\}\right]$. In this case,

$$\pi(\mathbf{U} | \boldsymbol{\Phi}, \boldsymbol{\Sigma}_\epsilon, \mathcal{Y}) \propto |\mathbf{U}|^{\frac{\nu + dp - 1}{2}} \exp\left[-\frac{1}{2}\text{tr}\{(\boldsymbol{\Phi}\boldsymbol{\Sigma}_\epsilon^{-1}\boldsymbol{\Phi}' + \mathbf{V}^{-1})\mathbf{U}\}\right]$$

which is $\text{Wishart}_{dp}((\boldsymbol{\Phi}\boldsymbol{\Sigma}_\epsilon^{-1}\boldsymbol{\Phi}' + \mathbf{V}^{-1})^{-1}, df = \nu + 2dp)$. Note that as long as we have $\boldsymbol{\nu} > -(dp + 1)$ the posterior of \mathbf{U} given $\boldsymbol{\Phi}, \boldsymbol{\Sigma}_\epsilon, \mathcal{Y}$ is proper.

Case 2: Bayesian Group Lasso—In this case as already discussed in Subsection 3.2, \mathbf{U}^{-1} has a block diagonal form $\text{BDiag}(\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_G)$ and $\boldsymbol{\tau}_g$'s are apriori independently distributed as Gamma with scale $(m_g + 1)/2$ and rate $\lambda_g^2/2$. Hence, the conditional distribution of $\boldsymbol{\tau}_g$ has the following form,

$$\frac{1}{\tau_g} | \Phi, \Sigma_\epsilon \mathcal{Y} \stackrel{ind}{\sim} \text{Inverse - Gaussian} \left(\mu_g = \frac{\lambda_g}{\| \Phi_{[g]} \Sigma_\epsilon^{-1/2} \|_F}, \lambda_g^2 \right).$$

Case 3: Multivariate t-distribution—By taking $\mathbf{U}^{-1} = \text{Diag}(\tau_1, \dots, \tau_{dp})$ and $1/\tau_i$ to be independently distributed as Gamma with scale α_i and rate $\lambda_i/2$ we have the multivariate t -distribution as the prior on Φ . In this case, the conditional distribution of τ_i has the following form

$$\frac{1}{\tau_i} | \Phi, \Sigma_\epsilon \mathcal{Y} \stackrel{ind}{\sim} \text{Gamma} \left(\alpha_i + \frac{dp}{2}, \frac{\| \Phi'_i \Sigma_\epsilon^{-1/2} \|_2^2 + \lambda_i}{2} \right).$$

Assumptions for posterior consistency—We now introduce regularity conditions to establish posterior consistency under the hierarchical prior model.

Assumption B1. The VAR(d) model given in (1) is stable.

Assumption B2. $\frac{1 + \mu_{\max}(\mathcal{A})}{\mu_{\min}(\mathcal{A})}$ is $O\left(\sqrt{\frac{n}{p}}\right)$ as $n \rightarrow \infty$.

Assumption B3. $0 < \inf_{n \geq 1} \lambda_{\min}(\mathbf{C}_X) \leq \sup_{n \geq 1} \lambda_{\max}(\mathbf{C}_X) < \infty$ and $0 < \lambda_{\max}(\Sigma_{\epsilon,0n}) = O(1)$.

Assumption B4. The singular values of the true parameter matrices $\{\Phi_{0n}\}_{n \geq 1}$ are uniformly bounded. Equivalently, the eigenvalues of $\{\Phi'_{0n} \Phi_{0n}\}_{n \geq 1}$ are uniformly bounded.

Assumption B5. $p = O\left(\frac{n}{\log n}\right)$.

Assumption B6. There exists (fixed and not-depending on n) $\alpha > 0$ such that $\liminf_{n \rightarrow \infty} \pi_{\text{scl},n}(\lambda_{\max}(\mathbf{U}) > \alpha) > 0$ and for every $\beta > 0$ we have $\lim_{n \rightarrow \infty} \pi_{\text{scl},n}(\lambda_{\max}(\mathbf{U}) > \beta n) = 0$.

We now discuss these assumptions and compare them to the assumptions for the non-hierarchical prior model.

- Assumptions B1 and B2 are identical to A1 and A2, while B3 is fairly similar to A3.
- One key difference is the permissible scaling of p as a function of the sample size n in Assumption B5, which is slightly more stringent than the permissible scaling for the nonhierarchical matrix normal prior in Assumption A5.
- Note that Assumption B6 is a mild one. For example, a sufficient condition for this assumption to be satisfied is that $\limsup_{n \rightarrow \infty} \max_{1 \leq i \leq p} \mathbb{E} \pi_{\text{scl},n}[\mathbf{U}_{ii}^\delta] < \infty$ for

some $\delta > 0$ and $\liminf_{n \rightarrow \infty} \pi_{\text{scl},n}(\mathbf{U}_{11} > \alpha) > 0$. It can be easily checked that this condition, and hence Assumption B6, is satisfied in the case of Wishart, Inverse-Wishart, Bayesian group lasso, multivariate t -distribution, Horseshoe (Carvalho et al. (2010)), Strawderman-Berger and generalized double Pareto (Armagan et al. (2013)) priors as long as the prior parameters do not depend on n .

- In the non-hierarchical prior case, assumptions regarding Φ_0 and \mathbf{U} (non-random) are simultaneously and exclusively provided in Assumption A4 through the conditions $\|\Phi_0^T \mathbf{U} \Phi_0\| = o(n)$ and $\|\mathbf{U} \Phi_0\| = O(n)$. For the hierarchical prior case, for clarity of exposition, we provide the assumptions regarding Φ_0 in Assumption B4, and those regarding the distribution of \mathbf{U} (random) in Assumption B6. Combining these two assumptions it can be easily shown that a priori $\|\Phi_0^T \mathbf{U} \Phi_0\|$ and $\|\mathbf{U} \Phi_0\|$ converge to zero in $\Pi_{\text{scl},n}$ -probability as $n \rightarrow \infty$. In that sense, the assumptions on (Φ_0, \mathbf{U}) in the hierarchical model are stronger than in the non-hierarchical model case.

With these assumptions in hand, we state our key consistency result, whose proof is delegated to Appendix C.6 of the supplementary material.

Theorem 3.5 (Posterior Consistency for Hierarchical Prior). *For any centered VAR(d) model with the hierarchical prior (8) on the transition matrix satisfying Assumptions B1-B6, the posterior consistency of the transition matrix can be achieved i.e. for every fixed $ce > 0$*

$$\mathbb{E}_0 \left[\Pi_n \left(\|\Phi - \Phi_0\| > \varepsilon \mid \mathcal{Y} = (X^0, \dots, X^n) \right) \right] \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

where Φ_0 is the true parameter matrix under the model (2).

4 Performance Evaluation

To illustrate the performance of our Bayesian modeling framework for VAR processes, we design three sets of numerical experiments involving: (a) Small VAR ($p = 10$), (b) Medium VAR ($p = 100$) and (c) Large VAR ($p = 500$) models, each with two lags - (i) $d = 1$ and (ii) $d = 2$.

In each setting, we use transition matrices \mathbf{A}_i 's with 10–30% non-zero entries that are generated from $\mathbf{U}(0,2) \mathbf{U}\mathbf{U}(-2,0)$ selected at random and rescaled to ensure that the process is stable with $\text{SNR} = 2$. For small VAR models, we generate $n = 40, 80, 120$ time points, for medium VAR models, $n = 400, 800, 1200$, while for large VAR models we use $n = 2000, 4000, 6000$. The hyper-parameters for the prior distributions are selected using the Deviance Information Criterion (DIC). Note that $\text{DIC} = 2\bar{D} - D(\bar{\Phi}, \bar{\Sigma}_e)$, where

$D(\Phi, \Sigma_e) := -2 \log L(\mathcal{Y} \mid \Phi, \Sigma_e) = n \log |\Sigma_e| + \text{tr} \left\{ \Sigma_e^{-1} (\Phi' \mathbf{X}' \mathbf{X} \Phi - 2\Phi' \mathbf{X}' \mathbf{Y} + \mathbf{Y}' \mathbf{Y}) \right\}$, \bar{D} is the posterior expectation of $D(\Phi, \Sigma_e)$ and $\bar{\Phi}$ and $\bar{\Sigma}_e$ are the posterior expectation of Φ and Σ_e respectively.

4.1 Non-hierarchical prior

We generate two different error processes using $\Sigma_\varepsilon = \sigma^2 \mathbf{I}_p$ and $\Sigma_\varepsilon = \sigma^2 ((\rho^{|i-j|}))_{ij}$ (Toeplitz form). For each of the small, medium and large VAR models, \mathbf{U} is taken to be a diagonal matrix, $\mathbf{c} \mathbf{I}_p$ where \mathbf{c} is chosen according to the minimum DIC value over the interval $[0, 10]$. In Table 1, for both the posterior mean (PM) and least squares estimator (LS), we report their relative estimation error ($\|\widehat{\Phi} - \Phi_0\|_2 / \|\Phi_0\|_2$) and the standard error of $\|\widehat{\Phi}\|_2$ within parenthesis averaged over $10 \times p$ replicates for small and medium VAR and 100 replications for large VAR ($p = 500$). Since the true parameter matrix Φ_0 is sparse, we identify entries whose 95% posterior credible intervals contain zero, and set them to zero in both parameter matrix estimates (PM and LS).

First, we assume the true error covariance matrix Σ_ε is diagonal i.e. $\sigma^2 \mathbf{I}_p$. Here % denotes percentage of non-zero entries in Φ and d represents lag length of the underlying VAR process. Recall that the sample sizes used for small VAR models are $n_1 = 40$, $n_2 = 80$, $n_3 = 120$, for medium VAR ones $n_1 = 400$, $n_2 = 800$, $n_3 = 1200$, and for large VAR ones $n_1 = 2000$, $n_2 = 4000$, $n_3 = 6000$.

It can be seen that the relative estimation error decreases with an increase in the number of time points (sample size) n for both lags $d = 1, 2$; further, its values are significantly larger in medium and large size VAR models than in small VAR ones. Moreover, the estimation error for lag 1 is uniformly smaller than that for lag 2, and the same holds true for their respective standard errors. Regarding the percentage of non-zero entries in the true transition matrices, the results show that for fixed n and p , the more true non-zero entries in $\mathbf{A}_1, \mathbf{A}_2$, the less accurate the posterior mean and the LS estimator are, while their variability as indicated by their standard errors also follows the same pattern. However, the posterior mean clearly outperforms the LS estimates, especially in settings with large p . This is to a large extent due to the fact that the true transition matrices $\mathbf{A}_1, \mathbf{A}_2$ exhibit weaker signal as p or the number of non-zero edges increases (this is to ensure stability of the underlying VAR model) and due to our choice of $\mathbf{U} = \mathbf{c} \mathbf{I}_{dp}$ the posterior mean is the ridge regression estimator which applies shrinkage on the coefficients.

Next, we introduce correlation in the error components by specifying Σ_ε to be of Toeplitz form. As discussed in Section 3.2.1 of Lütkepohl (2007) the generalized least squares estimate in this multivariate regression set up is the same as the ordinary one; i.e. $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$, a result due to Zellner (1962). In Table 2, we compare the performance of least squares and posterior (ridge) estimates with noise covariance $\Sigma_\varepsilon = \text{Toeplitz}(\rho = 0.8)$.

In this setting, the relative estimation error of both the least squares and ridge estimators increases compared to that with an uncorrelated error structure given in Table 1; in particular, the performance of the LS estimator deteriorates even further. However, with an increase in sample size, the accuracy of both estimates significantly improves. Further, as gleaned from the entries of the Table corresponding to lag 2, the relative error exhibits a further increase, a pattern consistent with the results in Table 1. This is quite expected as we

not only have Toeplitz type error covariance structure, but also the total number of unknown parameters has increased by p^2 .

Finally, we study the support recovery under both error processes. In Table 3, we provide the percentage of true positives recovered by using 95% posterior credible intervals based on the same sample sizes n_1, n_2 , and n_3 as used previously.

The above table indicates that support recovery is not sensitive to the sample size, or to the lag; however, it deteriorates for all VAR models and error covariance settings, as the density of non-zero entries increases and exhibits a small increase with model dimension.

4.2 Hierarchical Priors

As discussed in Section 3.2, three types of hierarchical priors (Wishart, group-lasso and multivariate t) are studied. Analogously to the non-hierarchical prior case, the performance of the LS estimator is not at all satisfactory in this set up as well. Thus, we only compare the relative accuracy of the three prior choices in this setting. We select $V = cI_p$ and $df = \nu = dp$ for the Wishart prior, $\lambda_i = \lambda$ for all $1 \leq i \leq dp$ for the group-lasso prior and $\alpha_i = 1, \lambda_i = \lambda$ for all $1 \leq i \leq dp$ for the multivariate t prior. The hyper-parameters c and λ are chosen using DIC. In Table 4, we report the relative estimation error $\left(\frac{\|\widehat{\Phi} - \Phi_0\|_2}{\|\Phi_0\|_2}, d = 1, 2 \right)$ of the three hierarchical estimators when the error process covariance is set to $\sigma^2 I_p$ and d represents lag length in the underlying VAR model.

In the following table, we present relative estimation errors with the same 3 hierarchical priors when $\Sigma_e = \text{Toeplitz}(\rho = 0.6)$.

All of our hierarchical estimates outperform the ridge estimator (Table 1 and Table 2) across all settings considered. This is again expected, since the A_i 's have sparse structure by construction and the group lasso prior favors sparsity. However, the above results are not conclusive whether the group lasso estimate exhibits better accuracy than the Wishart or multivariate t estimates. To gain some insight into this issue, we use a VAR(1) model with $p = 9$ and transition matrix A_1 in which the columns form three groups each containing three columns. The sparsity increases as we move from group 1 to group 3. Finally, we rescale the coefficient matrix so that the corresponding VAR process is stable with SNR=2. The structure of the resulting A_1 transition matrix is depicted next, where * indicates non-zero entries.

$$A_1 = \begin{pmatrix} * & * & * & * & * & 0 & 0 & 0 \\ * & * & * & * & * & 0 & 0 & 0 \\ * & * & * & * & * & 0 & 0 & 0 \\ * & * & * & 0 & 0 & 0 & 0 & 0 \\ * & * & * & 0 & * & 0 & 0 & 0 \\ * & * & * & 0 & 0 & * & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & * & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & * \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}_{9 \times 9}$$

We generate $n = 100$ observations from the above VAR(1) model using two white noise variances (1) $\Sigma_{\epsilon} = \sigma^2 \mathbf{I}_p$ and (2) $\Sigma_{\epsilon} = \text{Toeplitz}(\rho = 0.80)$ and report the relative estimation error ($\|\widehat{\mathbf{A}}_1 - \mathbf{A}_1\|_2 / \|\mathbf{A}_1\|_2$) of five different estimates - least squares (LS), posterior mean for non-hierarchical normal prior (Non H), hierarchical Wishart (W), group lasso (GL) and multivariate t prior (Mult. t), in Table 6.

The results show that the group lasso estimator exhibits the best performance, followed by the multivariate t one, whereas the LS estimator is the least accurate. The result is consistent with the structure of the underlying transition matrix, since the group lasso prior can capitalize on it.

In Appendix A.1 of the supplementary document we illustrate the posterior estimate of a VAR(1) model transition matrix \mathbf{A}_1 , 95% credible intervals and estimated posterior densities of several entries of \mathbf{A}_1 . We also look into the performance of $\widehat{\Sigma}_{\epsilon}$ when the true error covariance is Toeplitz($\rho = 0.8$). The relative error of $\widehat{\Phi}$ under all the four different priors using a new noise covariance matrix which is generated from a Wishart distribution with degrees of freedom $\nu = p$ and scale matrix \mathbf{I}_p is also given.

5 Application to Macroeconomic Data

We use the proposed Bayesian framework to understand the lead-lag relationships in the FRED-MD dataset containing $p = 137$ key macroeconomic variables for the period January 1973 to December 2014. VAR modeling for this task was strongly advocated by Sims (1980) and since then has become a standard tool for it, although usually the focus is on small models involving few macroeconomic indicators (e.g. consumer price index, employment index and the federal funds rate). However, recent work has advocated for larger VAR models (see Bernanke et al. (2005); Ba bura et al. (2010) and references therein), in order to improve forecastability and also avoid the presence of hard to interpret or even contradictory to economic theory relationships, due to the fact of not including an adequate number of variables for properly modeling the economic phenomenon under consideration. Before centering the data and estimating Σ_{ϵ} as discussed earlier in Section 3.2, we ensure stationarity by processing the variables according to the recommendations in Stock and Watson (2005). The specific transformations used for each time series are given in the supplementary documents. Analogously to Ba bura et al. (2010), we consider the following 3 different size VAR models:

- **SMALL:** This model contains $p = 4$ key variables - CPI, number of employees non-farm (PAYEMS), Federal Funds Rate (FEDFUNDS) and Unemployment Rate (UNRATE).
- **MEDIUM:** In addition to the four variables in the SMALL VAR model, this one contains an additional 16 variables (total $p = 20$) listed next - Reserves Of Depository Institutions (NONBORRES), Total Reserves of Depository Institutions (TOTRESNS), M2 (M2REAL), Real Personal Income (RPI), Real personal consumption expenditures (DPCERA), IP Index (INDPRO), Capacity Utilization: Manufacturing (CAPT), Housing Starts: Total New Privately Owned

(HOUST), Avg Hourly Earnings : Goods-Producing (CES), M1 (M1), S & P's Common Stock Price Index: Composite (S.P.), 10-Year Treasury Rate (GS10), Personal Cons. Expend.: Chain Index (PCEPI), Foreign Exchange Rate (EXS), Crude Oil, spliced WTI and Cushing (OIL) and Retail and Food Services Sales (RETAILx).

- LARGE: This specification has all $p = 134$ macroeconomic indicators (3 were excluded from further analysis due to the presence of a large number of missing values).

Based on initial exploratory work, we choose lag $d = 6$ according to the Bayesian information criterion (BIC) and the following distributions were used for prior specification to obtain the estimated parameter matrix Φ : (i) non-hierarchical normal (Non H), (ii) hierarchical Wishart (W), (iii) group lasso (GL) and (iv) multivariate t prior (Mult. t). Since with an increase in the lag length d the number of parameters increases linearly we suggest using BIC over the Akaike information criterion (AIC). For the non-hierarchical prior, we use $\mathbf{U} = \text{Bdiag}(\lambda_1, \dots, \lambda_d)$, while for the hierarchical Wishart, group-lasso and multivariate t priors on Φ , we use $\mathbf{V} = c_1 \mathbf{I}_{dp}$ and $\alpha_j = \alpha$ for all $1 \leq j \leq dp$. The values of c_1 and λ are chosen using the Deviance Information Criterion (DIC) which, as explained previously, is a hierarchical Bayesian modeling generalization of BIC. The respective posterior means were compared to the least squares (LS) estimate. For each of the estimates $\hat{\Phi}$, the residual norm ratio $(\|\mathbf{Y} - \mathbf{X}\hat{\Phi}\|_F / \|\mathbf{Y}\|_F)$ which measures the in-sample fit, is reported in Table 7.

Note that since the LS estimator is obtained by minimizing $\|\mathbf{Y} - \mathbf{X}\Phi\|_F$, it will always result in minimum relative residual norm as observed in the above table; i.e. the LS estimator is always the best one in terms of in-sample prediction accuracy.

Next, we investigate the 4 different Bayesian estimates based on their out-of-sample prediction performance with respect to the benchmark prior, analogously to the evaluation strategies discussed in Ba bura et al. (2010) and Stock and Watson (2005). We consider the following two benchmark priors:

1. Prior information is imposed exactly by setting $\mathbf{U}^{-1} = \mathbf{O}$ matrix (the zero matrix) and it corresponds to $\lambda = 0$ in the Minnesota prior. Ba bura et al. (2010) uses this specification as the benchmark prior, in which case the corresponding benchmark model becomes a random walk with drift; i.e. $X^t = \alpha + X^{t-1} + \epsilon^t$.
2. A uniform prior on Φ by setting $\mathbf{U} = \mathbf{O}$ which corresponds to $\lambda = \infty$ in the Minnesota prior. The posterior mean coincides with the least squared estimate (LS).

Let \hat{X}^{t+h} be the h -step ahead predicted value for X^{t+h} based on our posited Bayesian model and using the data upto time t . The corresponding forecast under the benchmark prior is \hat{X}_O^{t+h} . The mean squared forecast error relative to the benchmark (RMSFE) is defined to

$\frac{\sum_{t=T_0}^{T_1} \|X^{t+h} - \hat{X}^{t+h}\|_2}{\sum_{t=T_0}^{T_1} \|X^{t+h} - \hat{X}_O^{t+h}\|_2}$. Table 8 gives the RMSE results for three different choices of

forecasting horizons, $h = 1, 6, 12$, for the two benchmark priors considered, over the period $T_0 = \text{January 1978}$ to $T_1 = \text{December 2006}$. An RMSE value smaller than 1 implies the VAR model with the corresponding prior outperforms that with the naive/benchmark prior.

It can easily be seen that all four Bayesian methods not only outperform the LS estimate (uniform prior on Φ), but also exhibit substantially smaller relative error compared to the random walk with drift process (point-mass prior on Φ). Further, increasing the number of predictor variables improves forecasting performance, a point argued forcefully in favor of large VAR models by Sims (1980). On the other hand, forecasting performance deteriorates for larger values of h , an expected result. Nevertheless, even for $h = 12$ (one year ahead) the results are still very satisfactory. Further, for the SMALL and MEDIUM VAR models, the non-hierarchical normal and hierarchical Wishart priors result in better prediction, whereas for the LARGE VAR model the group lasso prior outperforms other forecasts.

Next, we examine closely the out-of-sample prediction performance of the following three macroeconomic variables - CPI, PAYEMS and FEDFUND under the hierarchical Wishart prior.

Note that Ba bura et al. (2010) only consider a random walk process as the naive prior. From the above table it can be seen that although for CPI and PAYEMS the LS estimate performs better than the Bayesian estimates in SMALL and MEDIUM VARs, overall the Wishart prior has better forecasting accuracy than both of the benchmark priors. As previously observed, adding information (i.e. including more variables) improves the accuracy of forecasts for all three variables. The fourth column (LARGE BGR) provides the numbers reported in Table III in Ba bura et al. (2010), where a Bayesian VAR model on the same 134 variables, with $d = 13$ lags was estimated using a normal-inverted Wishart distribution that leads to a ridge type posterior mean estimate for the parameters in Φ and based on data covering the period 1971–2003. Although the results are not directly comparable to those obtained by our methodology, they nevertheless provide a certain degree of calibration. It can be seen that our model is more parsimonious using only $d = 6$ lags and provides better forecasting performance for all three variables at all forecasting horizons.

Next, we examine in more detail the estimated transition matrix A_1 for the SMALL VAR model ($p = 4$) and the non-hierarchical normal and group lasso priors. Estimated posterior densities of the bold marked entries are shown next. It is worth noting that the non-hierarchical prior centers around a different value and exhibits a less smooth behavior, than the hierarchical one. This smoothness should be expected given the specification of the latter.

$$\widehat{\mathbf{A}}_1^{NonH} = \begin{matrix} & \begin{matrix} \text{CPI} \\ \text{PAYEMS} \\ \text{FEDFund} \\ \text{UnRate} \end{matrix} & \begin{pmatrix} \text{CPI} & \text{PAYEMS} & \text{FEDFund} & \text{UnRate} \\ -\mathbf{0.133} & -0.001 & 0.001 & 0.001 \\ -0.016 & 0.311 & 0.001 & -\mathbf{0.002} \\ -1.000 & \mathbf{10.200} & 0.498 & -0.185 \\ 1.217 & -\mathbf{23.300} & -0.035 & -0.105 \end{pmatrix} \end{matrix}$$

$$\widehat{\mathbf{A}}_1^{GL} = \begin{matrix} & \begin{matrix} \text{CPI} \\ \text{PAYEMS} \\ \text{FEDFund} \\ \text{UnRate} \end{matrix} & \begin{pmatrix} \text{CPI} & \text{PAYEMS} & \text{FEDFund} & \text{UnRate} \\ -\mathbf{0.167} & -0.005 & 0.001 & 0.001 \\ -0.021 & 0.560 & 0.001 & -\mathbf{0.001} \\ -1.614 & \mathbf{18.6079} & 0.486 & -0.134 \\ 1.609 & -\mathbf{41.818} & -0.015 & -0.107 \end{pmatrix} \end{matrix}$$

Further, we present the 95% posterior credible intervals (PCI) of \mathbf{A}_1 under the above two priors

$$\text{Non hierarchical:} = \begin{pmatrix} (-\mathbf{0.19}, -\mathbf{0.07}) & (-0.08, 0.07) & (+0, +0) & (-0, +0) \\ (-0.05, 0.01) & (0.27, 0.35) & (+0, +0) & (-\mathbf{0}, -\mathbf{0}) \\ (-7, 5.25) & (\mathbf{1.86}, \mathbf{19}) & (0.44, 0.55) & (-0.3, -0.06) \\ (-1.98, 4.45) & (-\mathbf{27.8}, -\mathbf{18.7}) & (-0.06, -0.01) & (-0.07, 0.06) \end{pmatrix}$$

$$\text{Group lasso:} = \begin{pmatrix} (-\mathbf{0.22}, -\mathbf{0.11}) & (-0.10, 0.1) & (+0, +0) & (-0, +0) \\ (-0.05, 0.01) & (0.51, 0.61) & (+0, +0) & (-\mathbf{0}, -\mathbf{0}) \\ (-8.58, 5.7) & (\mathbf{6.5}, \mathbf{31.01}) & (0.44, 0.55) & (-0.26, -0) \\ (-1.75, 4.86) & (-\mathbf{47.72}, -\mathbf{36.37}) & (-0.04, 0.01) & (-0.07, 0.06) \end{pmatrix}$$

Next, in Figure 2 we depict the estimated networks for the MEDIUM VAR based on the first lag transition matrix produced by: (a) least squares and (b) a non-hierarchical normal prior, where for ease of representation the nodes of the network are abbreviated; the full list of the variable names is given in Table A1 of Appendix A in the supplementary material.

As expected, for most variables their previous lag value influences the current value. Further, for the LS based network, there is a high degree of connectivity, whereas the non-hierarchical based one exhibits a sparser structure. For the latter, of interest is that the employment index (PAYEMS), the personal consumer expenditures (GS10) and CPI are influenced by many other variables. On the other hand, the Federal Funds Rate influences the broad stock market (SP500) as expected based on finance theory and GS10. In general the sparser result provided by the non-hierarchical prior, in addition to better forecasting also aids in interpretation, vis-a-vis the LS estimate.

6 Discussion

In this paper, we investigate posterior consistency in Bayesian VAR(d) models with both nonhierarchical and hierarchical matrix normal prior distributions on the transition matrices under a Gaussian assumption for the temporal evolution of the time series under consideration and in the presence of a general covariance matrix that captures additional contemporary dependence between them. We establish posterior consistency for both of these priors under high-dimensional scaling. To obtain the desired results, some novel

concentration inequalities are provided that are of independent interest. The methodology is illustrated on synthetic and real macroeconomic data. The proposed priors provide better forecasts than the LS estimates for periods up to one year ahead, while leading to sparser and potentially easier to interpret relationships, especially for large scale models.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

* The authors gratefully acknowledge support from NSF grants DMS-1511945 (KK) and IIS-1632730 and CNS-1422078 (GM) and NIH grant R01 5R01GM11402902.

References

- Armagan A, Clyde M, and Dunson DB (2011). Generalized beta mixtures of gaussians In Shawe-Taylor J, Zemel RS, Bartlett PL, Pereira F, and Weinberger KQ (Eds.), *Advances in Neural Information Processing Systems 24*, pp. 523–531. Curran Associates, Inc. [PubMed: 25364213]
- Armagan A, Dunson DB, and Lee J (2013). Generalized double pareto shrinkage. *Statistica Sinica* 23, 119–143. [PubMed: 24478567]
- Armagan A, Dunson DB, Lee J, Bajwa WU, and Strawn N (2013). Posterior consistency in linear models under shrinkage priors. *Biometrika* 100(4), 1011–1018.
- Ba bura M, Giannone D, and Reichlin L (2010). Large bayesian vector auto regressions. *Journal of Applied Econometrics* 25(1), 71–92.
- Basu S and Michailidis G (2015, 08). Regularized estimation in sparse high-dimensional time series models. *Ann. Statist* 43(4), 1535–1567.
- Basu S, Shojaie A, and Michailidis G (2015). Network granger causality with inherent grouping structure. *Journal of Machine Learning Research* 16, 417–453.
- Bernanke BS, Boivin J, and Elias P (2005). Measuring the effects of monetary policy: a factor-augmented vector autoregressive (FAVAR) approach. *The Quarterly Journal of Economics* 120(1), 387–422.
- Billio M, Getmansky M, Lo AW, and Pelizzon L (2012). Econometric measures of connectedness and systemic risk in the finance and insurance sectors. *Journal of Financial Economics* 104(3), 535–559.
- Bontemps D (2011, 10). Bernstein-von mises theorems for gaussian regression with increasing number of regressors. *Ann. Statist* 39f, 2557–2584.
- Carvalho CM, Polson NG, and Scott JG (2010). The horseshoe estimator for sparse signals. *Biometrika* 97(2), 465–480.
- Doan T, Litterman R, and Sims C (1984). Forecasting and conditional projection using realistic prior distributions. *Econometric reviews* 3(1), 1–100.
- Ghoshal S (1999). Asymptotic normality of posterior distributions in high-dimensional linear models. *Bernoulli* 5(2), 315–331.
- Griffin JE and Brown PJ (2010). Inference with normal-gamma prior distributions in regression problems. *Bayesian Anal.* 5(1), 171–188.
- Korobilis D (2013). Var forecasting using bayesian variable selection. *Journal of Applied Econometrics* 28, 204–230.
- Kyung M, Gill J, Ghosh M, and Casella G (2010, 06). Penalized regression, standard errors, and bayesian lassos. *Bayesian Anal.* 5(2), 369–411.
- Lee J and Oh H-S (2013). Bayesian regression based on principal components for high-dimensional data. *Journal of Multivariate Analysis* 117, 175–192.
- Lin J and Michailidis G (2017). Regularized estimation and testing for high-dimensional multiblock vector-autoregressive models. *Journal of Machine Learning Research* 18(117), 1–49.

- Litterman RB et al. (1979). Techniques of forecasting using vector autoregressions. Technical report.
- Lütkepohl H (2007). *New Introduction to Multiple Time Series Analysis*. Springer.
- Melnyk I and Banerjee A (2016). Estimating structured vector autoregressive models. *Proceedings of The 33rd International Conference on Machine Learning*, 830–839.
- Mol CD, Giannone D, and Reichlin L (2008). Forecasting using a large number of predictors: Is bayesian shrinkage a valid alternative to principal components? *Journal of Econometrics* 146(2), 318–328. Honoring the research contributions of Charles R. Nelson.
- Park T and Casella G (2008). The bayesian lasso. *Journal of the American Statistical Association* 103, 681–686.
- Raskutti G and Yuan M (2015, 12). *Convex Regularization for High-Dimensional Tensor Regression*. ArXiv e-prints.
- Schweinberger M, Babkin S, and Ensor KB (2017). High-dimensional multivariate time series with additional structure. *Journal of Computational and Graphical Statistics* 26(3), 610–622.
- Seth AK, Barrett AB, and Barnett L (2015). Granger causality analysis in neuroscience and neuroimaging. *Journal of Neuroscience* 35(8), 3293–3297. [PubMed: 25716830]
- Sims CA (1980). Macroeconomics and reality. *Econometrica: Journal of the Econometric Society* 48(1), 1–48.
- Sparks DK, Khare K, and Ghosh M (2015). Necessary and sufficient conditions for high-dimensional posterior consistency under g -priors. *Bayesian Anal.* 10(3), 627–664.
- Stock JH and Watson MW (2005). An empirical comparison of methods for forecasting using many predictors. Manuscript, Princeton University.
- Zellner A (1962). An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias. *Journal of the American Statistical Association* 57(298), 348–368.

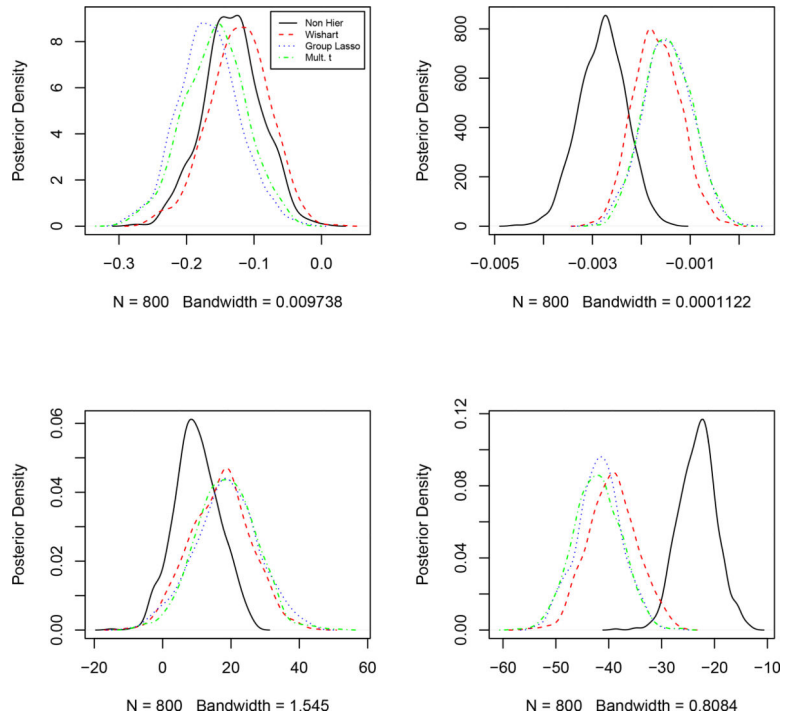


Figure 1: Posterior densities of entries (11), (24), (32) and (42) in \mathbf{A}_1 under 4 different priors

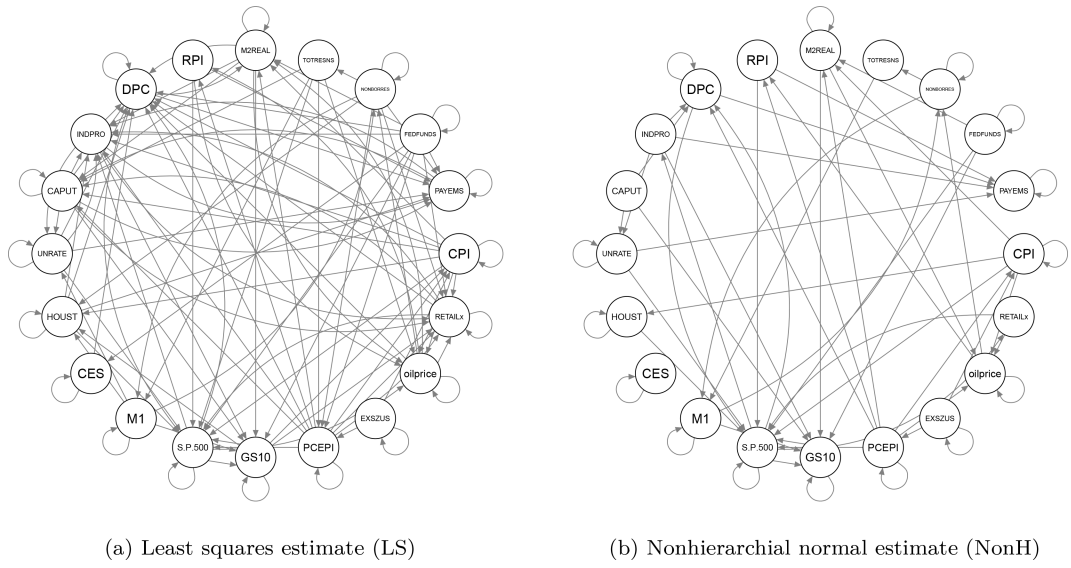


Figure 2:
Network representation of the transition matrix (A_1)

Table 1:

Relative error in VAR ($d = 1, 2$) with $\Sigma_\varepsilon = \sigma^2 I_p$ where % denotes percentage of non-zero entries in Φ_0 .

		Lag $d = 1$						Lag $d = 2$					
		n_1		n_2		n_3		n_1		n_2		n_3	
	%	LS	PM	LS	PM	LS	PM	LS	PM	LS	PM	LS	PM
Small VAR	10	0.93 (0.21)	0.83 (0.12)	0.82 (0.12)	0.70 (0.08)	0.70 (0.07)	0.55 (0.05)	1.26 (0.24)	1.09 (0.17)	1.11 (0.17)	0.96 (0.07)	1.00 (0.10)	0.76 (0.08)
	20	1.06 (0.34)	1.00 (0.24)	1.00 (0.26)	0.87 (0.16)	0.83 (0.16)	0.71 (0.11)	1.39 (0.35)	1.22 (0.30)	1.28 (0.27)	1.07 (0.21)	1.14 (0.18)	0.94 (0.15)
	30	1.23 (0.45)	1.15 (0.37)	1.12 (0.37)	0.97 (0.27)	0.99 (0.27)	0.81 (0.19)	1.57 (0.47)	1.40 (0.40)	1.45 (0.37)	1.28 (0.32)	1.30 (0.34)	1.15 (0.22)
Medium VAR	10	1.81 (0.40)	1.64 (0.21)	1.69 (0.31)	1.53 (0.12)	1.56 (0.23)	1.40 (0.07)	2.45 (0.46)	2.12 (0.32)	2.31 (0.37)	2.01 (0.24)	2.24 (0.31)	1.85 (0.14)
	20	1.95 (0.50)	1.81 (0.32)	1.85 (0.42)	1.70 (0.25)	1.74 (0.35)	1.56 (0.17)	2.60 (0.57)	2.26 (0.42)	2.50 (0.51)	2.15 (0.33)	2.43 (0.40)	2.01 (0.29)
	30	2.10 (0.64)	1.94 (0.43)	1.98 (0.52)	1.84 (0.38)	1.87 (0.45)	1.74 (0.27)	2.74 (0.69)	2.46 (0.54)	2.62 (0.60)	2.31 (0.46)	2.52 (0.56)	2.20 (0.39)
Large VAR	10	2.70 (0.58)	2.47 (0.30)	2.55 (0.50)	2.34 (0.23)	2.46 (0.40)	2.22 (0.16)	3.66 (0.67)	3.18 (0.45)	3.50 (0.57)	3.04 (0.35)	3.38 (0.49)	2.92 (0.28)
	20	2.84 (0.70)	2.62 (0.42)	2.75 (0.60)	2.52 (0.34)	2.60 (0.54)	2.35 (0.25)	3.81 (0.78)	3.33 (0.56)	3.69 (0.70)	3.23 (0.50)	3.54 (0.61)	3.05 (0.40)
	30	3.03 (0.82)	2.78 (0.54)	2.90 (0.70)	2.65 (0.44)	2.71 (0.63)	2.49 (0.35)	3.96 (0.88)	3.49 (0.69)	3.89 (0.84)	3.32 (0.63)	3.78 (0.76)	3.25 (0.51)

Table 2:

Relative error in VAR ($d = 1, 2$) with $\Sigma_\varepsilon = \text{Toeplitz}(\rho = 0.80)$ where % denotes percentage of non-zero entries in Φ_0 .

		Lag $d = 1$						Lag $d = 2$					
		n_1		n_2		n_3		n_1		n_2		n_3	
	%	LS	PM	LS	PM	LS	PM	LS	PM	LS	PM	LS	PM
Small VAR	10	1.03 (0.27)	0.95 (0.18)	0.95 (0.20)	0.82 (0.09)	0.82 (0.10)	0.64 (0.05)	1.45 (0.37)	1.24 (0.27)	1.31 (0.30)	1.08 (0.19)	1.19 (0.18)	0.98 (0.09)
	20	1.20 (0.40)	1.10 (0.30)	1.06 (0.29)	0.94 (0.21)	0.98 (0.25)	0.83 (0.17)	1.60 (0.48)	1.41 (0.37)	1.49 (0.40)	1.30 (0.29)	1.35 (0.31)	1.17 (0.19)
	30	1.36 (0.50)	1.27 (0.41)	1.22 (0.44)	1.11 (0.32)	1.12 (0.32)	0.99 (0.25)	1.76 (0.60)	1.61 (0.51)	1.66 (0.50)	1.45 (0.40)	1.50 (0.46)	1.26 (0.30)
Medium VAR	10	2.02 (0.52)	1.87 (0.34)	1.94 (0.41)	1.72 (0.27)	1.75 (0.37)	1.61 (0.16)	2.80 (0.69)	2.46 (0.50)	2.70 (0.62)	2.35 (0.43)	2.53 (0.54)	2.21 (0.33)
	20	2.20 (0.63)	2.01 (0.45)	2.10 (0.55)	1.87 (0.38)	1.92 (0.45)	1.80 (0.33)	2.95 (0.81)	2.62 (0.62)	2.86 (0.71)	2.52 (0.52)	2.74 (0.65)	2.34 (0.42)
	30	2.36 (0.73)	2.20 (0.56)	2.19 (0.67)	2.06 (0.48)	2.06 (0.56)	1.91 (0.41)	3.12 (0.94)	2.78 (0.72)	2.99 (0.83)	2.67 (0.67)	2.91 (0.77)	2.50 (0.57)
Large VAR	10	3.03 (0.75)	2.80 (0.49)	2.88 (0.67)	2.63 (0.41)	2.83 (0.60)	2.55 (0.33)	4.19 (1.03)	3.67 (0.73)	4.09 (0.97)	3.55 (0.62)	3.93 (0.87)	3.44 (0.58)
	20	3.17 (0.86)	2.94 (0.60)	3.07 (0.79)	2.79 (0.51)	2.89 (0.69)	2.69 (0.45)	4.34 (1.14)	3.85 (0.84)	4.25 (1.06)	3.70 (0.76)	4.11 (1.00)	3.594 (0.66)
	30	3.33 (0.98)	3.10 (0.73)	3.27 (0.90)	3.00 (0.63)	3.08 (0.82)	2.80 (0.56)	4.53 (1.26)	4.00 (0.95)	4.37 (1.17)	3.89 (0.87)	4.26 (1.11)	3.70 (0.79)

Table 3:

Percentage of true positive non-zero entries recovered in $\hat{\Phi}$.

		Lag $d = 1$						Lag $d = 2$						
		$\Sigma_e = \sigma^2 I_p$			$\Sigma_e = \text{Toep}$			$\Sigma_e = \sigma^2 I_p$			$\Sigma_e = \text{Toep}$			
		%	n_1	n_2	n_3	n_1	n_2	n_3	n_1	n_2	n_3	n_1	n_2	n_3
Small VAR	10	85.0	85.0	82.0	85.0	85.0	84.0	84.7	84.5	81.6	84.6	84.6	83.6	
	20	80.0	80.0	81.0	80.0	78.0	80.0	79.7	79.6	80.5	79.5	77.6	79.8	
	30	77.0	75.0	77.0	77.0	73.0	73.0	76.6	74.6	76.5	76.7	72.6	72.6	
Medium VAR	10	89.3	90.0	89.3	89.3	89.8	88.3	88.8	89.5	89.0	88.9	89.4	87.9	
	20	85.3	85.5	85.5	85.3	84.5	85.0	84.9	85.3	85.3	84.8	84.1	84.8	
	30	81.0	80.8	81.3	81.0	79.8	78.8	80.8	80.5	81.0	80.6	79.3	78.3	
Large VAR	10	92.0	92.1	92.2	92.0	92.1	91.8	91.8	91.6	91.8	91.7	91.7	91.4	
	20	88.1	88.1	88.1	88.1	87.9	87.9	87.7	87.8	87.7	87.7	87.5	87.5	
	30	83.1	83.3	83.4	83.1	82.8	82.9	82.7	82.8	82.9	82.7	82.4	82.6	

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 4:

Relative error in VAR ($d = 1, 2$) with $\Sigma_\epsilon = \sigma^2 I_p$, where % denotes percentage of non-zero entries in Φ_0 .

Lag $d = 1$	%	Wishart			Group lasso			Multivariate t		
		n_1	n_2	n_3	n_1	n_2	n_3	n_1	n_2	n_3
Small VAR	10	0.84	0.75	0.65	0.84	0.75	0.64	0.85	0.74	0.58
	20	0.99	0.91	0.73	1.00	0.89	0.79	1.00	0.90	0.76
	30	1.14	1.09	0.92	1.16	1.02	0.94	1.15	1.04	0.93
Medium VAR	10	1.63	1.53	1.42	1.62	1.54	1.43	1.63	1.53	1.39
	20	1.78	1.71	1.53	1.76	1.69	1.56	1.79	1.66	1.55
	30	1.93	1.83	1.70	1.95	1.86	1.72	1.95	1.82	1.69
Large VAR	10	2.39	2.26	2.16	2.41	2.31	2.22	2.42	2.34	2.22
	20	2.59	2.42	2.31	2.58	2.47	2.39	2.59	2.46	2.32
	30	2.77	2.61	2.52	2.70	2.60	2.53	2.74	2.63	2.54
Lag $d = 2$										
Small VAR	10	0.88	0.74	0.60	0.87	0.80	0.71	0.89	0.78	0.63
	20	1.05	0.91	0.83	1.04	0.94	0.85	1.08	0.94	0.86
	30	1.25	1.13	0.97	1.23	1.10	0.99	1.27	1.13	1.03
Medium VAR	10	1.71	1.58	1.46	1.71	1.58	1.53	1.70	1.60	1.50
	20	1.85	1.72	1.60	1.85	1.76	1.71	1.89	1.82	1.66
	30	2.05	1.96	1.77	2.01	1.95	1.82	2.04	1.97	1.83
Large VAR	10	2.52	2.37	2.26	2.51	2.41	2.29	2.53	2.44	2.25
	20	2.69	2.58	2.47	2.67	2.59	2.47	2.69	2.59	2.51
	30	2.88	2.75	2.60	2.86	2.75	2.62	2.87	2.74	2.64

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 5:

Relative error in VAR ($d = 1, 2$) with $\Sigma_e = \text{Toeplitz}(\rho = 0.80)$, where % denotes percentage of non-zero entries in Φ_0 .

Lag $d = 1$	%	Wishart			Group lasso			Multivariate t		
		n_1	n_2	n_3	n_1	n_2	n_3	n_1	n_2	n_3
Small VAR	10	0.90	0.76	0.65	0.90	0.83	0.69	0.89	0.79	0.70
	20	1.06	0.94	0.77	1.03	0.92	0.82	1.07	0.96	0.81
	30	1.23	1.13	0.99	1.21	1.09	0.99	1.23	1.14	0.97
Medium VAR	10	1.75	1.61	1.50	1.73	1.63	1.53	1.74	1.61	1.56
	20	1.92	1.80	1.67	1.92	1.78	1.73	1.91	1.79	1.72
	30	2.08	1.98	1.82	2.02	1.98	1.88	2.07	1.94	1.84
Large VAR	10	2.58	2.49	2.33	2.60	2.50	2.38	2.60	2.49	2.33
	20	2.76	2.65	2.51	2.75	2.65	2.54	2.77	2.62	2.57
	30	2.93	2.81	2.71	2.90	2.81	2.74	2.93	2.82	2.74
Lag $d = 2$										
Small VAR	10	1.16	1.02	0.96	1.15	1.06	0.98	1.18	1.06	0.89
	20	1.31	1.24	1.12	1.31	1.22	1.11	1.34	1.24	1.10
	30	1.49	1.41	1.23	1.49	1.41	1.27	1.49	1.40	1.26
Medium VAR	10	2.25	2.13	2.04	2.25	2.17	2.05	2.27	2.16	2.01
	20	2.43	2.30	2.16	2.41	2.32	2.24	2.45	2.32	2.26
	30	2.59	2.47	2.33	2.59	2.49	2.38	2.59	2.48	2.36
Large VAR	10	3.37	3.24	3.10	3.36	3.25	3.18	3.36	3.24	3.11
	20	3.53	3.40	3.29	3.51	3.43	3.31	3.54	3.46	3.36
	30	3.72	3.62	3.41	3.65	3.57	3.49	3.74	3.59	3.48

Table 6:

Relative error

Estimator	$\Sigma_e = \sigma^2 I_p$	$\Sigma_e = \text{Toeplitz}$
LS	0.604	1.384
Non H	0.583	0.614
W	0.462	0.544
Mult. t	0.430	0.414
GL	0.321	0.362

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 7:

In-sample prediction error

	SMALL ($p = 4$)	MEDIUM ($p = 20$)	LARGE ($p = 134$)
LS	0.847	0.863	0.673
Non-H	0.852	0.864	0.674
W	0.845	0.863	0.675
Mult. t	0.877	0.885	0.685
GL	0.847	0.873	0.674

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 8:

Out of sample relative prediction error

		Uniform prior			Random Walk		
		SMALL	MEDIUM	LARGE	SMALL	MEDIUM	LARGE
		$p = 4$	$p = 20$	$p = 134$	$p = 4$	$p = 20$	$p = 134$
Non Hierarchical	$h = 1$	0.88	0.72	0.62	0.49	0.40	0.33
	$h = 6$	0.90	0.78	0.62	0.43	0.42	0.37
	$h = 12$	0.95	0.84	0.74	0.43	0.41	0.37
Wishart	$h = 1$	0.88	0.81	0.68	0.49	0.45	0.32
	$h = 6$	0.86	0.86	0.68	0.42	0.40	0.36
	$h = 12$	0.93	0.92	0.71	0.43	0.41	0.38
Group Lasso	$h = 1$	0.90	0.86	0.60	0.51	0.45	0.30
	$h = 6$	0.88	0.88	0.63	0.46	0.42	0.34
	$h = 12$	0.93	0.92	0.67	0.47	0.45	0.37
Mult. t	$h = 1$	0.91	0.89	0.71	0.50	0.46	0.31
	$h = 6$	0.87	0.93	0.77	0.49	0.44	0.34
	$h = 12$	0.92	0.94	0.81	0.45	0.41	0.38

Table 9:

Out-of-sample relative prediction error for CPI, PAYEMS and Fed- Fund for the three VAR model specifications considered. The column LARGE BVAR corresponds to the entries of Table III in Bafibura et al. (2010) for a Bayesian VAR model with a normal-inverted Wishart prior distribution and $d=13$ lags, based on the same set of variables, but covering the period 1971–2003.)

		SMALL	MEDIUM	LARGE	LARGE BVAR
	Uniform prior	$p = 4$	$p = 20$	$p = 134$	$p = 134$
$h = 1$	CPI	1.05	0.91	0.44	-
	PAYEMS	1.21	1.04	0.91	-
	FFUND	0.78	0.75	0.68	-
$h = 6$	CPI	1.03	0.97	0.38	-
	PAYEMS	1.08	1.06	0.48	-
	FFUND	0.92	0.82	0.68	-
$h = 12$	CPI	0.98	0.96	0.42	-
	PAYEMS	0.93	0.91	0.73	-
	FFUND	0.92	0.88	0.72	-
Random Walk					
$h = 1$	CPI	0.43	0.41	0.34	0.50
	PAYEMS	0.45	0.43	0.39	0.46
	FFUND	0.50	0.45	0.36	0.75
$h = 6$	CPI	0.38	0.34	0.28	0.40
	PAYEMS	0.53	0.48	0.39	0.50
	FFUND	0.41	0.37	0.36	1.29
$h = 12$	CPI	0.51	0.45	0.42	0.44
	PAYEMS	0.51	0.88	0.73	0.78
	FFUND	0.33	0.31	0.28	1.93