# Sparse Bayesian Vector Autoregressions in Huge Dimensions

Gregor Kastner

Institute for Statistics and Mathematics

Department of Finance, Accounting and Statistics

WU Vienna University of Economics and Business

and

Florian Huber

Institute for Macroeconomics

Department of Economics

WU Vienna University of Economics and Business

June 28, 2018

**Abstract**

We develop a Bayesian vector autoregressive (VAR) model with multivariate stochastic volatility that is capable of handling vast dimensional information sets. Three features are introduced to permit reliable estimation of the model. First, we assume that the reduced-form errors in the VAR feature a factor stochastic volatility structure, allowing for conditional equation-by-equation estimation. Second, we apply recently developed global-local shrinkage priors to the VAR coefficients to cure the curse of dimensionality. Third, we utilize recent innovations to efficiently sample from high-dimensional multivariate Gaussian distributions. This makes simulation-based fully Bayesian inference feasible when the dimensionality is large but the time series length is moderate. We demonstrate the merits of our approach in an extensive simulation study and apply the model to US macroeconomic data to evaluate its forecasting capabilities.

*Keywords:* factor stochastic volatility, curse of dimensionality, shrinkage, Dirichlet-Laplace prior, Normal-Gamma prior, efficient MCMC.

# 1 Introduction

Previous research has identified two important features that macroeconometric models should possess: the ability to exploit high dimensional information sets (Bańbura et al., 2010; Stock and Watson, 2011) and the possibility to capture non-linear features of the underlying time series (Cogley and Sargent, 2002; Primiceri, 2005; Clark, 2011; Clark and Ravazzolo, 2015). While the literature suggests several paths to estimate large models, the majority of such approaches implies that once non-linearities are taken into account, analytical solutions are no longer available and the computational burden becomes prohibitive.[1] This implies that high dimensional non-linear models can practically be estimated only under strong (and often unrealistic) restrictions on the dynamics of the model. However, especially in forecasting applications or in structural analysis, recent literature suggests that successful models should be able to exploit lots of information and also control for breaks in the autoregressive parameters or, more importantly, changes in the volatility of economic shocks (Primiceri, 2005; Sims and Zha, 2006; Koop et al., 2009).

Two reasons limit the use of large (or even huge) non-linear models. The first reason is statistical. Since the number of parameters in a standard vector autoregression rises quadratically with the number of time series included and commonly used macroeconomic time series are rather short, in-sample overfitting proves to be a serious issue. As a solution, the Bayesian literature on VAR modeling (Doan et al., 1984; Litterman, 1986; Sims and Zha, 1998; George et al., 2008; Bańbura et al., 2010; Koop, 2013; Clark, 2011; Clark and Ravazzolo, 2015; Korobilis and Pettenuzzo, 2016; Follett and Yu, 2017; Huber and Feldkircher, 2018) suggests shrinkage priors that push the parameter space towards some stylized prior model like a multivariate random walk. On the other hand, Ahelegbey et al. (2016) suggest to view VARs as graphical models and perform model selection drawing from the literature on sparse directed acyclic graphs. This typically leads to much improved forecasting properties and more meaningful structural inference. Moreover, much of the literature on Bayesian VARs imposes conjugate priors on the autoregressive parameters, allowing for analytical posterior solutions and thus avoiding simulation based techniques like Markov chain Monte Carlo (MCMC). Frequentist approaches often consider multi-step approaches

---

[1]One exception is Koop and Korobilis (2013).

(e.g. Davis et al., 2016).

The second reason is computational. Since non-linear Bayesian models typically have to be estimated by means of MCMC, computational intensity increases vastly if the number of variables included becomes large. The increase in computational complexity stems from the fact that standard algorithms for multivariate regression models call for the inversion of large covariance matrices. Especially for large systems, this can quickly turn prohibitive since the inverse of the posterior variance-covariance matrix on the coefficients has to be computed for each sweep of the MCMC algorithm. For natural conjugate models, this step can be vastly simplified since the likelihood features a convenient Kronecker structure, implying that all equations in the VAR feature the same set of explanatory variables. This speeds up computation by large margins but restricts the flexibility of the model. Carriero et al. (2016), for instance, exploit this fact and introduce a simplified stochastic volatility specification. Another strand of the literature augment each equation of the VAR by including the residuals of the preceding equations (Carriero et al., 2015) which also provides significant improvements in terms of computational speed. Finally, in a recent contribution, Koop et al. (2016) reduce the dimensionality of the problem at hand by randomly compressing the lagged endogenous variables in the VAR.

All papers mentioned hitherto focus on capturing cross-variable correlation in the conditional mean through the VAR part and the co-movement in volatilities is captured by a rich specification of the error variance (Primiceri, 2005) or by a single factor (Carriero et al., 2016). Another strand of the literature, typically used in financial econometrics, utilizes factor models to provide a parsimonious representation of a covariance matrix, focusing exclusively on the second moment of the predictive density. For instance, Pitt and Shephard (1999) and Aguilar and West (2000) assume that the variance-covariance matrix of a broad panel of time series might be described by a lower dimensional matrix of latent static factors featuring stochastic volatility and a variable-specific idiosyncratic stochastic volatility process.

The present paper combines the virtues of exploiting large information sets and allowing for movements in the error variance. The overfitting issue mentioned above is solved as follows. First, we use a Dirichlet-Laplace (DL) prior specification (see Bhattacharya

et al., 2015) on the VAR coefficients. This prior is a global-local shrinkage prior in the spirit of Polson and Scott (2011); it enables us to heavily shrink the parameter space but at the same time provides enough flexibility to allow for non-zero regression coefficients if necessary. Second, a factor stochastic volatility model on the VAR errors grants a parsimonious representation of the time-varying error variance-covariance matrix of the VAR. To deal with the computational complexity, we exploit the fact that, conditionally on the latent factors and their loadings, equation by equation estimation becomes possible within each MCMC iteration. Moreover, we apply recent advances for fast sampling from high dimensional multivariate Gaussian distributions (Bhattacharya et al., 2016) that enables us to estimate models with hundreds of thousands of autoregressive parameters and an error covariance matrix with tens of thousands nontrivial time-varying elements on a quarterly US dataset in a reasonable amount of time. In a careful analysis we show to what extent our proposed method improves upon a set of standard algorithms typically used to simulate from the joint posterior distribution of large dimensional Bayesian VARs.

We first assess the merits of our approach in an extensive simulation study based on a range of different data generating processes. Relative to a set of competing benchmark specifications we show that, in terms of point estimates, the proposed global-local shrinkage prior yields precise parameter estimates and successfully introduces shrinkage in the modeling framework, without overshrinking significant signals.

In an empirical application we use a modified version of the quarterly dataset proposed by Stock and Watson (2011) and McCracken and Ng (2016). To illustrate the out-of-sample performance of our model, we forecast important economic indicators such as output, consumer price inflation and short-term interest rates, amongst others. The proposed model is benchmarked against several alternatives. Our findings suggests that the proposed model performs well, outperforming all benchmarks in terms of one-step-ahead predictive likelihoods. In addition, investigating the time profile of the cumulative log predictive likelihood reveals that allowing for large information sets in combination with the factor structure especially pays off in times of economic stress.

The remainder of this paper is structured as follows. Section 2 introduces the econometric framework. Section 3 details the Bayesian estimation approach, including an elaborated

account of the prior setup adopted and the corresponding conditional posterior distributions. Section 4 provides an analysis of the computational gains of our proposed algorithm relative to a wide set of different algorithms. Section 5 presents the results of an extensive simulation study comparing the performance of a set of carefully selected shrinkage priors for different time series lengths and model dimensions within several (sparse and dense) data generating scenarios. Section 6, after giving a brief overview of the dataset used along with the model specification, illustrates our modeling approach by fitting a single factor model to US data. Moreover, we perform a forecasting exercise to assess the predictive performance of our approach and discuss the choice of the number of latent factors. Finally, the last section concludes.

## 2 Econometric framework

Suppose interest centers on modeling an $m \times 1$ vector of time series denoted by $\boldsymbol{y}_t$ with $t = 1, \ldots, T$. We assume that $\boldsymbol{y}_t$ follows a heteroscedastic VAR($p$) process,[2]

$$\boldsymbol{y}_t = \boldsymbol{A}_1 \boldsymbol{y}_{t-1} + \cdots + \boldsymbol{A}_p \boldsymbol{y}_{t-p} + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}_m(\boldsymbol{0}, \boldsymbol{\Omega}_t). \tag{1}$$

Each $\boldsymbol{A}_j$ ($j = 1, \ldots, p$) is an $m \times m$ matrix of autoregressive coefficients and the error term is assumed to follow a multivariate Gaussian distribution with time-varying variance-covariance matrix $\boldsymbol{\Omega}_t$. To permit reliable and parsimonious estimation when $m$ is large, we decompose the residual covariance matrix into

$$\boldsymbol{\Omega}_t = \boldsymbol{\Lambda} \boldsymbol{V}_t \boldsymbol{\Lambda} + \boldsymbol{\Sigma}_t, \tag{2}$$

where both $\boldsymbol{\Sigma}_t = \mathrm{diag}(\sigma_{1t}^2, \ldots, \sigma_{mt}^2)$ and $\boldsymbol{V}_t = \mathrm{diag}(e^{h_{1t}}, \ldots, e^{h_{qt}})$ are diagonal matrices with dimension $m$ and $q$, respectively, and $\boldsymbol{\Lambda}$ denotes an $m \times q$ matrix with typical element $\lambda_{ij}$ ($i = 1, \ldots, m; j = 1, \ldots, q$). The logarithms of the diagonal elements of $\boldsymbol{\Sigma}_t$ and $\boldsymbol{V}_t$ follow AR(1) processes,

$$h_{jt} = \rho_{hj} h_{j,t-1} + e_{hj,t}, \quad j = 1, \ldots, q, \tag{3}$$

$$\log \sigma_{it}^2 = \mu_{\sigma i} + \rho_{\sigma i}(\log \sigma_{i,t-1}^2 - \mu_{\sigma i}) + e_{\sigma i,t}, \quad i = 1, \ldots, m. \tag{4}$$

---

[2]For simplicity of exposition we omit the intercept term in the following discussion (which we nonetheless include in the empirical application).

To identify the scaling of the elements of $\boldsymbol{\Lambda}$, the process specified in (3) is assumed to have mean zero while $\mu_{\sigma j}$ in (4) is the unconditional mean of the log-elements of $\boldsymbol{\Sigma}_t$ to be estimated from the data (cf. Kastner et al., 2017). The parameters $\rho_{hj}$ and $\rho_{\sigma i}$ are a priori restricted to the interval $(-1, 1)$ and denote the persistences of the latent log variances. The error terms $e_{hj,t}$ and $e_{\sigma i,t}$ denote independent zero mean innovations with variances $\varsigma_{hj}^2$ and $\varsigma_{\sigma i}^2$, respectively. This specification implies that the volatilities are mean reverting and thus bounded in the limit.

This error structure is known as the factor stochastic volatility model (see e.g. Pitt and Shephard, 1999; Aguilar and West, 2000). It can be equivalently written by introducing $q$ conditionally independent latent factors $\boldsymbol{f}_t \sim \mathcal{N}_q(\boldsymbol{0}, \boldsymbol{V}_t)$ and rewriting the error term in (1) as

$$\boldsymbol{\varepsilon}_t = \boldsymbol{\Lambda} \boldsymbol{f}_t + \boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \sim \mathcal{N}_m(\boldsymbol{0}, \boldsymbol{\Sigma}_t). \tag{5}$$

Note that off-diagonal entries of $\boldsymbol{\Omega}_t$ exclusively stem from the volatilities of the $q$ factors while the diagonal entries of $\boldsymbol{\Omega}_t$ are allowed to feature idiosyncratic deviations driven by the elements of $\boldsymbol{\Sigma}_t$. This specification reduces the number of free elements in $\boldsymbol{\Omega}_t$ from $m(m+1)/2$ to $mq$, where the latter quantity is typically much smaller than the former. In addition, by conditioning on the latent factors, this representation enables us to derive an efficient Gibbs sampler that allows for conditional equation-by-equation estimation. As will be discussed in more detail in Section 3.2, this constitutes a key feature for computationally feasible Bayesian inference when the dimensionality $m$ becomes large.

The model described by Eqs. (1) to (2) is related to several alternative specifications commonly used in the literature. For instance, assuming that $\boldsymbol{V}_t = \boldsymbol{I}$ and $\boldsymbol{\Sigma}_t \equiv \boldsymbol{\Sigma}$ for all $t$ leads to the specification adopted in Stock and Watson (2005). Setting $q = 1$ and $\boldsymbol{\Sigma}_t \equiv \boldsymbol{\Sigma}$ yields a specification that is similar to the one stipulated in Carriero et al. (2016), with the difference that our model imposes restrictions on the covariances whereas Carriero et al. (2016) estimate a full (but constant) covariance matrix and our model implies that the stochastic volatility enters $\boldsymbol{\Omega}_t$ in an additive fashion.

Before proceeding to the next subsection it is worth summarizing the key features of the model given by Eqs. (1) to (2). First, we capture cross-variable movements in the conditional mean through the VAR block of the model and assume that co-movement in

conditional variances is captured by the factor model in (5). Second, the model introduces stochastic volatility by assuming that a large panel of volatilities may be efficiently summarized through a set of latent factors. This choice is more flexible than a single factor model for the volatility, effectively providing a parsimonious representation of $\boldsymbol{\Omega}_t$ that is flexible enough to replicate the dynamic behavior of the variances of a broad set of macroeconomic quantities.

# 3   Shrinkage in large dimensional VAR models

Our approach to estimation and inference is Bayesian. This implies that we have to specify suitable prior distributions on the parameters of the model described in the previous subsection and combine the prior distributions with the likelihood to obtain the corresponding posterior distributions.

## 3.1   A global-local shrinkage prior

For prior implementation, it proves to be convenient to define a $k \times 1$ vector of predictors $\boldsymbol{x}_t = (\boldsymbol{y}'_{t-1}, \ldots, \boldsymbol{y}'_{t-p})'$ and an $m \times k$ coefficient matrix $\boldsymbol{B} = (\boldsymbol{A}_1, \ldots, \boldsymbol{A}_p)$ with $k = mp$ to rewrite the model in (1) more compactly as $\boldsymbol{y}_t = \boldsymbol{B}\boldsymbol{x}_t + \boldsymbol{\varepsilon}_t$. Stacking the rows of $\boldsymbol{y}_t$, $\boldsymbol{x}_t$, and $\boldsymbol{\varepsilon}_t$ yields

$$\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{B}' + \boldsymbol{E}, \tag{6}$$

where $\boldsymbol{Y} = (\boldsymbol{y}_1, \ldots, \boldsymbol{y}_T)'$, $\boldsymbol{X} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_T)'$ and $\boldsymbol{E} = (\boldsymbol{\varepsilon}_1, \ldots, \boldsymbol{\varepsilon}_T)'$ denote the corresponding full data matrices.

Typically, the matrix $\boldsymbol{B}$ is a sparse matrix with non-zero elements mainly located on the main diagonal of $\boldsymbol{A}_1$. In fact, existing priors in the Minnesota tradition tend to strongly push the system towards the prior model in high dimensions. However, especially in large models an extremely tight prior on $\boldsymbol{B}$ might lead to severe overshrinking, effectively zeroing out coefficients that might be important to explain $\boldsymbol{y}_t$. If the matrix $\boldsymbol{B}$ is characterized by a relatively low number of non-zero regression coefficients, a possible solution is a global-local shrinkage prior (Polson and Scott, 2011).

A recent variant that falls within the class of global-local shrinkage priors is the Dirichlet-Laplace (DL) prior put forward in Bhattacharya et al. (2015). This prior possesses convenient shrinkage properties in the presence of a large degree of sparsity of the parameter vector $\boldsymbol{b} = \text{vec}(\boldsymbol{B})$. In what follows, we impose the DL prior on each of the $K = mk$ elements of $\boldsymbol{b}$, denoted as $b_j$ for $j = 1, \ldots, K$,

$$b_j \sim \mathcal{DE}(\vartheta_j \zeta) \quad \Leftrightarrow \quad b_j \sim \mathcal{N}(0, \psi_j \vartheta_j^2 \zeta^2), \quad \psi_j \sim \mathcal{E}(1/2), \tag{7}$$

where $\mathcal{DE}$ denotes the double exponential (Laplace) and $\mathcal{E}$ the exponential distribution, $\psi_j$ is an auxiliary scaling parameters to achieve conditional normality, and the elements of $\boldsymbol{\vartheta} = (\vartheta_1, \ldots, \vartheta_K)'$ are local auxiliary scaling parameters that are bounded to the $(K-1)$-dimensional simplex $\mathcal{S}^{K-1} = \{\boldsymbol{\vartheta} : \vartheta_j \geq 0, \sum_{j=1}^{n} \vartheta_j = 1\}$. A natural prior choice for $\vartheta_j$ is the (symmetric) Dirichlet distribution with hyperparameter $a$, $\vartheta_j \sim \mathcal{D}(a, \ldots, a)$. In addition, $\zeta$ is a global shrinkage parameter that pushes all elements in $\boldsymbol{B}$ towards zero and exhibits an important role in determining the tail behavior of the marginal prior distribution on $b_j$, obtained after integrating out the $\vartheta_j$s. Thus, we follow Bhattacharya et al. (2015) and adopt a fully Bayesian approach by specifying a Gamma distributed prior on $\zeta \sim \mathcal{G}(Ka, 1/2)$. It is noteworthy that this prior setup has at least two convenient features that appear to be of prime importance for VAR modeling. First, it exerts a strong degree of shrinkage on all elements of $\boldsymbol{B}$ but still provides additional flexibility such that non-zero regression coefficients are permitted. This critical property is a feature which a large class of global-local shrinkage priors share (Griffin and Brown, 2010; Carvalho et al., 2010; Polson and Scott, 2011) and has been recently adopted in a VAR framework by Huber and Feldkircher (2018) and within the general context of state space models by Bitto and Frühwirth-Schnatter (2018). Second, implementation is simple and requires relatively little additional input from the researcher. In fact, the prior heavily relies on a single structural hyperparameter that has to be specified with care, namely $a$.

The hyperparameter $a$ influences the empirical properties of the proposed shrinkage prior along several important dimensions. Smaller values of $a$ lead to heavy shrinkage on all elements of $\boldsymbol{B}$. To see this, note that lower values of $a$ imply that more prior mass is placed on small values of $\zeta$ a priori. Similarly, when $a$ is small, the Dirichlet prior places more mass on values of $\vartheta_j$ close to zero. Since lower values of $\zeta$ translate into thicker tails

of the marginal prior on $b_j$, the specific choice of $a$ not only influences the overall degree of shrinkage but also the tail behavior of the prior. Bhattacharya et al. (2015) show that if $a$ is specified as $K^{-(1+\Delta)}$ for any $\Delta > 0$ to be small, the DL prior displays excellent posterior contraction rates, and Pati et al. (2014) discuss the shrinkage properties of the proposed prior within the context of factor models.

For the factor loadings we independently use a standard normally distributed prior on each element $\lambda_{ij} \sim \mathcal{N}(0,1)$ for $i = 1, \ldots, m$ and $j = 1, \ldots, q$. Likewise, we impose a normally distributed prior on the mean of the log-volatility $\mu_{\sigma j} \sim \mathcal{N}(0, M_\mu)$ with $M_\mu$ denoting the prior variance. Furthermore, we place the commonly employed Beta distributed prior on the transformed persistence parameter of the log-volatility $\frac{\rho_{sj}+1}{2} \sim \mathcal{B}(a_0, b_0)$ for $s \in \{h, \sigma\}$ and $a_0, b_0 \in \mathbb{R}^+$ to ensure stationarity. Finally, we use a restricted Gamma prior on the innovation variances in Eqs. (3) and (4), $\varsigma_{sj}^2 \sim \mathcal{G}(\frac{1}{2}, \frac{1}{2\xi})$. Here, $\xi$ is a hyperparameter used to control the tightness of the prior. This choice, motivated in Frühwirth-Schnatter and Wagner (2010) implies that if the data is not informative on the degree of time variation of the log volatilities then we do not bound $\varsigma_{sj}^2$ artificially away from zero, effectively applying more shrinkage than the standard inverted Gamma prior.

## 3.2   Full conditional posterior distributions

Conditional on the latent factors and the corresponding loadings, the model in (1) can be cast as a system of $m$ unrelated regression models for the elements in $\boldsymbol{z}_t = \boldsymbol{y}_t - \boldsymbol{\Lambda}\boldsymbol{f}_t$, labeled $z_{it}$, with heteroscedastic errors,

$$z_{it} = \boldsymbol{B}_{i\bullet}\boldsymbol{x}_t + \eta_{it}, \text{ for } i = 1, \ldots, m. \tag{8}$$

Here we let $\boldsymbol{B}_{i\bullet}$ denote the $i$th row of $\boldsymbol{B}$ and $\eta_{it}$ is the $i$th element of $\boldsymbol{\eta}_t$. The corresponding posterior distribution of $\boldsymbol{B}'_{i\bullet}$ is a $k$-dimensional multivariate Gaussian distribution,

$$\boldsymbol{B}'_{i\bullet}|\bullet \sim \mathcal{N}(\boldsymbol{b}_i, \boldsymbol{Q}_i), \tag{9}$$

with $\bullet$ indicating that we condition on the remaining parameters and latent quantities of the model. The posterior variance and mean are given by

$$\boldsymbol{Q}_i = (\tilde{\boldsymbol{X}}_i'\tilde{\boldsymbol{X}}_i + \boldsymbol{\Phi}_i^{-1})^{-1}, \tag{10}$$

$$\boldsymbol{b}_i = \boldsymbol{Q}_i(\tilde{\boldsymbol{X}}_i'\tilde{\boldsymbol{z}}_i). \tag{11}$$

The diagonal prior covariance matrix of the coefficients related to the $i$th equation is given by $\boldsymbol{\Phi}_i$, the respective $k \times k$ diagonal submatrix of $\boldsymbol{\Phi} = \zeta \times \text{diag}(\psi_1\vartheta_1^2, \ldots, \psi_K\vartheta_K^2)$. Moreover, $\tilde{\boldsymbol{X}}_i$ is a $T \times k$ matrix with typical row $t$ given by $\boldsymbol{X}_t/\sigma_{it}$ and $\tilde{\boldsymbol{z}}_i$ is a $T$-dimensional vector with the $t$th element given by $z_{it}/\sigma_{it}$. This normalization renders (8) conditionally homoscedastic with standard normally distributed white noise errors.

The full conditional posterior distribution of $\psi_j$ is inverse Gaussian distributed,

$$\psi_j|\bullet \sim iG(\vartheta_j\zeta/|b_j|, 1) \text{ for } j = 1, \ldots, K. \tag{12}$$

For the global shrinkage parameter $\zeta$ the conditional posterior follows a generalized inverted Gaussian (GIG) distribution,

$$\zeta|\bullet \sim \mathcal{GIG}\left(K(a-1), 1, 2\sum_{j=1}^{K}|b_j|/\vartheta_j\right). \tag{13}$$

To draw from this distribution, we use the R-package `GIGrvg` (Leydold and Hörmann, 2017) implementing the efficient algorithm of Hörmann and Leydold (2013). Moreover, we sample the scaling parameters $\vartheta_j$ by first sampling $L_j$ from $L_j|\bullet \sim \mathcal{GIG}(a-1, 1, 2|b_j|)$, and then setting $\vartheta_j = L_j/\sum_{i=1}^{K} L_i$.

The conditional posterior distributions of the factors are Gaussian and thus straightforward to draw from. The factor loadings are sampled using "deep interweaving" (see Kastner et al., 2017), and the parameters in (3) and (4) along the full histories of the latent log-volatilities are sampled as in Kastner and Frühwirth-Schnatter (2014) using the R-packages `factorstochvol` (Kastner, 2017) and `stochvol` (Kastner, 2016).

Our MCMC algorithm iteratively draws from the conditional posterior distributions outlined above and discards the first $J$ draws as burn-in. In terms of computational requirements, the single most intensive step is the simulation from the joint posterior of the autoregressive coefficients in $\boldsymbol{B}$. Because this step is implemented on an equation-by-equation basis, speed improvements relative to the standard approach are already quite

substantial. However, note that if $k$ is large (i.e. of the order of several thousands), even the commonly employed equation-by-equation sampling fails to deliver a sufficient amount of draws within a reasonable time window. Consequently, we outline an alternative algorithm to draw from a high-dimensional multivariate Gaussian distribution under a Bayesian prior that features a diagonal prior variance-covariance matrix in the upcoming section.

# 4   Computational aspects

The typical approach to sampling from (9) is based on the full system and simultaneously samples from the full conditional posterior of $\boldsymbol{B}$, implying that the corresponding posterior distribution is a $K$-dimensional Gaussian distribution with a $K \times K$ dimensional variance-covariance matrix. Under a non-conjugate prior the computational difficulties arise from the need to invert the $K \times K$ variance-covariance matrix which requires operations of order $O(m^6 p^3)$ under Gaussian elimination.

If a conjugate prior in combination with a constant (or vastly simplified heteroscedastic, see Carriero et al., 2015) specification of $\boldsymbol{\Omega}_t$ is used, the corresponding variance-covariance features a Kronecker structure which is computationally cheaper to invert and scales better in large dimensions. Specifically, the manipulations of the corresponding covariance matrix are of order $O(m^3 + k^3)$, a significant gain relatively to the standard approach. However, this comes at a cost since all equations have to feature the same set of variables, the prior on the VAR coefficients has to be symmetric and any stochastic volatility specification is necessary overly simplistic to preserve conditional conjugacy.

By contrast, recent studies emphasize the computational gains that arise from utilizing a framework that is based on equation-by-equation estimation. Carriero et al. (2015) and Koop et al. (2016) augment each equation of the system by either contemporaneous values of the endogenous variables of the preceding equations or the residuals from the previous equations. Here, our approach renders the equations of the system conditionally independent by conditioning on the factors. From a computational perspective, the differences between using a factor model to disentangle the equations and an approach based on augmenting specific equations by quantities that aim to approximate covariance parameters are negligible. If we sample from (9) directly, the computations involved are of order

11

$O(mk^3) = O(m^4p^3)$. This already poses significant improvements relative to full system estimation.

One contribution of the present paper is the application of the algorithm proposed by Bhattacharya et al. (2016) and developed for univariate regression models under a global-local shrinkage prior. This algorithm is applied to each equation in the system and cycles through the following steps:

1. Sample independently $\boldsymbol{u}_i \sim \mathcal{N}(\mathbf{0}_k, \boldsymbol{\Phi}_i)$ and $\boldsymbol{\delta}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I}_T)$.

2. Use $\boldsymbol{u}_i$ and $\boldsymbol{\delta}_i$ to construct $\boldsymbol{v}_i = \tilde{\boldsymbol{X}}_i \boldsymbol{u}_i + \boldsymbol{\delta}_i$.

3. Solve $(\tilde{\boldsymbol{X}}_i \boldsymbol{\Phi}_i \tilde{\boldsymbol{X}}_i' + \boldsymbol{I}_T) \boldsymbol{w}_i = (\tilde{\boldsymbol{z}}_i - \boldsymbol{v}_i)$ for $\boldsymbol{w}_i$.

4. Set $\boldsymbol{B}_{i\bullet}' = \boldsymbol{u}_i + \boldsymbol{\Phi}_i \tilde{\boldsymbol{X}}_i' \boldsymbol{w}_i$.

This algorithm outperforms all competing variants discussed previously in situations where $k \gg T$, a situation commonly encountered when dealing with large VAR models. In such cases, steps (1) to (4) can be carried out using $O(pm^2T^2)$ floating point operations. In situations where $k \approx T$, the computational advantages relative to the standard equation-by-equation algorithm mentioned above are modest or even negative. However, note that the cost is quadratic in $m$ and linear in $p$ and thus scales much better when the number of endogenous variables and/or lags thereof is increased. More information on the empirical performance of our algorithm can be found in Section 6.4.

# 5 Simulation Study

This section aims at comparing the performance of the DL prior with a range of commonly used alternatives. We investigate sparse, intermediate, and dense data generating processes (DGPs) where $T \in \{50, 100, 150, 200, 250\}$ and $m \in \{10, 20, 50, 100\}$. The probability of an off-diagonal entry to be non-zero is 0.01, 0.1, and 0.8 in each of the respective scenarios. In all scenarios, each intercept entry has a 0.1 probability of being non-zero and all diagonal elements are non-zero with probability 0.8. The non-zero elements are randomly generated from Gaussian distributions roughly tuned to yield stable VARs. More concretely, both the mean $\mu_I$ and the standard deviation $\sigma_I$ of the intercept are set to 0.01, whereas mean

and standard deviation of the diagonal ($D$) and the off-diagonal ($O$) elements are chosen as follows:

- Dense: 80% offdiagonal density level, and $\mu_D = \sigma_D = 0.15$ and $\mu_O = \sigma_O = 0.01$.
- Intermediate: 10% offdiagonal density level, and $\mu_D = \sigma_D = 0.15$ and $\mu_O = \sigma_O = 0.1$.
- Sparse: 1% offdiagonal density level, and $\mu_D = \sigma_D = \mu_O = \sigma_O = 0.3$.

Concerning the errors, we use a single factor SV specification. The factor loadings are generated from $\mathcal{N}(0.001, 0.001^2)$ to roughly match the above scaling. The AR(1) processes driving the idiosyncratic log-variances are assumed to have mean $\mu_{\sigma i} = -12$ with persistences $\rho_{\sigma i}$ ranging from 0.85 to 0.98 and innovation standard deviations $\varsigma_{\sigma i}$ from 0.3 to 0.1. The process driving the factor log variance is assumed to be highly persistent with $\rho_{h1} = 0.99$ and $\varsigma_{h1} = 0.1$. For each of the 60 settings, we simulate 10 data sets. For each of these, we run our MCMC algorithm to obtain 2000 posterior draws after a burn-in of 1000. Consequently, the posterior means are compared to the true values and root mean squared errors are computed. Finally, the median of each of these is reported in Table 1. Alongside the DL prior with weak ($a_{DL} = a = 1/2$) and strong ($a_{DL} = a = 1/K$) shrinkage, we also consider the Normal-Gamma (NG) prior with a single global shrinkage parameter (see Huber and Feldkircher, 2018, for the exact specification) and a standard conjugate Minnesota prior with a single shrinkage parameter $a_M$, implemented by using dummy observations. For the NG prior we specify the prior on the global shrinkage parameter to induce heavy shrinkage (i.e. by setting both hyperparameters of the Gamma prior equal to 0.01) and the prior controlling the excess kurtosis $a_{NG}$ is set equal to 1, corresponding to the Bayesian Lasso (see Park and Casella, 2008), and $a_{NG} = 0.1$. The latter choice places significant prior mass around zero but at the same time leads to a heavy tailed marginal prior. Finally, we report RMSEs of the OLS estimator (if it exists).

As is to be expected, Table 1 reveals strong to severe overfitting of OLS (corresponding to a flat prior), which can be mitigated to a certain extent when the Minnesota prior with strong shrinkage ($a_M = 0.001$) is employed instead. Similarly, the DL prior with weak shrinkage ($a_{DL} = 1/2$) displays a tendency to overfit, in particular when $m$ is large. By contrast, the DL prior with $a_{DL} = 1/K$ performs superior in the medium and high-dimensional settings, leading to a sparse solution while reliably detecting non-zero regression coefficients.

Table 1: Median RMSEs relative to those under the DL prior with strong shrinkage ($a_{DL} = 1/K$). Results stem from 10 simulations per setting, where numbers smaller than one mean that the corresponding prior outperforms the DL($1/K$) prior for a given scenario. Shading:

| 0.25 | 0.33 | 0.44 | 0.57 | 0.76 | 1.00 | 1.32 | 1.74 | 2.30 | 3.03 | 4.00 |

| $m$ / $T$ | sparse | | | | intermediate | | | | dense | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 10 | 20 | 50 | 100 | 10 | 20 | 50 | 100 | 10 | 20 | 50 | 100 |
| **DL ($a_{DL} = 1/2$)** | | | | | | | | | | | | |
| 50 | 1.339 | 1.810 | 2.312 | 2.652 | 1.208 | 1.616 | 1.808 | 1.989 | 1.395 | 1.936 | 2.657 | 3.141 |
| 100 | 1.308 | 1.671 | 2.106 | 2.549 | 1.145 | 1.326 | 1.485 | 1.647 | 1.325 | 1.695 | 2.182 | 2.669 |
| 150 | 1.240 | 1.649 | 1.994 | 2.455 | 1.055 | 1.293 | 1.429 | 1.491 | 1.253 | 1.703 | 2.079 | 2.359 |
| 200 | 1.270 | 1.593 | 2.024 | 2.346 | 1.072 | 1.212 | 1.313 | 1.429 | 1.285 | 1.534 | 1.879 | 2.182 |
| 250 | 1.190 | 1.599 | 2.063 | 2.401 | 1.069 | 1.176 | 1.295 | 1.388 | 1.224 | 1.559 | 1.839 | 2.073 |
| **NG ($a_{NG} = 1$)** | | | | | | | | | | | | |
| 50 | 1.070 | 1.117 | 1.207 | 1.263 | 1.005 | 1.028 | 1.012 | 1.037 | 1.081 | 0.990 | 0.929 | 0.931 |
| 100 | 1.188 | 1.259 | 1.374 | 1.448 | 1.025 | 1.060 | 1.076 | 1.069 | 1.129 | 1.126 | 1.040 | 0.995 |
| 150 | 1.233 | 1.314 | 1.367 | 1.471 | 1.059 | 1.097 | 1.122 | 1.106 | 1.146 | 1.198 | 1.134 | 1.032 |
| 200 | 1.219 | 1.334 | 1.460 | 1.528 | 1.069 | 1.117 | 1.126 | 1.133 | 1.191 | 1.196 | 1.130 | 1.068 |
| 250 | 1.177 | 1.363 | 1.498 | 1.608 | 1.070 | 1.098 | 1.136 | 1.168 | 1.207 | 1.267 | 1.156 | 1.075 |
| **NG ($a_{NG} = 0.1$)** | | | | | | | | | | | | |
| 50 | 1.100 | 1.448 | 1.949 | 2.701 | 0.985 | 1.314 | 1.560 | 2.037 | 1.073 | 1.496 | 2.203 | 3.208 |
| 100 | 1.057 | 1.296 | 1.726 | 2.157 | 1.006 | 1.133 | 1.301 | 1.483 | 1.091 | 1.327 | 1.806 | 2.271 |
| 150 | 1.030 | 1.348 | 1.645 | 2.042 | 0.934 | 1.147 | 1.268 | 1.318 | 1.070 | 1.423 | 1.722 | 1.988 |
| 200 | 1.058 | 1.301 | 1.653 | 1.940 | 0.923 | 1.064 | 1.168 | 1.260 | 1.023 | 1.282 | 1.569 | 1.860 |
| 250 | 0.969 | 1.322 | 1.680 | 1.987 | 0.945 | 1.029 | 1.146 | 1.231 | 1.056 | 1.288 | 1.548 | 1.754 |
| **Minnesota ($a_M = 0.001$)** | | | | | | | | | | | | |
| 50 | 2.303 | 3.161 | 4.458 | 3.158 | 1.945 | 2.688 | 3.181 | 2.273 | 2.367 | 3.416 | 5.023 | 3.671 |
| 100 | 2.122 | 3.119 | 4.243 | 5.936 | 1.798 | 2.311 | 2.771 | 3.244 | 2.199 | 3.170 | 4.396 | 5.691 |
| 150 | 2.266 | 2.951 | 3.873 | 5.448 | 1.768 | 2.129 | 2.551 | 2.896 | 2.281 | 2.915 | 4.006 | 4.938 |
| 200 | 2.189 | 2.858 | 3.884 | 4.978 | 1.695 | 2.004 | 2.305 | 2.697 | 2.187 | 2.731 | 3.568 | 4.495 |
| 250 | 1.824 | 2.912 | 3.948 | 4.985 | 1.621 | 1.966 | 2.255 | 2.570 | 1.942 | 2.808 | 3.486 | 4.013 |
| **Minnesota ($a_M = 0.0001$)** | | | | | | | | | | | | |
| 50 | 1.142 | 1.186 | 1.309 | 1.383 | 1.062 | 1.035 | 1.022 | 0.996 | 1.125 | 0.892 | 0.863 | 0.743 |
| 100 | 1.466 | 1.492 | 1.760 | 1.963 | 1.307 | 1.151 | 1.188 | 1.186 | 1.382 | 1.093 | 1.043 | 0.966 |
| 150 | 1.764 | 1.731 | 1.992 | 2.328 | 1.526 | 1.272 | 1.331 | 1.304 | 1.601 | 1.230 | 1.226 | 1.036 |
| 200 | 1.909 | 1.909 | 2.265 | 2.591 | 1.596 | 1.381 | 1.374 | 1.429 | 1.692 | 1.297 | 1.250 | 1.138 |
| 250 | 1.825 | 2.122 | 2.548 | 2.930 | 1.625 | 1.413 | 1.478 | 1.527 | 1.689 | 1.464 | 1.366 | 1.188 |
| **OLS (if exists)** | | | | | | | | | | | | |
| 50 | 2.706 | 4.726 | DNE | DNE | 2.286 | 4.050 | DNE | DNE | 2.892 | 5.038 | DNE | DNE |
| 100 | 2.454 | 3.799 | 6.141 | DNE | 2.047 | 2.799 | 3.975 | DNE | 2.572 | 3.878 | 6.238 | DNE |
| 150 | 2.540 | 3.497 | 4.961 | 8.279 | 2.020 | 2.525 | 3.217 | 4.780 | 2.560 | 3.519 | 5.126 | 8.101 |
| 200 | 2.513 | 3.143 | 4.520 | 6.350 | 1.946 | 2.239 | 2.690 | 3.688 | 2.499 | 3.031 | 4.209 | 6.001 |
| 250 | 2.055 | 3.294 | 4.726 | 6.495 | 1.808 | 2.171 | 2.717 | 3.269 | 2.143 | 3.133 | 4.161 | 5.434 |

Turning towards the NG prior with $a_{NG} = 1$ we tend observe slightly inferior overall performance due to overshrinking, particularly in the sparse setting. Choosing $a$ more extreme ($a_{NG} = 0.1$) is a good choice when the dimensionality is small whereas for $m \gtrsim 20$ this leads to posterior instability stemming from the extreme prior excess kurtosis. The Minnesota prior with $a_M = 0.0001$ yields an extreme degree of shrinkage, translating into estimates of autoregressive coefficients that are very close to zero, irrespectively of the contribution from the likelihood. In that sense, it overshrinks most of the nonzero coefficients. Nevertheless, in scenarios with extremely low signal-to-noise ratios (such as the dense scenario with $T = 50$ and $m = 100$), this can be beneficial for the overall performance. For further illustration, we showcase four exemplary scenarios in Figures 5 to 8 to be found in the Appendix.

# 6    Empirical forecasting application

In Section 6.1 we first summarize the data set adopted and outline specification choices made. The section that follows (Section 6.2) estimates a simple one factor model to outline the virtues of our proposed framework. Section 6.3 presents the main findings of our forecasting exercise and discusses the choice of the number of factors used for modeling the error covariance structure.

## 6.1    Data, model specification and selection issues

In the empirical application, we forecast a set of key US macroeconomic quantities. To this end, we use the quarterly dataset provided by McCracken and Ng (2016), a variant of the well-known Stock and Watson (2011) dataset for the US. The data spans the period ranging from 1959:Q1 to 2015:Q4. We include $m = 215$ quarterly time series, capturing information on 14 important segments of the economy and follow McCracken and Ng (2016) in transforming the data to be approximately stationary. Furthermore, we standardize each component series to have zero mean and variance one. In the empirical examples

we include $p = 1$ lags of the endogenous variables.[3] The hyperparameters are chosen as follows: $M_\mu = 10$, $a_0 = 20$, $b_0 = 1.5$, $\xi = 1$.

## 6.2  Some empirical key features of the model

To provide some intuition on how our modeling approach works in practice, we start by estimating a simple one factor model (i.e. $q = 1$) and investigate several features of our empirical model. In the next section we will perform an extensive forecasting exercise and discuss the optimal number of factors in terms of forecasting accuracy.

We start by inspecting the posterior distribution of $\boldsymbol{\Lambda}$ and assess what variables appear to load heavily on the latent factor. It is worth emphasizing that most quantities[4] associated with real activity (i.e. industrial production and its components, GDP growth, employment measures) load heavily on the factor. Moreover, expectation measures, housing markets, equity prices and spreads also load heavily on the joint factor.

To assess whether spikes in the volatility associated with the factor coincide with major economic events, the bottom panel of Figure 1 depicts the evolution of the posterior distribution of factor volatility over time. A few findings are worth mentioning. First, volatility spikes sharply during the midst of the 1970s, a period characterized by the first oil price shock and the bankruptcy of Franklin National Bank in 1974. After declining markedly during the second half of the 1970s, the shift in US monetary policy towards aggressively fighting inflation and the second oil price shock again translate into higher macroeconomic uncertainty. Note that from the mid 1980s onward, we observe a general decline in macroeconomic volatility that lasts until the beginning of the 1990s. There we observe a slight increase in volatility possibly caused by the events surrounding the first gulf war. The remaining years up to the beginning of the 2000s has been relatively unspectacular, with volatility levels being muted most of the time. In 2000/2001, volatility again increases due to the burst of the dot-com bubble and the 9/11 attacks. Finally, we observe marked spikes in volatility during recessionary episodes like the recent financial crisis in 2008.

Finally, we assess how well the DL prior with $a = 1/K$ performs in shrinking the

---

[3]We have also experimented with higher lag orders but found only limited evidence of signals beyond lag one for the dataset at hand. Moreover, also out-of-sample predictive studies favored one lag only.

[4]Hereby we refer to the one-step-ahead forecast error related to a given time series.
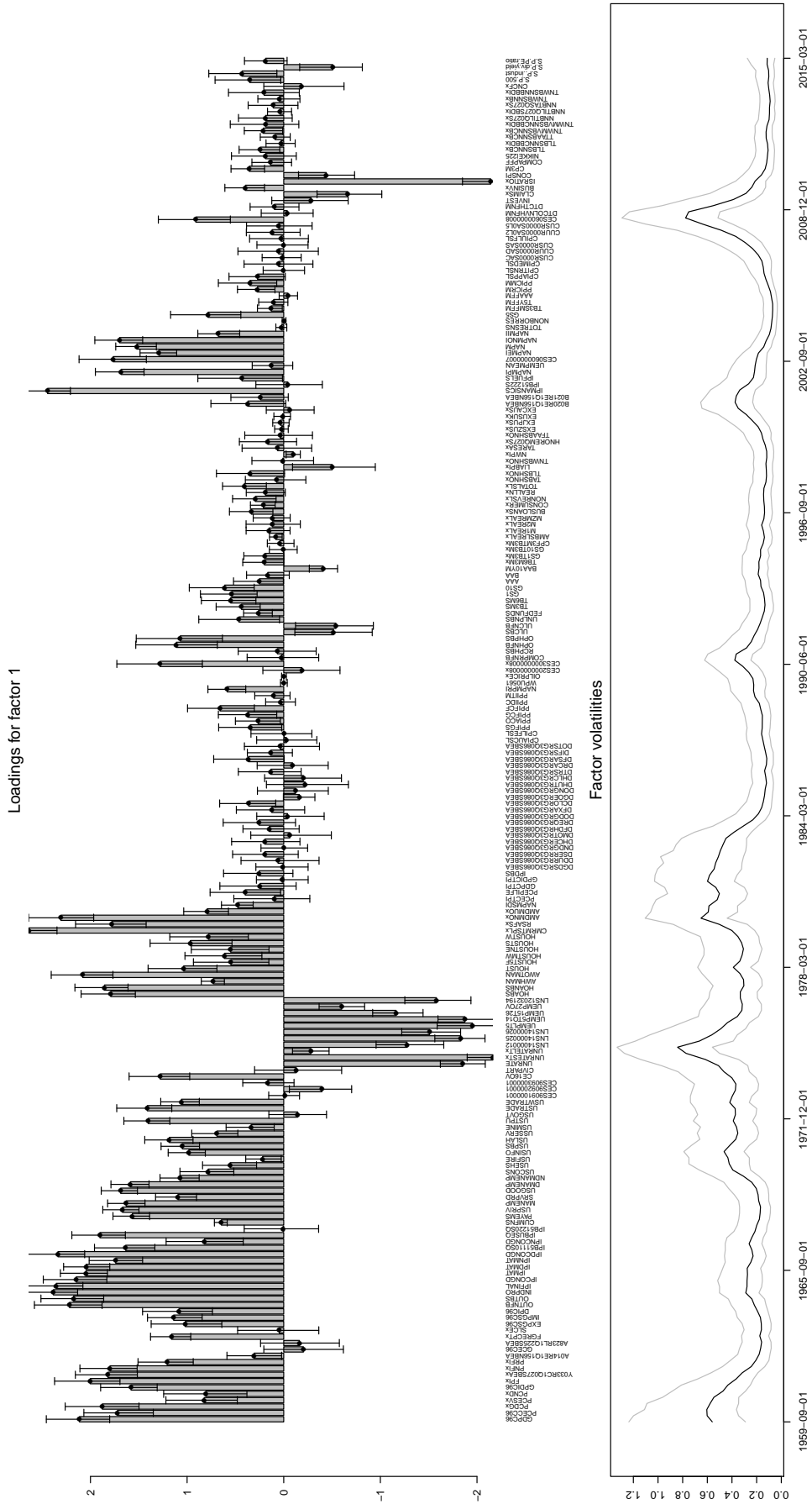
Figure 1: 5th, 50th, and 95th posterior percentiles of factor loadings (upper panel) and factor volatility (lower panel).

coefficients in $\boldsymbol{B}$ to zero. The top panel of Figure 2 depicts a heatmap that gives a rough feeling on the size of each regression coefficient based on the posterior median of $\boldsymbol{B}$. The bottom panel of Figure 2 depicts the posterior interquartile range, providing some evidence on posterior uncertainty.[5] The DL prior apparently succeeds in shrinking the vast majority of the approximately 50 000 coefficients towards zero.

Interestingly, for selected time series measuring inflation (both consumer and producer price inflation) we find that lags of monetary aggregates are allowed to load on the respective inflation series. This result points towards a big advantage of our proposed prior relative to standard VAR priors in the Minnesota tradition: while these priors have been shown to work relatively well in huge dimensions (see Bańbura et al., 2010), they also display a tendency to overshrink when the overall tightness of the prior is integrated out in a Bayesian framework, effectively pushing the posterior distribution of $\boldsymbol{B}$ towards the prior mean and thus ruling out patterns observed under the DL prior.

Inspection of the interquartile range also indicates that the proposed shrinkage prior succeeds in reducing posterior uncertainty markedly. Note that the pattern found for the posterior median of $\boldsymbol{B}$ can also be found in terms of the posterior dispersion. We again observe that the coefficients associated with the first, own lag of a given variable are allowed to be non-zero whereas in most other cases the associated posterior is strongly concentrated around zero. For comparison, we provide heatmaps of posterior medians and interquartile ranges for $a = 1/2$ in the Appendix, see Figure 9.

## 6.3  Predictive evidence

We focus on forecasting gross domestic product (GDPC96), industrial production (IN-PRO), total nonfarm payroll (PAYEMS), civilian unemployment rate (UNRATE), new privately owned housing units started (HOUST), consumer price index inflation (CPI-AUCSL), producer price index for finished goods inflation (PPIFGS), effective federal funds rate (FEDFUNDS), 10-year treasury constant maturity rate (GS10), U.S./U.K. exchange rate (EXUSUKx), and the S&P 500 (S.P.500). This choice includes the variables investi-

---

[5]Since the corresponding posterior distribution is quite heavy-tailed, using posterior standard deviations, while providing a qualitatively similar picture, tend to be slightly exaggerated.

Figure 2: Posterior medians (top) and posterior interquartile ranges (bottom) of VAR coefficients, $a = 1/K = 1/46440$.

gates by Koop et al. (2016) and some additional important macroeconomic indicators that are commonly monitored by practitioners, resulting in a total of eleven series.

To assess the forecasting performance of our model, we conduct a pseudo out-of-sample forecasting exercise with initial estimation sample ranging from 1959:Q3 to 1990:Q2. Based on this estimation period, we compute one-quarter-ahead predictive densities for the first period in the hold-out (i.e. 1990:Q3). After obtaining the corresponding predictive densities and evaluating the corresponding log predictive likelihoods, we expand the estimation period and re-estimate the model. This procedure is repeated 100 times until the final point of the full sample is reached. The quarterly scores obtained this way are then accumulated.

Our model with factors $q \in \{0, 1 \ldots, 4\}$ is benchmarked against the prior model, a pure factor stochastic volatility (FSV) model with conditional mean equal to zero (i.e. $\boldsymbol{B} = \boldsymbol{0}_{m \times k}$). In what follows we label this specification FSV 0. To assess the merits of the proposed shrinkage prior vis-á-vis a Minnesota prior and a NG shrinkage prior we also include the models described in Section 5. Moreover, we include two models that impose the restriction that $\boldsymbol{A}_1 = \boldsymbol{I}_m$ and $\boldsymbol{A}_1 = 0.8 \times \boldsymbol{I}_m$ while $\boldsymbol{A}_j$ for $j > 1$ are set equal to zero matrices in both cases. The first model, labeled FSV 1, assumes that the conditional mean of $\boldsymbol{y}_t$ follows a random walk process and the second specification, denoted as FSV 0.8, imposes the restriction that the variables in $\boldsymbol{y}_t$ feature a rather strong degree of persistence but are stationary. The exercise serves to evaluate whether it pays off to impose a VAR structure on the first moment of the joint density of our data and to assess how many factors are needed to obtain precise multivariate density predictions for our eleven variables of interest.

Overall log predictive scores are summarized in Table 2. An immediate finding is that ignoring the error covariance structure (using zero factors) produces rather inaccurate forecasts for all models considered. While a single factor model improves predictive accuracy by a large margin, allowing for more factors (i.e. even more flexible modeling of the covariance structure) further increases the forecasting performance. For this specific exercise, we identify two factors to be a reasonable choice for most models when the joint log predictive scores of the aforementioned variables are considered. We would like to stress that this choice critically depends on the number of variables we include in our prediction set. If we focus attention on the marginal predictive densities (i.e. the univariate predictive densities

Table 2: Overall log predictive scores for the number of factors $q \in \{0, 1, \ldots, 4\}$ for the VAR-FSVs as well as the competing FSV models. Larger numbers indicate better joint predictive density performance for 11 variables of interest.

|  | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| VAR-FSV DL(1/2) | -1482 | -1353 | -1310 | -1240 | -1210 |
| VAR-FSV DL($1/K$) | -1049 | -992 | -912 | -926 | -935 |
| VAR-FSV NG(1) | -1068 | -1034 | -976 | -977 | -978 |
| VAR-FSV NG(0.1) | -1319 | -1207 | -1134 | -1106 | -1112 |
| VAR-FSV Min(0.01) | -1101 | -1039 | -984 | -1007 | -1018 |
| VAR-FSV Min(0.0001) | -1215 | -1121 | -1103 | -1080 | -1090 |
| FSV 0 | -1208 | -1108 | -1102 | -1066 | -1061 |
| FSV 0.8 | -1178 | -1132 | -1106 | -1087 | -1096 |
| FSV 1 | -1171 | -1128 | -1095 | -1095 | -1107 |

obtained after integrating out the remaining elements in $\boldsymbol{y}_t$) we find that fewer or even no factors receive more support (see Table 3), whereas in the case of higher dimensional prediction sets more factors lead to more accurate density predictions.

Considering forecasting accuracy across models reveals that our proposed VAR-FSV with a DL prior (with $a_{DL} = 1/K$) displays excellent forecasting capabilities, outperforming all competitors by rather large margins. Notice that the Bayesian Lasso (i.e. VAR-FSV NG(1)) also performs rather well, irrespective of the fact that it tends to overshrink significant signals. A similar finding carries over to the VAR coupled with a Minnesota prior with $a_M = 0.01$. This simplistic prior specification also tends to perform well (conditional on setting $q = 2$).

Comparing the differences between the benchmark FSV 0 model and the VAR models considered, we find that explicitly modeling the conditional mean improves the forecasting accuracy in most cases and conditional on including at least a single factor. We conjecture that the inclusion of the factors has pronounced effects on the estimates of $\boldsymbol{B}$ and this could potentially lead to an increase in predictive accuracy. Intuitively, this might be explained by the fact that if the underlying data generating process suggests that $\boldsymbol{B}$ is time-varying, inclusion of the factors might alleviate issues associated with model misspecification by, at least to a certain extent, controlling for structural breaks in $\boldsymbol{B}$.

To investigate whether forecasting performance is homogenous over time, Figure 3 vi-
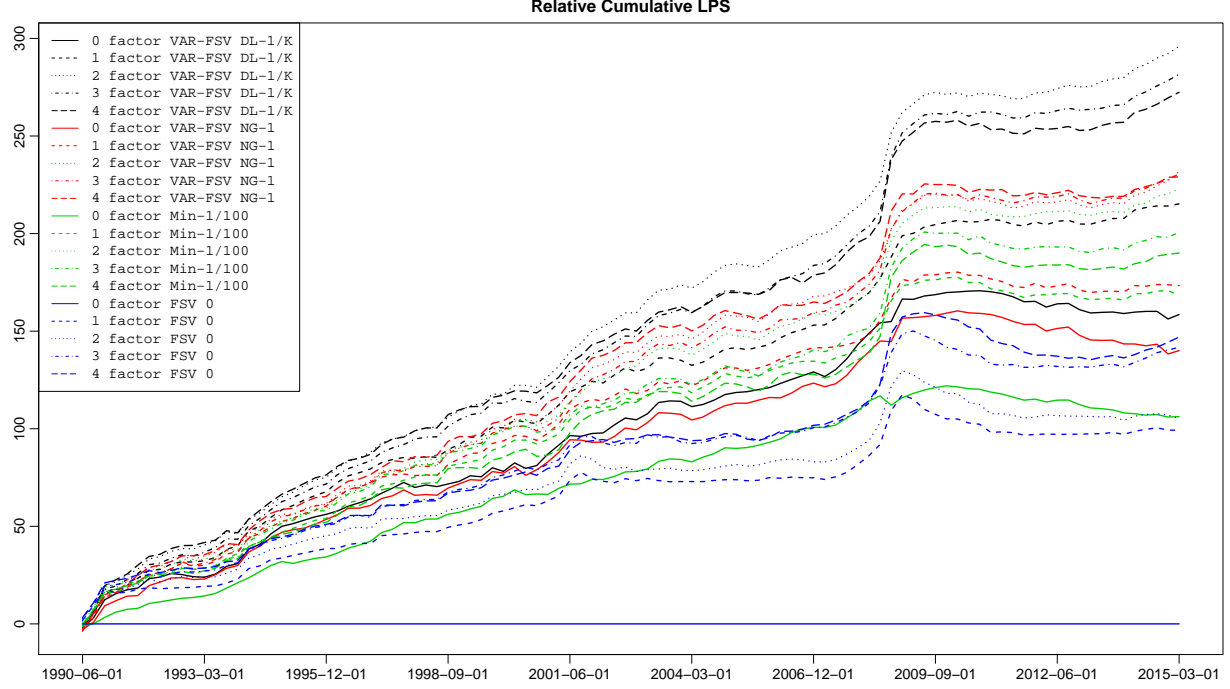
**Relative Cumulative LPS**



Figure 3: Cumulative log predictive scores, relative to a zero-mean model with independent stochastic volatility components for all component series. Higher values correspond to better one-quarter-ahead density predictions up to the corresponding point in time.

sualizes the cumulative log predictive scores (LPS) relative to the zero-factor FSV model over time. The benefit of the flexible SV structure in the VAR residuals is particularly pronounced during the 2008 financial crisis which can be seen by comparing the solid lines to the broken lines. During this period, time-varying covariance modeling appears to be of great importance and the performance of models that ignore contemporaneous dependence deteriorates. This finding is in line with, e.g., Kastner (2018), who reports analogous results for US asset returns. The increase in predictive accuracy can be traced back to the fact that within an economic downturn, the correlation structure of our dataset changes markedly, with most indicators that measure real activity sharply declining in lockstep. A model that takes contemporaneous cross-variable linkages seriously is thus able to fully exploit such behavior which in turn improves predictions.

Up to this point, we focused exclusively on the joint performance of our model for the specific set of variables considered. To gain a deeper understand on how our model performs for relevant selected quantities, Table 3 displays marginal LPSs for the three top

Table 3: Univariate log predictive scores for inflation (CPIAUCSL), short-term interest rates (FEDFUNDS), and output growth (GDPC96), with $q \in \{0, 1, 2\}$ factors.

| | CPIAUCSL | | | FEDFUNDS | | | GDPC96 | | |
|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 0 | 1 | 2 | 0 | 1 | 2 |
| VAR-FSV DL(1/$K$) | -101 | -110 | -108 | -127 | -128 | -121 | -3 | -7 | -9 |
| VAR-FSV NG(1) | -100 | -111 | -112 | -129 | -127 | -126 | -23 | -25 | -28 |
| VAR-FSV Min(0.01) | -101 | -118 | -118 | -127 | -127 | -127 | -25 | -35 | -39 |

performing models and variables. The variables we consider are inflation (CPIAUCSL), short-term interest rates (FEDFUNDS), and output growth (GDPC96).

Compared to the findings based on joint LPS, we observe that models without a factor structure perform better than models that set $q > 0$. This finding corroborates our conjecture stated above, implying that if the set of focus variables is subsequently enlarged, more factors are necessary in order to obtain precise density predictions. Here, we only focus on marginal model performance, implying that for each variable, contemporaneous relations between the elements in $\boldsymbol{y}_t$ are integrated out. This, in turn, implies that the additional gain in model flexibility is offset by the comparatively larger number of parameters. This finding holds true for inflation and output growth. For short-term interest rates, a different pattern emerges. Considering the marginal LPS for FEDFUNDS across different $q$ suggests almost no differences between choosing $q = 0$ and $q = 2$.

## 6.4 A note on the computational burden

Even though the efficient sampling schemes outlined in this paper help to overcome absolutely prohibitive computational burdens, the CPU time needed to perform fully Bayesian inference in a model of this size can still be considered substantial. In what follows we shed light on the estimation time required and how it is related to the length of the time series $T$, the lag length $p$ and to the number of latent factors $q \in \{0, 50\}$. Figure 4 shows the time needed to perform a single draw from the joint posterior distribution of the $215 + 215^2 p$ coefficients and their corresponding $2(215 + 215^2 p) + 1$ auxiliary shrinkage quantities, the $qT$ factor realizations and the associated $215q$ loadings, alongside $(T + 1)(215 + q)$ latent volatilities with their corresponding $645 + 2q$ parameters. This amounts to $166\,841$ random
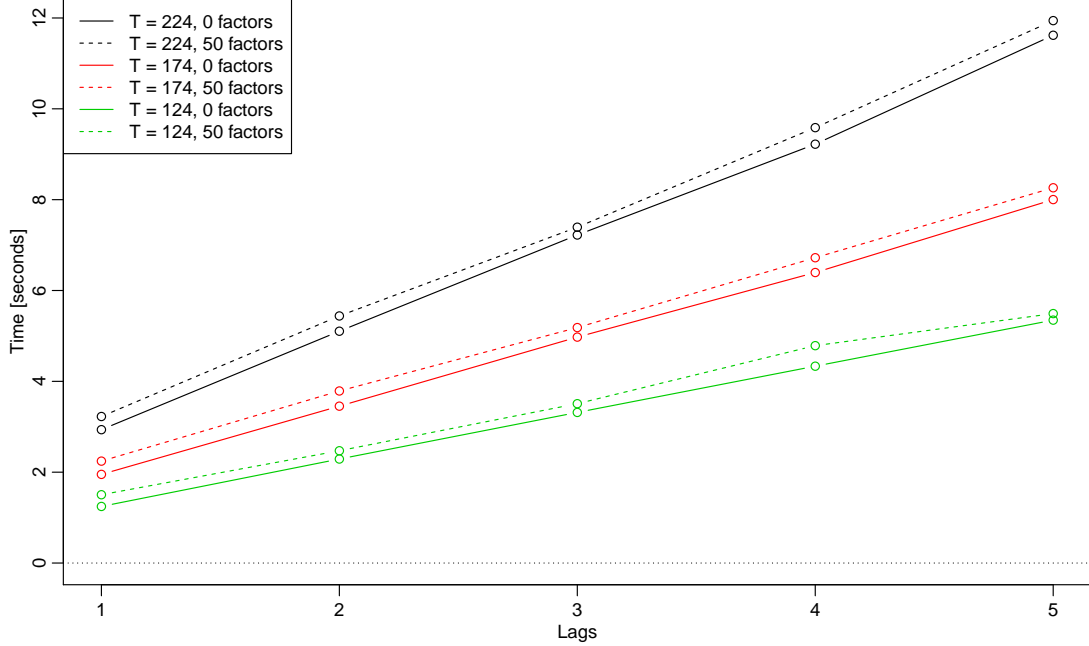
Figure 4: Empirical CPU times for each MCMC iteration on a standard laptop computer using one core. Time series lengths are $T \in \{124, 174, 224\}$; the numbers of latent factors are $q \in \{0, 50\}$.

draws for the smallest model considered (one lag, no factors, $T = 124$) and $776\,341$ random draws for the largest model (5 lags, 50 factors, $T = 224$) at each MCMC iteration.

As mentioned above, the computation time rises approximately linearly with the number of lags included. Dotted lines indicate the time in seconds needed to perform a single draw from a model with 50 factors included while solid lines refer to the time needed to estimate a model without factors and a diagonal time-varying variance-covariance matrix $\Sigma_t$. Interestingly, the additional complexity when moving from a model without factors to a highly parameterized model with 50 factors appears to be negligible, increasing the time needed by a fraction of a second on average. The important role of the length of the sample can be seen by comparing the green, red and black lines. The time necessary to perform a simple MCMC draw quickly rises with the length of our sample, consistent with the statements made in Section 2.4. This feature of our algorithm, however, is convenient especially when researchers are interested in combining many short time series or performing recursive forecasting based on a tiny initial estimation sample.

24

# 7  Closing remarks

In this paper we propose an alternative route to estimate huge dimensional VAR models that allow for time-variation in the error variances. The Dirichlet-Laplace prior, a recent variant of a global-local shrinkage prior, enables us to heavily shrink the parameter space towards the prior model while providing enough flexibility that individual regression coefficients are allowed to be unrestricted. This prior setup alleviates overfitting issues generally associated with large VAR models. To cope with computational issues we assume that the one-step-ahead forecast errors of the VAR feature a factor stochastic volatility structure that enables us to perform equation-by-equation estimation, conditional on the loadings and the factors. Since posterior simulation of each equation's autoregressive parameters involves manipulating large matrices, we implement an alternative recent algorithm that improves upon existing methods by large margins, rendering a fully fledged Bayesian estimation of truly huge systems possible.

In an empirical application we first present various key features of our approach based on a single factor model. This single factor which summarizes the joint dynamics of the VAR errors can be interpreted as an uncertainty measure that closely tracks observed factors such as the volatility index. The question whether such a simplistic structure proves to be an adequate representation of the time-varying covariance matrix naturally arises and we thus provide a detailed forecasting exercise to evaluate the merits of our approach relative to the prior model and a set of competing models with a different number of latent factors in the errors.

Finally, two potential generalizations are worth mentioning. First, note that it is trivial to relax the assumption of symmetry for the DL components. In the context of VARs, this might be of particular interest for distinguishing diagonal ($a_{\mathrm{D}}$ large) from off-diagonal ($a_{\mathrm{O}}$ small) elements in the spirit of the Minnesota prior or increasing the amount of shrinkage with increasing lag order (cf. Huber and Feldkircher, 2018, for a similar setup in the context of the Normal-Gamma shrinkage prior). Second, we would like to stress that our approach could also be used to estimate huge dimensional time-varying parameter VAR models with stochastic volatility. To cope with the computational difficulties associated with the vast state space, a possible approach could be to rely on an additional layer of hierarchy that

imposes a dynamic factor structure on the time-varying autoregressive coefficients and thus reduce the computational burden considerably.

# References

Aguilar, O. and M. West (2000). Bayesian dynamic factor models and portfolio allocation. *Journal of Business & Economic Statistics 18*(3), 338–357.

Ahelegbey, D. F., M. Billio, and R. Casarin (2016). Sparse graphical vector autoregression: A Bayesian approach. *Annals of Economics and Statistics* (123/124), 333–361.

Bańbura, M., D. Giannone, and L. Reichlin (2010). Large Bayesian vector auto regressions. *Journal of Applied Econometrics 25*(1), 71–92.

Bhattacharya, A., A. Chakraborty, and B. K. Mallick (2016). Fast sampling with Gaussian scale-mixture priors in high-dimensional regression. *Biometrika 4*(103), 985–991.

Bhattacharya, A., D. Pati, N. S. Pillai, and D. B. Dunson (2015). Dirichlet–Laplace priors for optimal shrinkage. *Journal of the American Statistical Association 110*(512), 1479–1490.

Bitto, A. and S. Frühwirth-Schnatter (2018). Achieving shrinkage in a time-varying parameter model framework. *Journal of Econometrics forthcoming.* arXiv pre-print 1611.01310.

Carriero, A., T. Clark, and M. Marcellino (2015). Large vector autoregressions with asymmetric priors and time varying volatilities. Working Papers 759, Queen Mary University of London, School of Economics and Finance.

Carriero, A., T. E. Clark, and M. Marcellino (2016). Common drifting volatility in large Bayesian VARs. *Journal of Business & Economic Statistics 34*(3), 375–390.

Carvalho, C. M., N. G. Polson, and J. G. Scott (2010). The horseshoe estimator for sparse signals. *Biometrika 97*(2), 465–480.

Clark, T. E. (2011). Real-time density forecasts from Bayesian vector autoregressions with stochastic volatility. *Journal of Business & Economic Statistics 29*(3), 327–341.

Clark, T. E. and F. Ravazzolo (2015). Macroeconomic forecasting performance under alternative specifications of time-varying volatility. *Journal of Applied Econometrics 30*(4), 551–575.

Cogley, T. and T. J. Sargent (2002, May). Evolving post-World War II U.S. inflation dynamics. In B. S. Bernanke and K. Rogoff (Eds.), *NBER Macroeconomics Annual 2001*, pp. 331–388. National Bureau of Economic Research, Inc.

Davis, R. A., P. Zang, and T. Zheng (2016). Sparse vector autoregressive modeling. *Journal of Computational and Graphical Statistics 25*(4), 1077–1096.

Doan, T. R., B. R. Litterman, and C. A. Sims (1984). Forecasting and conditional projection using realistic prior distributions. *Econometric Reviews 3*(1), 1–100.

Follett, L. and C. Yu (2017). Achieving parsimony in Bayesian VARs with the Horseshoe prior. arXiv pre-print 1709.07524.

Frühwirth-Schnatter, S. and H. Wagner (2010). Stochastic model specification search for Gaussian and partial non-Gaussian state space models. *Journal of Econometrics 154*(1), 85–100.

George, E. I., D. Sun, and S. Ni (2008). Bayesian stochastic search for VAR model restrictions. *Journal of Econometrics 142*(1), 553–580.

Griffin, J. E. and P. J. Brown (2010). Inference with normal-gamma prior distributions in regression problems. *Bayesian Analysis 5*(1), 171–188.

Hörmann, W. and J. Leydold (2013). Generating generalized inverse Gaussian random variates. *Statistics and Computing 24*(4), 1–11.

Huber, F. and M. Feldkircher (2018). Adaptive shrinkage in Bayesian vector autoregressive models. *Journal of Business & Economic Statistics forthcoming.*

Kastner, G. (2016). Dealing with stochastic volatility in time series using the R package stochvol. *Journal of Statistical Software 69*(5), 1–30.

Kastner, G. (2017). *factorstochvol: Bayesian estimation of (sparse) latent factor stochastic volatility models.* R package version 0.8.4.

Kastner, G. (2018). Sparse Bayesian time-varying covariance estimation in many dimensions. *Journal of Econometrics forthcoming.* arXiv preprint arXiv:1608.08468.

Kastner, G. and S. Frühwirth-Schnatter (2014). Ancillarity-sufficiency interweaving strategy (asis) for boosting MCMC estimation of stochastic volatility models. *Computational Statistics & Data Analysis 76*, 408–423.

Kastner, G., S. Frühwirth-Schnatter, and H. F. Lopes (2017). Efficient Bayesian inference for multivariate factor stochastic volatility models. *Journal of Computational and Graphical Statistics* (26), 905–917.

Koop, G. and D. Korobilis (2013). Large time-varying parameter VARs. *Journal of Econometrics 177*(2), 185–198.

Koop, G., D. Korobilis, and D. Pettenuzzo (2016). Bayesian compressed vector autoregressions. Working Papers 103, Brandeis University, Department of Economics and International Businesss School.

Koop, G., R. Leon-Gonzalez, and R. W. Strachan (2009). On the evolution of the monetary policy transmission mechanism. *Journal of Economic Dynamics and Control 33*(4), 997–1017.

Koop, G. M. (2013). Forecasting with medium and large Bayesian VARs. *Journal of*

*Applied Econometrics 28*(2), 177–203.

Korobilis, D. and D. Pettenuzzo (2016). Adaptive Minnesota prior for high-dimensional vector autoregressions. Essex finance centre working papers, University of Essex, Essex Business School.

Leydold, J. and W. Hörmann (2017). *GIGrvg: Random variate generator for the GIG distribution.* R package version 0.5.

Litterman, R. (1986). Forecasting with Bayesian vector autoregressions – Five years of experience. *Journal of Business and Economic Statistics 4*(1), 25–38.

McCracken, M. W. and S. Ng (2016). FRED-MD: A monthly database for macroeconomic research. *Journal of Business & Economic Statistics 34*(4), 574–589.

Park, T. and G. Casella (2008). The Bayesian Lasso. *103*(452), 681–686.

Pati, D., A. Bhattacharya, N. S. Pillai, and D. Dunson (2014). Posterior contraction in sparse Bayesian factor models for massive covariance matrices. *Annals of Statistics 42*(3), 1102–1130.

Pitt, M. K. and N. Shephard (1999). Time-varying covariances: A factor stochastic volatility approach. In J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith (Eds.), *Bayesian Statistics 6 – Proceedings of the Sixth Valencia International Meeting*, pp. 547–570. Oxford University Press.

Polson, N. G. and J. G. Scott (2011). Shrink globally, act locally: Sparse Bayesian regularization and prediction. In J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, and M. West (Eds.), *Bayesian Statistics 9 – Proceedings of the Ninth Valencia International Meeting*, pp. 501–538. Oxford University Press.

Primiceri, G. E. (2005). Time varying structural vector autoregressions and monetary policy. *The Review of Economic Studies 72*(3), 821–852.

Sims, C. A. and T. Zha (1998). Bayesian methods for dynamic multivariate models. *International Economic Review 39*(4), 949–968.

Sims, C. A. and T. Zha (2006). Were there regime switches in U.S. monetary policy? *American Economic Review 96*(1), 54–81.

Stock, J. H. and M. W. Watson (2005). Understanding changes in international business cycle dynamics. *Journal of the European Economic Association 3*(5), 968–1006.

Stock, J. H. and M. W. Watson (2011). Dynamic factor models. In M. P. Clements and D. F. Henry (Eds.), *Oxford Handbook of Economic Forecasting*, pp. 35–59. Oxford: Oxford University Press.

# A    Additional results
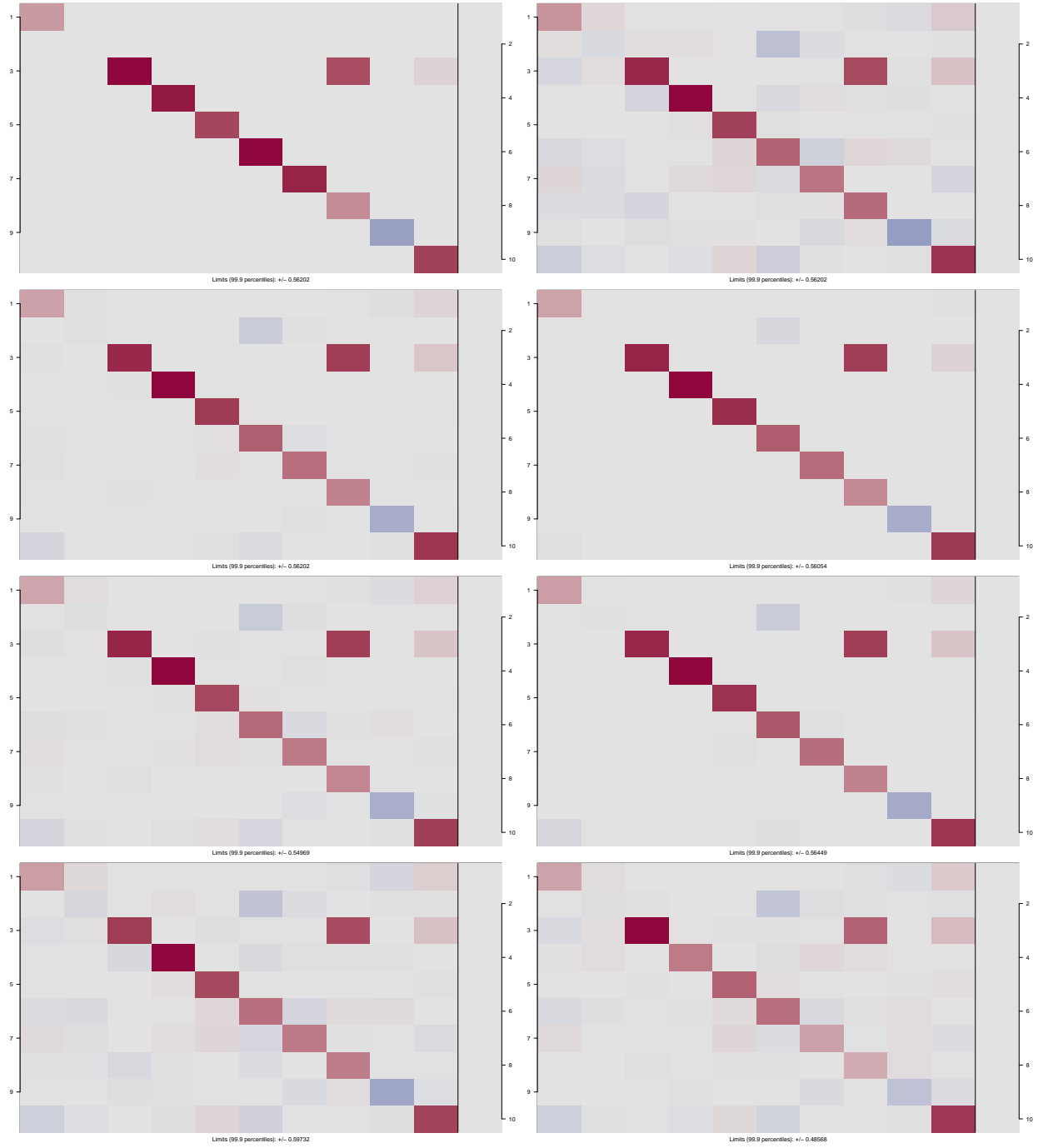
This appendix contains some additional results.

Figure 5: Exemplary visualization of the true and estimated VAR coefficients in the *sparse* scenario where $T = 250$ and $m = 10$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/K = 1/110$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/10$ (left) and $a_M = 1/1000$ (right).
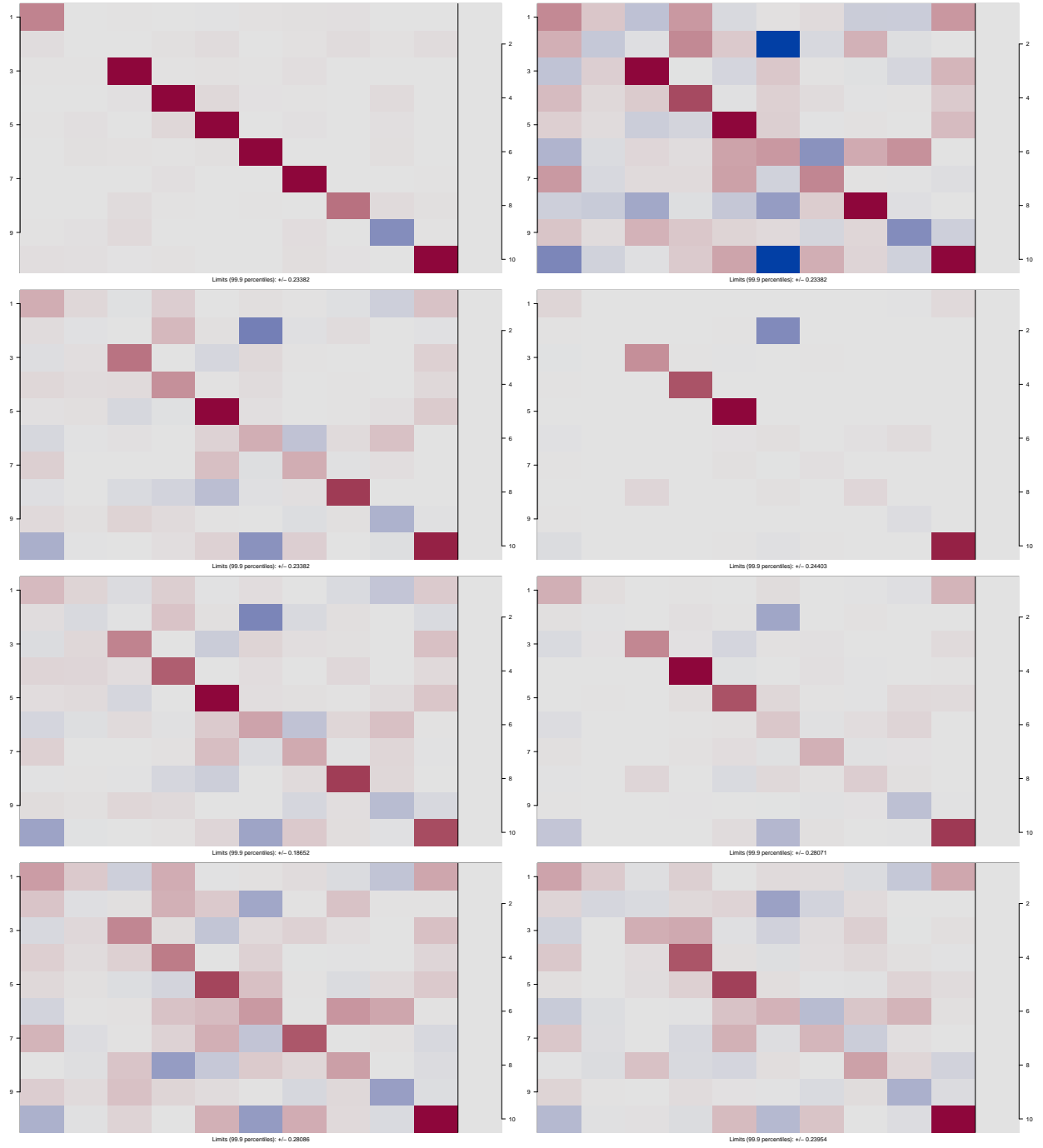
Figure 6: Exemplary visualization of the true and estimated VAR coefficients in the *dense* scenario where $T = 250$ and $m = 10$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/K = 1/110$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/10$ (left) and $a_M = 1/1000$ (right).
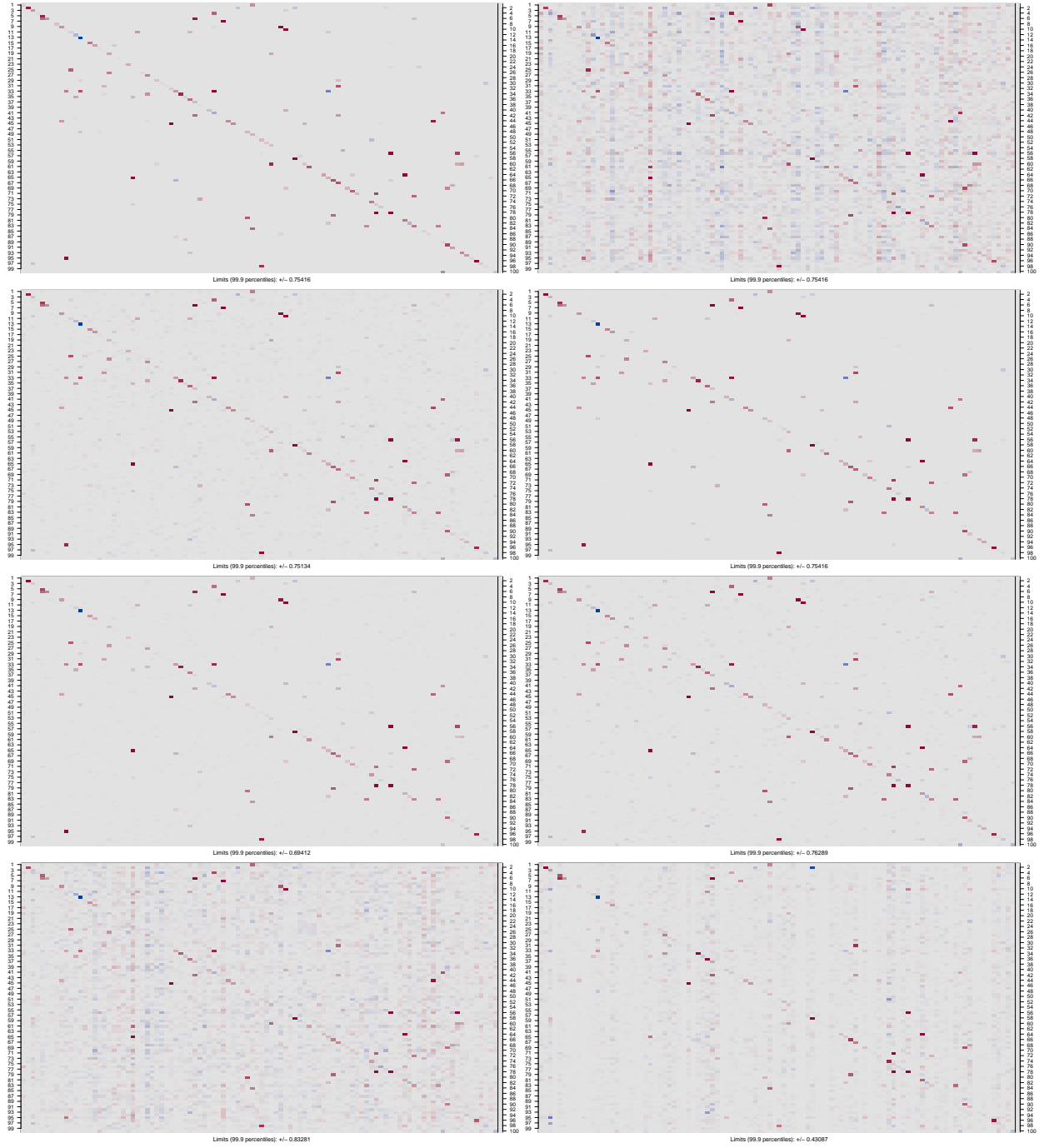
Figure 7: Exemplary visualization of the true and estimated VAR coefficients in the *sparse* scenario where $T = 250$ and $m = 100$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/K = 1/10100$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/10$ (left) and $a_M = 1/1000$ (right).
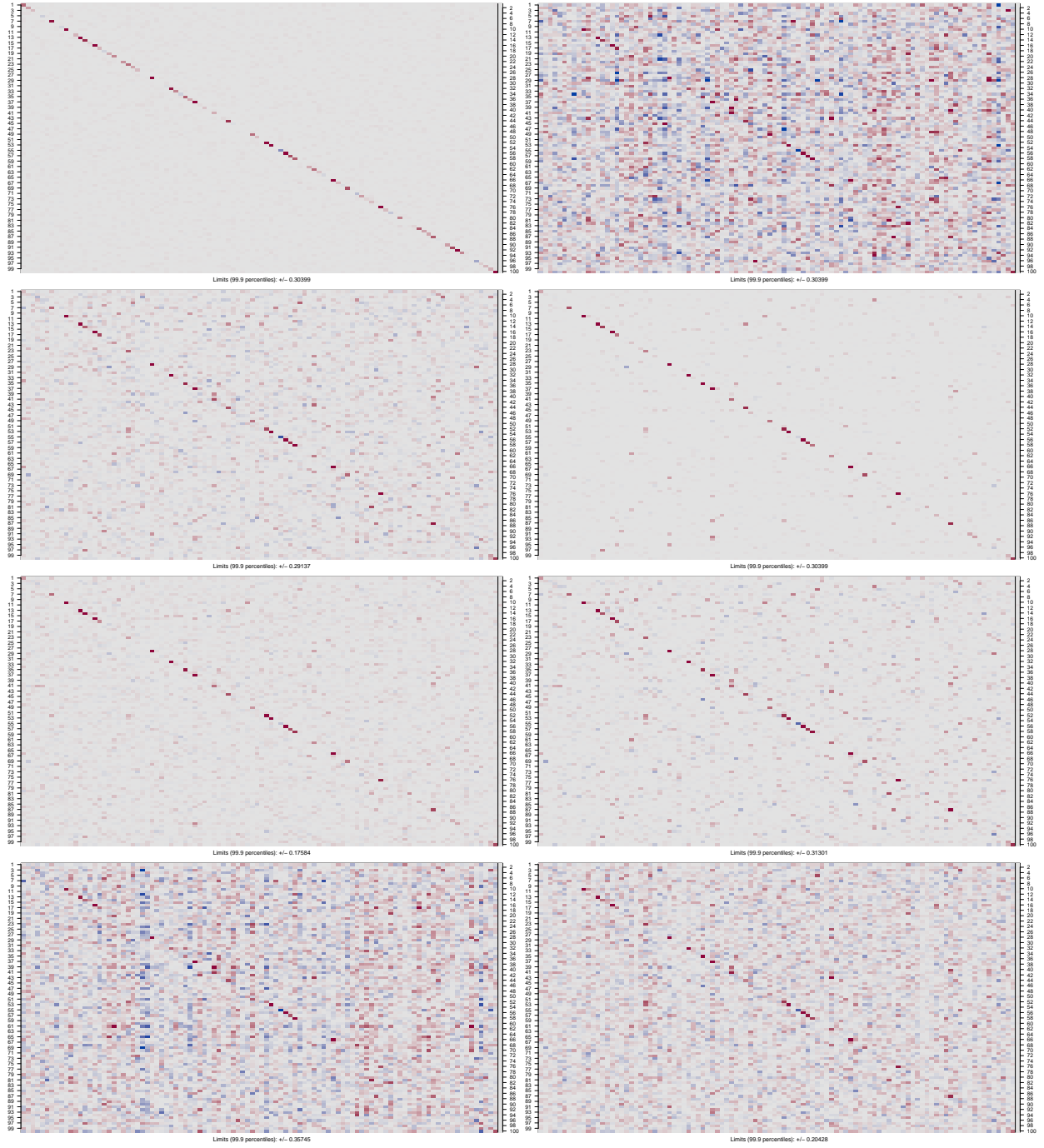
Figure 8: Exemplary visualization of the true and estimated VAR coefficients in the *dense* scenario where $T = 250$ and $m = 100$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/K = 1/10100$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/10$ (left) and $a_M = 1/1000$ (right).

The top panel of Figure 9 displays the posterior median estimates for the real data example discussed in Section 6 of the main paper when the shrinkage parameter $a$ is chosen to be $1/2$ (cf. Bhattacharya et al., 2015, for a discussion of this choice). While $a = 1/2$ appears to provide a fair amount of shrinkage in other applications, for our huge dimensional example this prior exerts only relatively little shrinkage and appears to lead to overfitting. The diagonal pattern in the first two lags known from Figure 2 in the main paper appears here as well, but there is a considerable amount of nonzero medians elsewhere. Correspondingly, the interquartile ranges visualized in the bottom panel of Figure 9 are also very large compared to those obtained under the much more tight prior used in the main paper.

Figure 9: Posterior medians (top) and posterior interquartile ranges (bottom) of VAR coefficients, $a = 1/2$.