

Variational Bayes Inference for Large Vector Autoregressions

Reza Hajargasht

Tomasz Wozniak

University of Melbourne

Melbourne, July 2016

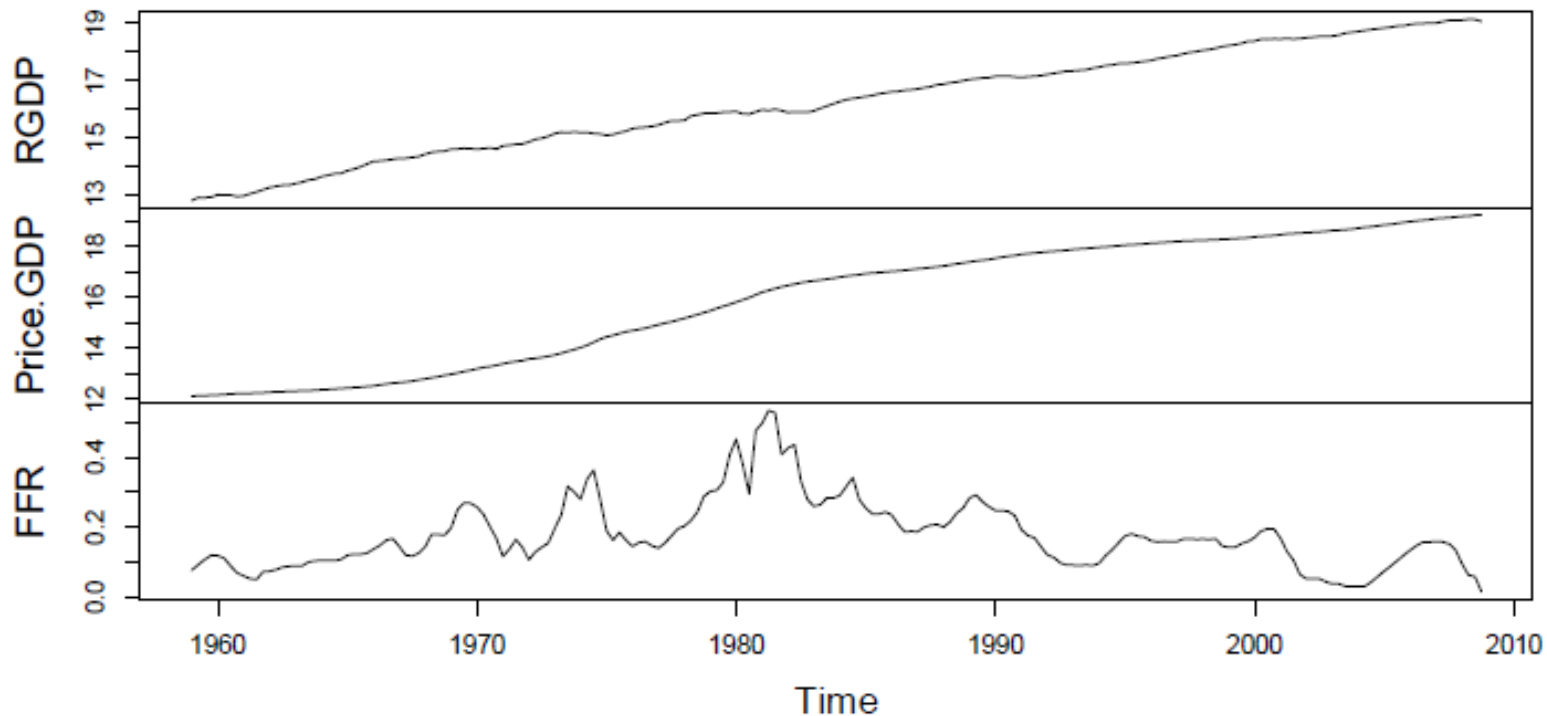
VAR Models

- Understanding the *dynamics of macro variables*, their *relationship* and *forecasting*.
- VAR models popularized by Chris SIMS in 80s is one the most important tools
- A typical small VAR model

$$\begin{cases} gdp_t = \alpha_1 + \alpha_{11}gdp_{t-1} + \dots + \alpha_{14}gdp_{t-4} + \beta_{11}p_{t-1} + \dots + \beta_{14}p_{t-4} + \gamma_{11}r_{t-1} + \dots + \gamma_{14}r_{t-4} + \varepsilon_{1,t} \\ P_t = \alpha_2 + \alpha_{21}gdp_{t-1} + \dots + \alpha_{24}gdp_{t-4} + \beta_{21}p_{t-1} + \dots + \beta_{24}p_{t-4} + \gamma_{22}r_{t-1} + \dots + \gamma_{24}r_{t-4} + \varepsilon_{2,t} \\ r_t = \alpha_3 + \alpha_{31}gdp_{t-1} + \dots + \alpha_{34}gdp_{t-4} + \beta_{31}p_{t-1} + \dots + \beta_{34}p_{t-4} + \gamma_{31}r_{t-1} + \dots + \gamma_{34}r_{t-2} + \varepsilon_{3,t} \end{cases}$$

$$\text{cov}(\boldsymbol{\varepsilon}_t) = \boldsymbol{\Sigma}_{3 \times 3}$$

UK Macroeconomics Series



We consider 3 cases with 3, 7 and 22 variables.
Source: Giannone, Lenza and Primiceri (2012) and
Stock and Watson (2008)

Bayesian VAR

large number of parameters

leading to imprecise estimation and forecasting

- Bayesians address this by specifying **Shrinkage Priors**

Our Contribution

1. Derivation of algorithms:

Variational Bayes estimation

Forecasting

for various VAR models [Conjugate and Independent Priors]

2. Optimal hyper-parameters

3. Providing an accurate method for computation of Marginal Likelihood using VB

Variational Bayes

- VB is an alternative to MCMC developed in Machine Learning finding its way into statistics and Econometrics
- The basic idea is to approximate the posterior $p(\boldsymbol{\theta} | \mathbf{y})$ with another tractable density function $q(\boldsymbol{\theta})$

Variational Bayes

- .

Problem setting.

Approximate the posterior distribution, $p(\theta|\mathbf{y})$, with another distribution, $q(\theta)$, that is **tractable**.

Characterisation.

	Variational Bayes	simulations
Calculations:	approximate	exact
Convergence:	deterministic	stochastic
Methods:	mean field theory calculus of variations	MCMC

Variational Bayes

- Minimize K-L Distance

$$\text{Min}_q KL(q \parallel p) = \int q(\boldsymbol{\theta}) \log \left\{ \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta} | \mathbf{y})} \right\} d\boldsymbol{\theta}$$

subject to $q(\boldsymbol{\theta}) = \prod_{k=1}^M q_k(\boldsymbol{\theta}_k)$

where $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K\}$ is a partition of $\boldsymbol{\theta}$

Variational Bayes

- This is a calculus of variation problem. It can be shown:
- Initialize with $q_2^*(\boldsymbol{\theta}_2), \dots, q_M^*(\boldsymbol{\theta}_M)$ and cycle

$$\left\{ \begin{array}{l} q_1^*(\boldsymbol{\theta}_1) \leftarrow \frac{\exp(E_{\boldsymbol{\theta}_{-1}} \ln p(\mathbf{y}, \boldsymbol{\theta}))}{\int \exp(E_{\boldsymbol{\theta}_{-1}} \ln p(\mathbf{y}, \boldsymbol{\theta})) d\boldsymbol{\theta}_1} \\ \vdots \\ q_M^*(\boldsymbol{\theta}_M) \leftarrow \frac{\exp(E_{\boldsymbol{\theta}_{-M}} \ln p(\mathbf{y}, \boldsymbol{\theta}))}{\int \exp(E_{\boldsymbol{\theta}_{-M}} \ln p(\mathbf{y}, \boldsymbol{\theta})) d\boldsymbol{\theta}_M} \end{array} \right.$$

Lower Bound for Marginal Likelihood

- It has been shown that ML or MDD

$$\ln ML = E_q \ln p(\mathbf{y}, \boldsymbol{\theta}) - E_q \ln q(\boldsymbol{\theta}) + KL(q, p)$$

- Define

$$\ln \underline{ML} = E_q \ln p(\mathbf{y}, \boldsymbol{\theta}) - E_q \ln q(\boldsymbol{\theta})$$

- Since KL is always positive, \underline{ML} provides a lower bound to marginal likelihood

VAR Model

M = number of endogenous variables

T = number of time periods

d = number of lags

$$\mathbf{y}_t = \boldsymbol{\gamma}_0 + \sum_{j=1}^d \boldsymbol{\Gamma}_j \mathbf{y}_{t-j} + \boldsymbol{\varepsilon}_t$$

$$\text{cov}(\boldsymbol{\varepsilon}_t) = \boldsymbol{\Sigma}_{M \times M}$$

Matrix Form

$$\mathbf{Y}_{T \times M} = \mathbf{X}_{T \times K} \boldsymbol{\Gamma}_{K \times M} + \mathbf{E}_{T \times M}$$

Conjugate Priors

$$P(\mathbf{\Gamma}, \mathbf{\Sigma}^{-1}) = P(\mathbf{\Gamma} | \mathbf{\Sigma}^{-1})P(\mathbf{\Sigma}^{-1})$$

where

$$P(\mathbf{\Gamma} | \mathbf{\Sigma}^{-1}) = MN(\underline{\mathbf{\Gamma}}, \mathbf{\Sigma}, \underline{\mathbf{V}})$$

$$P(\mathbf{\Sigma}^{-1}) = W(\underline{\mathbf{S}}^{-1}, \underline{\nu})$$

- 1) Normal-Wishart
- 2) Minnesota (proposed by Litterman (1980) and been used in many studies such as BGR 2010)
- 3) Conjugate SSVS (Koop 2012)
- 4) Hierarchical (Giannone, Lenza and Primiceri, 2012)

Minnesota

- For $\underline{\gamma} = \text{vec}(\underline{\Gamma})$

$\underline{\gamma}_i = 0$ except for first own lag which is 0.9 or 1

- $\underline{\mathbf{V}}$ is a diagonal matrix with (e.g. size 13)

$$\left\{ \begin{array}{l} \text{Off-diagonal elements of } \underline{\mathbf{V}} \text{ set to zeroes} \\ \underline{\mathbf{V}}_1 = \lambda_0 \text{ for intercept} \\ \underline{\mathbf{V}}_j = \frac{\lambda_1^2}{r^2 \sigma_{jj}^2} \text{ corresponding to lag } r \\ \sigma_{ii} = s_i^2 \end{array} \right.$$

Posterior without Factorization

- Exact Posterior

$$\begin{cases} \boldsymbol{\Gamma} \mid \boldsymbol{\Sigma}^{-1}, \mathbf{Y} \sim MN\{\bar{\boldsymbol{\Gamma}}, \boldsymbol{\Sigma}, \bar{\mathbf{V}}\} \\ \boldsymbol{\Sigma}^{-1} \mid \mathbf{Y} \sim W(\bar{\mathbf{S}}^{-1}, \bar{\nu}) \end{cases}$$

where

$$\bar{\mathbf{V}} = [\underline{\mathbf{V}}^{-1} + \mathbf{X}'\mathbf{X}]^{-1} \quad \bar{\boldsymbol{\Gamma}} = \bar{\mathbf{V}}[\underline{\mathbf{V}}^{-1}\underline{\boldsymbol{\Gamma}} + \mathbf{X}'\mathbf{Y}]$$

$$\bar{\mathbf{S}} = (\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\Gamma}})'(\mathbf{Y} - \mathbf{X}\bar{\boldsymbol{\Gamma}}) + \underline{\mathbf{S}} + (\bar{\boldsymbol{\Gamma}}' - \underline{\boldsymbol{\Gamma}}')\underline{\mathbf{V}}^{-1}(\bar{\boldsymbol{\Gamma}}' - \underline{\boldsymbol{\Gamma}}) \quad \bar{\nu} = T + \underline{\nu}$$

- Posterior Moments

$$E(\boldsymbol{\Gamma}) = \bar{\boldsymbol{\Gamma}}$$

$$Var(\boldsymbol{\gamma}) = \frac{1}{T + \underline{\nu} - p - 1} \bar{\mathbf{S}} \otimes \bar{\mathbf{V}}$$

$$E(\boldsymbol{\Sigma}^{-1}) = (T + \underline{\nu})\bar{\mathbf{S}}^{-1}$$

$$Var(\sigma_{ij}^+) = \bar{\nu} (s_{ij}^{+2} + s_{ii}^+ s_{jj}^+)$$

Posterior with Factorization

- Factorization $q(\mathbf{\Gamma}, \mathbf{\Sigma}^{-1}) = q(\mathbf{\Gamma})q(\mathbf{\Sigma}^{-1})$

- VB Posterior

$$\begin{cases} \mathbf{\Gamma} \sim MN\{\bar{\mathbf{\Gamma}}, \bar{\mathbf{\Sigma}}^{-1}, \bar{\mathbf{V}}\} \\ \mathbf{\Sigma}^{-1} \sim W(\bar{\mathbf{S}}_{VB}^{-1}, \bar{\nu}_{VB}) \end{cases}$$

where

$$\bar{\mathbf{V}} = [\underline{\mathbf{V}}^{-1} + \mathbf{X}'\mathbf{X}]^{-1} \quad \bar{\mathbf{\Gamma}} = \bar{\mathbf{V}}[\underline{\mathbf{V}}^{-1}\underline{\mathbf{\Gamma}} + \mathbf{X}'\mathbf{Y}]$$

$$\bar{\mathbf{S}}_{VB} = \frac{\bar{\nu}_{VB}}{\bar{\nu}} \bar{\mathbf{S}} \quad \bar{\nu}_{VB} = T + p + \underline{\nu}$$

- Posterior Means

$$E(\mathbf{\Gamma}_{VB}) = \bar{\mathbf{\Gamma}}$$

$$Var(\boldsymbol{\gamma}_{VB}) = \frac{1}{T + \underline{\nu}} \bar{\mathbf{S}} \otimes \bar{\mathbf{V}}$$

$$E(\mathbf{\Sigma}^{-1}) = (T + \underline{\nu}) \bar{\mathbf{S}}^{-1}$$

$$Var(\sigma_{ij}^+) = \frac{\bar{\nu}_G^2}{\bar{\nu}_{VB}} (s_{ij}^{+2} + s_{ii}^+ s_{jj}^+)$$

VB vs Full Bayes

- Posterior Means for Γ and Σ^{-1} are the same
- Posterior SE for factorized model are smaller

$$\frac{\text{Var}(\gamma_{VB})}{\text{Var}(\gamma_G)} = \frac{T + \underline{\nu} - p - 1}{T + \underline{\nu}}$$

$$\frac{\text{Var}(\sigma_{ijVB}^+)}{\text{Var}(\sigma_{ijG}^+)} = \frac{T + \underline{\nu}}{T + \underline{\nu} + p}$$

- KL Distance

$$KL = -\frac{Mp}{2} \ln 2e + \frac{M}{2} \ln \left(\frac{\bar{\nu}_{VB} \bar{\nu}_{VB}}{\bar{\nu} \bar{\nu}} \right) - \ln \left(\frac{\Gamma_M(\bar{\nu}_{VB}/2)}{\Gamma_M(\bar{\nu}/2)} \right)$$

- KL increases with increase in M and d and decreases with increase in T and $\underline{\nu}$. It does not depend on data and other hyper-parameters!

VB vs Full Bayes

$$\lim_{T \rightarrow \infty} KL \rightarrow \frac{1}{2} \sum_{j=1}^m \ln(1) = 0$$

$$\lim_{\underline{v} \rightarrow \infty} KL \rightarrow \frac{1}{2} \sum_{j=1}^m \ln(1) = 0$$

$$\lim_{p \rightarrow \infty} KL \rightarrow \frac{1}{2} \sum_{j=1}^m \ln(0) = -\infty$$

$$\lim_{p \rightarrow 0} KL \rightarrow \frac{1}{2} \sum_{j=1}^m \ln(1) = 0$$

Another way of writing VAR

Another way of writing VAR

$$\begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_T \end{pmatrix} = \begin{pmatrix} \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_T \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta}_1 \\ \vdots \\ \boldsymbol{\beta}_M \end{pmatrix} + \boldsymbol{\varepsilon} \quad \text{where} \quad \mathbf{Z}_t = \begin{pmatrix} z_{1t} & 0 & \cdots & 0 \\ 0 & z_{2t} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & z_{Mt} \end{pmatrix}$$

$$z_t = (y_{1,t}, \dots, y_{m,t}, y_{1,t-1}, \dots, y_{m,t-1}, \dots, y_{1,t-d}, \dots, y_{m,t-d})$$

$$\mathbf{Y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} \sim N(0, \mathbf{I}_T \otimes \boldsymbol{\Sigma}_M)$$

Independent Priors

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{-1}) = p(\boldsymbol{\beta}) p(\boldsymbol{\Sigma}^{-1})$$

where

$$P(\boldsymbol{\beta}) = N(\boldsymbol{\beta}, \underline{\mathbf{V}})$$

$$P(\boldsymbol{\Sigma}^{-1}) = W(\underline{\mathbf{S}}^{-1}, \underline{\nu})$$

1- Normal-Wishart

2- Minnesota [Kadiyala and Karlsson, 1997]

3- Lasso Type Priors [Korobilis 2012, Gefang 2012]

4- SSVS [George, Sun and Ni 2008, Korobilis 2013]

Minnesota

- For $\underline{\beta}$

$\underline{\beta}_i = 0$ except for first own lag which is 0.9 or 1

- \underline{V} is a diagonal matrix of size p e.g. 39

$$\left\{ \begin{array}{l} \text{Off-diagonal elements of } \underline{V} \text{ set to zeroes} \\ \underline{V}_1 = \lambda_0^2 \sigma_{ii}^2 \quad \text{for intercepts} \\ \underline{V}_{i,jj} = \frac{\lambda_1^2 \sigma_{ii}^2}{r^2} \quad \text{for own lag } r \\ \underline{V}_{i,jj} = \frac{\lambda_2^2 \sigma_{ii}^2}{r^2 \sigma_{jj}^2} \quad \text{for other variables lag} \end{array} \right.$$

Exact Posterior

Gibbs Sampler

$$\begin{cases} \boldsymbol{\beta} \mid \boldsymbol{\Sigma}^{-1}, \mathbf{Y} \sim N(\bar{\boldsymbol{\beta}}, \bar{\mathbf{V}}) \\ \boldsymbol{\Sigma}^{-1} \mid \boldsymbol{\beta}, \mathbf{Y} \sim W(\bar{\mathbf{S}}^{-1}, \bar{\nu}) \end{cases}$$

where

$$\bar{\mathbf{V}} = [\underline{\mathbf{V}}^{-1} + \sum_{t=1}^T \mathbf{Z}_t' \boldsymbol{\Sigma}^{-1} \mathbf{Z}_t]^{-1} \quad \bar{\boldsymbol{\beta}} = \bar{\mathbf{V}} [\underline{\mathbf{V}}^{-1} \underline{\boldsymbol{\beta}} + \sum_{t=1}^T \mathbf{Z}_t' \boldsymbol{\Sigma}^{-1} \mathbf{y}_t]$$

$$\bar{\mathbf{S}} = \underline{\mathbf{S}} + \sum_{t=1}^T (\mathbf{y}_t - \mathbf{Z}_t \boldsymbol{\beta})(\mathbf{y}_t - \mathbf{Z}_t \boldsymbol{\beta})' \quad \bar{\nu} = T + \underline{\nu}$$

VB Posterior

- VB Posterior
$$\begin{cases} \boldsymbol{\beta} \sim N(\bar{\boldsymbol{\beta}}, \bar{\mathbf{V}}) \\ \boldsymbol{\Sigma}^{-1} \sim W(\bar{\mathbf{S}}^{-1}, \bar{\nu}) \end{cases}$$
- Initialize with $\bar{\boldsymbol{\Sigma}}^{-1}$ and cycle

$$\bar{\mathbf{V}} \leftarrow [\underline{\mathbf{V}}^{-1} + \sum_{t=1}^T \mathbf{Z}_t' \bar{\boldsymbol{\Sigma}}^{-1} \mathbf{Z}_t]^{-1}$$

$$\bar{\boldsymbol{\beta}} \leftarrow \bar{\mathbf{V}} [\underline{\mathbf{V}}^{-1} \underline{\boldsymbol{\beta}} + \sum_{t=1}^T \mathbf{Z}_t' \bar{\boldsymbol{\Sigma}}^{-1} \mathbf{y}_t]$$

$$\bar{\mathbf{S}} \leftarrow \underline{\mathbf{S}} + \sum_{t=1}^T \{ (\mathbf{y}_t - \mathbf{Z}_t \bar{\boldsymbol{\beta}})(\mathbf{y}_t - \mathbf{Z}_t \bar{\boldsymbol{\beta}})' + \mathbf{Z}_t \bar{\mathbf{V}} \mathbf{Z}_t' \}$$

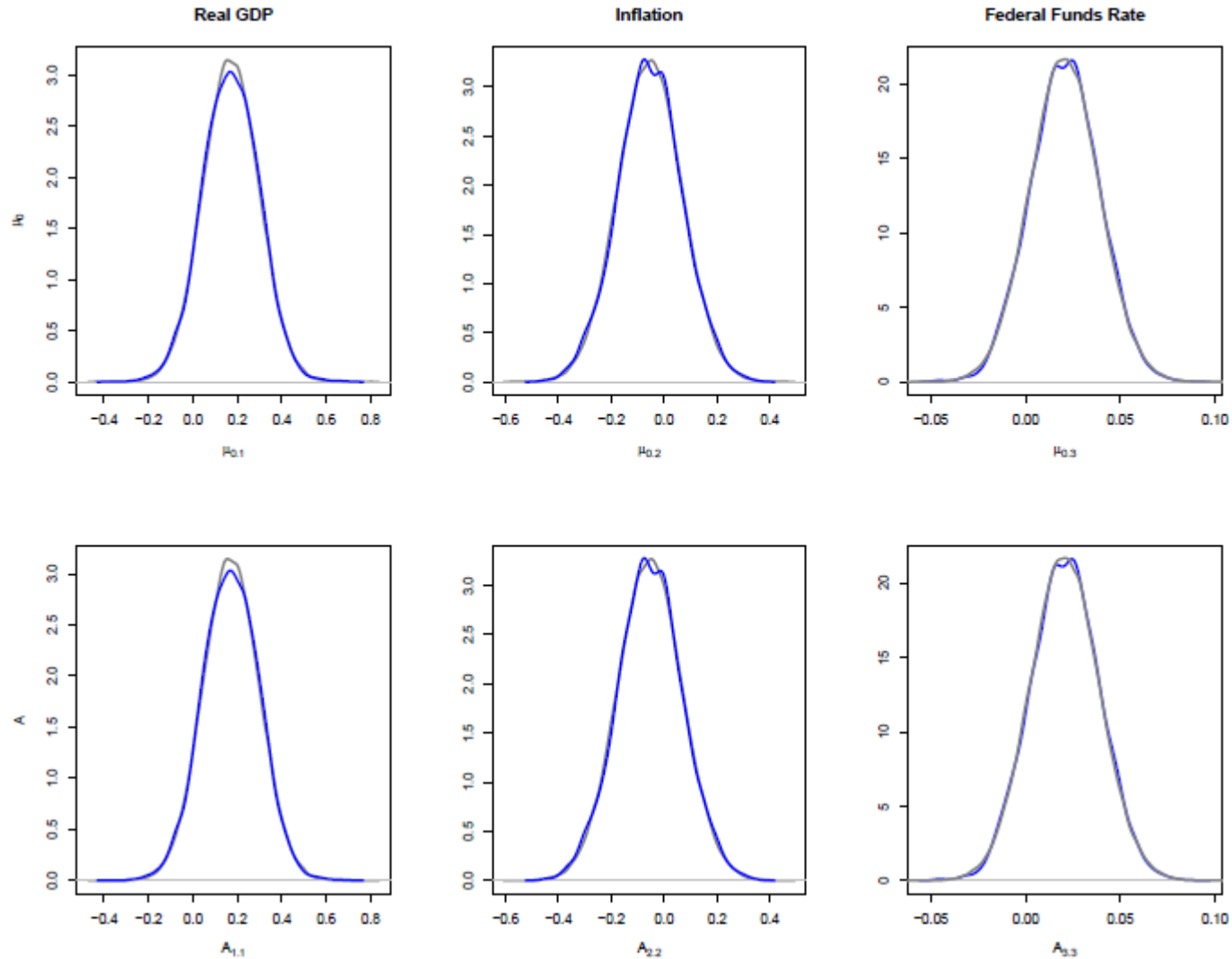
$$\bar{\nu} \leftarrow T + \underline{\nu}$$

$$\bar{\boldsymbol{\Sigma}}^{-1} \leftarrow \bar{\nu} \bar{\mathbf{S}}^{-1}$$

- Lower bound for Log of Marginal Likelihood

$$\begin{aligned} \ln ML = & \frac{p}{2} - \frac{MT}{2} \ln \pi + \ln \Gamma_M(\bar{\nu}/2) - \ln \Gamma_M(\underline{\nu}/2) + \frac{1}{2} (\ln |\bar{\mathbf{V}}| - \ln |\underline{\mathbf{V}}|) \\ & + \frac{1}{2} (\bar{\nu} \ln |\bar{\mathbf{S}}^{-1}| - \underline{\nu} \ln |\underline{\mathbf{S}}^{-1}|) - \frac{1}{2} \text{tr} \{ \underline{\mathbf{V}}^{-1} [(\bar{\boldsymbol{\beta}} - \underline{\boldsymbol{\beta}})(\bar{\boldsymbol{\beta}} - \underline{\boldsymbol{\beta}})' + \bar{\mathbf{V}}] \} \end{aligned}$$

Comparison



Optimal Hyper-parameter Selection

- Let λ be the vector of hyper-parameters and define

$$\ln \underline{ML}(\lambda) = \ln \underline{ML}(\phi | \lambda)$$

- Then optimal hyper-parameters are obtained

$$\lambda^* = \arg \max_{\lambda} \ln \underline{ML}(\lambda)$$

Other available techniques.

- Maximisation of $\ln MDD$ - available only for natural-conjugate prior distribution
- Point forecasting performance measures (Bańbura, Giannone, Reichlin, 2010, JAE)
- Hierarchical prior structures (see e.g. Giannone, Lenza, Primiceri, 2013, NBER)

Point Forecasting

- We are interested in

$$E(y_{i,\tau+h} | Data_{\tau})$$

- Even with VB posterior this does not seem to have an analytic form [with more than one step ahead forecast].
- We can use simulation though: simulate

$$\Sigma^{-1(j)} \sim W(\bar{S}^{-1}, \bar{V}) \quad \beta^{(j)} \sim N(\bar{\beta}, V_{\beta}) \quad \mathbf{u}_{\tau+h}^{(j)} \sim N(0, \Sigma)$$

and calculate recursively for $h=1, \dots, H$

$$\tilde{\mathbf{y}}_{\tau+h}^{(j)} = \mathbf{a}_0^{(j)} + \sum_{i=1}^{h-1} \mathbf{A}_i^{(j)} \tilde{\mathbf{y}}_{\tau+h-i}^{(j)} + \sum_{i=h}^{p_1} \mathbf{A}_i^{(j)} \mathbf{y}_{\tau+h-i}^{(j)} + \mathbf{u}_{\tau+h}^{(j)}$$

Point Forecasting

- But note this does not increase running time very significantly because:
 - 1- The draws are independent, therefore a small number of draws e.g. 1000 should be enough
 - 2- This doesn't have to be done in a loop
- Forecasting Metrics

$$MSFE = \frac{\sum_{\tau=\tau_0}^{T-h} [y_{i,\tau+h}^0 - E(y_{i,\tau+h} | Data_{\tau})]^2}{T - h - \tau_0 + 1}$$

Overall Forecasting

- Predictive likelihood

$$p[y_{i,\tau+h}^0 | Data_\tau] = \int p[y_{i,\tau+h}^0 | Data_\tau, \boldsymbol{\theta}] \pi(\boldsymbol{\theta} | Data_\tau) d\boldsymbol{\theta}$$

where $\pi(\boldsymbol{\theta} | Data_\tau)$ is the posterior.

- Even with VB posterior this doesn't seem analytically tractable but we can draw from the VB posterior and calculate the integral
- A quadratic approximation

$$\ln p[y_{i,\tau+h}^0 | Data_\tau] \sim -\frac{1}{2} \left\{ \ln 2\pi + \ln V_{i,\tau+h|\tau} + (y_{i,\tau+h}^0 - \bar{y}_{i,\tau+h|\tau})^2 / V_{i,\tau+h|\tau} \right\}$$

- Forecasting Metrics

$$\sum_{t=\tau_0}^{T-h} \ln p[y_{i,\tau+h}^0 | Data_\tau]$$

Forecasting Performance

- ▶ Models estimated with VB methods forecast **as accurately** as those estimated with the exact methods
 - ▶ for the natural-conjugate prior
benchmark: analytical solutions
 - ▶ for the independent prior
benchmark: Gibbs sampling output
- ▶ Hyper-parameters optimisation improves the forecasting performance
 - ▶ mixed evidence for point forecasts at some horizons
- ▶ Hyper-parameters optimization using InMDD_{VBLB} gives exactly the same results as using InMDD
- ▶ Nearly no improvements for independent priors

Accurate Computation of ML

- VB can be used for accurate computation of Marginal Likelihood under a reciprocal importance sampling framework e.g.

$$ML = M \left[\sum_{m=1}^M \frac{q(\boldsymbol{\theta}^m | \mathbf{y})}{p(\mathbf{y} | \boldsymbol{\theta}^m) p(\boldsymbol{\theta}^m)} \right]^{-1}$$

- Where $\boldsymbol{\theta}^m$ s are draws from MCMC and $q(\boldsymbol{\theta} | \mathbf{y})$ is the posterior from VB.
- This method has good properties if the candidate density has narrower tail which is usually the case with VB.

Marginal Likelihood

.

		VB_LB	GD_Gew	GD_VB	Bridge_VB	Chib	IS_VB	VB_UB
VAR_S 3 variable	Mean	635.25	635.27	635.42	635.42	635.42	635.42	635.62
	SD	NA	0.0462	0.0077	0.0056	0.0079	0.0107	0.0081
VAR_M 7 variable	Mean	961.63	958.23	962.77	962.77	962.71	962.76	964.09
	SD	NA	0.2940	0.0438	0.0164	0.0162	0.0528	0.0292

Conclusion

1. Derivation of fast algorithms:

VB Estimation

VB Forecasting

for Bayesian VARs with conjugate, independent and hierarchical priors

2. Selection of optimal values for hyper-parameters of the model.

3- Accurate Computation of Marginal Likelihood using VB Posterior in par with Chib method

4- Can be extended to more elaborate models