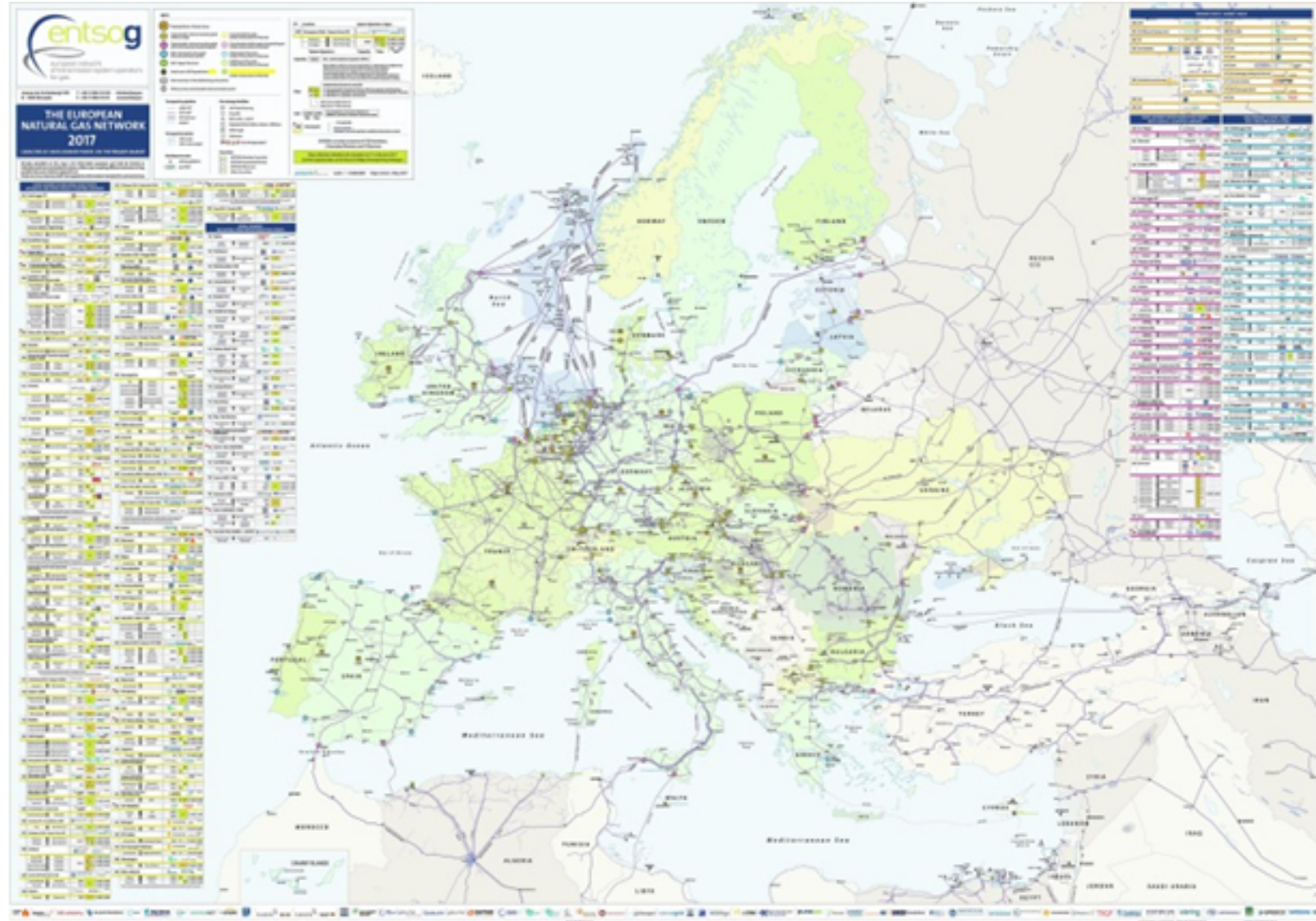


# Variational Bayesian Inference for Network Autoregression Models

Wei-Ting Lai, Ray-Bing Chen, Ying Chen, Thorsten Koch

## Introduction

- ▶ Networks have emerged and become widely available in various fields such as energy transmission, logistics, transportation, and financial systems. Networks are dynamic in terms of their temporal dependence, and often have large scales.
- ▶ Tremendous growth and heterogeneity in both nodes/edges and dependence over time are the key characteristics of such networks.



- ▶ The temporal dependence of network can be represented in the vector autoregression (VAR) modeling framework. Also call the network autoregression model (NAR).
- ▶ Unlike univariate time series, the temporal dependence of a multivariate series consists of not only the serial dependence within each marginal series but also the interdependence across different marginal series.

## Structure Selection in Network Autoregression Model

- ▶ Consider the following NAR/VAR model with lag- $p$ ,

$$\mathbf{Y}_t = \mathbf{Y}_{t-1}\mathbf{B}_1 + \dots + \mathbf{Y}_{t-p}\mathbf{B}_p + \varepsilon_t, \quad (1)$$

where  $\mathbf{Y}_t = (y_{t,1}, \dots, y_{t,m})$  is an  $1 \times m$  vector,  $\mathbf{B}_l$  is a  $m \times m$  coefficient matrix for  $l = 1, \dots, p$ , and  $\varepsilon_t, t = 1, \dots, T$ , are i.i.d.  $MN_{1 \times m}(\mathbf{0}, \Sigma)$  random vectors.

- ▶ The total number of coefficient is  $\mathcal{O}(m^2p)$ .
- ▶ If the nodes  $m$  be large or the temporal dependence  $p$  increases, the NAR/VAR model may have an overparameterization problem.
- ▶ To overcome the overparameterization, structure sparse assumption is adopted here.
- ▶ Here we consider a Bayesian analysis approach for structured NAR/VAR models.
- ▶ Instead of MCMC algorithms which take huge computational cost, we develop a variational Bayesian method for structure selection in the NAR/VAR model.
- ▶ Following Song and Bickel [6], three structures are considered here,

- ▶ Structure 1: Universal grouping (UG) in  $\mathbf{B}_l$

$$\begin{pmatrix} & \mathbf{M3} & \mathbf{M4} & \mathbf{M5} & \mathbf{M6} \\ \mathbf{M3} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{M4} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{M5} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{M6} & \bullet & \bullet & \bullet & \bullet \end{pmatrix} \rightarrow \begin{pmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{pmatrix} + \begin{pmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix}$$

(Total)                      (Own)                      (Others)

- ▶ Structure 2: Segmentation grouping (SG) in  $\mathbf{B}_l$

$$\begin{pmatrix} & \mathbf{M3} & \mathbf{M4} & \mathbf{I1} & \mathbf{I5} \\ \mathbf{M3} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{M4} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{I1} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{I5} & \bullet & \bullet & \bullet & \bullet \end{pmatrix} \rightarrow \begin{pmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{pmatrix} + \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix}$$

(Total)                      (Own)                      (Others)

- ▶ Structure 3: No grouping (NG) in  $\mathbf{B}_l$

$$\begin{pmatrix} & \mathbf{M3} & \mathbf{M4} & \mathbf{B1} & \mathbf{O1} \\ \mathbf{M3} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{M4} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{B1} & \bullet & \bullet & \bullet & \bullet \\ \mathbf{O1} & \bullet & \bullet & \bullet & \bullet \end{pmatrix} \rightarrow \begin{pmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{pmatrix} + \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix}$$

(Total)                      (Own)                      (Others)

- ▶ Among these three structures, it is clear that we can treat no grouping and universal grouping structures as special cases of the segmented grouping structure.

## Variational Bayesian Inference

- ▶ Instead of directly generating the posterior samples, the variational Bayesian (VB) approach is used to identify the best approximation distribution of the true posterior for the further Bayesian inference.
- ▶ The Kullback-Leibler divergence (KL-divergence) is used to measure the difference among  $\mathbf{P}$  (true posterior distribution) and  $\mathbf{Q}$  (approximation posterior distribution). Thus given a proper assumption of  $\mathbf{Q}$ , in the variational Bayesian approach, we solve the following minimization problem,

$$\min_{\mathbf{Q}} KL(\mathbf{Q}||\mathbf{P}) = \min_{\mathbf{Q}} E_{\mathbf{Q}} \left[ \log \frac{\mathbf{Q}}{\mathbf{P}} \right].$$

- ▶ This optimization problem can be solved via Expectation-Maximization (EM) type method.

## Notation

- ▶ Take the segmented grouping structure as an illustration.
- ▶ Let  $\mathbf{S} = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_g\}$  as the overall index set for the columns in each  $\mathbf{B}_\ell$ , and let  $\mathbf{s}_k$  denote the index set of the  $k$ -th segment with size  $n(\mathbf{s}_k)$ ,  $k = 1, 2, \dots, g$  and  $\sum_{k=1}^g n(\mathbf{s}_k) = m$ .
- ▶ Two indicators  $\gamma_{\ell,i,i}$  and  $\eta_{\ell,i,\tilde{\mathbf{s}}_k}$  to identify the active structures in  $\mathbf{B}_\ell$ . Here  $\tilde{\mathbf{s}}_k = \mathbf{s}_k \setminus \{i\}$

$$\begin{aligned} \gamma_{\ell,1,1} = 1 &\rightarrow \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} &\leftarrow \eta_{\ell,1,\tilde{\mathbf{s}}_1} = 0, \eta_{\ell,1,\tilde{\mathbf{s}}_2} = 0, \eta_{\ell,1,\tilde{\mathbf{s}}_3} = 0 \\ \gamma_{\ell,2,2} = 1 &\rightarrow \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} &\leftarrow \eta_{\ell,2,\tilde{\mathbf{s}}_1} = 1, \eta_{\ell,2,\tilde{\mathbf{s}}_2} = 0, \eta_{\ell,2,\tilde{\mathbf{s}}_3} = 0 \\ & &\leftarrow \eta_{\ell,3,\tilde{\mathbf{s}}_1} = 0, \eta_{\ell,3,\tilde{\mathbf{s}}_2} = 1, \eta_{\ell,3,\tilde{\mathbf{s}}_3} = 0 \\ \gamma_{\ell,5,5} = 0 &\rightarrow \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \end{aligned}$$

- ▶ Spike-and-slab prior:

$$\begin{aligned} \mathbf{B}_{\ell,i,i} | \gamma_{\ell,i,i}, \sigma_{\mathbf{B}}^2 &\sim \gamma_{\ell,i,i} \mathbf{N}(\mathbf{0}, \sigma_{\mathbf{B}}^2) + (1 - \gamma_{\ell,i,i}) \delta_0, \\ \mathbf{B}_{\ell,i,\tilde{\mathbf{s}}_k} | \eta_{\ell,i,\tilde{\mathbf{s}}_k}, \sigma_{\mathbf{B}}^2 &\sim \eta_{\ell,i,\tilde{\mathbf{s}}_k} MN_{1 \times n(\tilde{\mathbf{s}}_k)}(\mathbf{0}, \sigma_{\mathbf{B}}^2 \mathbf{I}_{n(\tilde{\mathbf{s}}_k)}) + (1 - \eta_{\ell,i,\tilde{\mathbf{s}}_k}) \delta_0, \\ \gamma_{\ell,i,i} &\sim \text{Ber}(\pi_1), \eta_{\ell,i,\tilde{\mathbf{s}}_k} \sim \text{Ber}(\pi_2). \end{aligned}$$

## The Proposed Variational Bayesian Approach

- ▶ To simplify the model structure, we reparametrize the coefficient as the product of the coefficient and the indicator.

$$\begin{aligned} \mathbf{B}_{\ell,i,i} &= \gamma_{\ell,i,i} \tilde{\mathbf{B}}_{\ell,i,i}, \tilde{\mathbf{B}}_{\ell,i,i} | \sigma_{\mathbf{B}}^2 \sim \mathbf{N}(\mathbf{0}, \sigma_{\mathbf{B}}^2) \\ \mathbf{B}_{\ell,i,\tilde{\mathbf{s}}_k} &= \eta_{\ell,i,\tilde{\mathbf{s}}_k} \tilde{\mathbf{B}}_{\ell,i,\tilde{\mathbf{s}}_k}, \tilde{\mathbf{B}}_{\ell,i,\tilde{\mathbf{s}}_k} | \sigma_{\mathbf{B}}^2 \sim MN_{n(\tilde{\mathbf{s}}_k)}(\mathbf{0}, \sigma_{\mathbf{B}}^2 \mathbf{I}_{n(\tilde{\mathbf{s}}_k)}) \end{aligned}$$

- ▶ The product forms  $\eta_{\ell,i,\tilde{\mathbf{s}}_k} \tilde{\mathbf{B}}_{\ell,i,\tilde{\mathbf{s}}_k}$  and  $\gamma_{\ell,i,i} \tilde{\mathbf{B}}_{\ell,i,i}$  are the same as the spike and slab priors.

- ▶  $\log(\mathbf{P}(\mathbf{Y})) = \mathbf{L}(\mathbf{q}) + KL(\mathbf{q}||\mathbf{P})$ , where  $\mathbf{L}(\mathbf{q}) = \int \sum_{\gamma} \sum_{\eta} -\mathbf{q}(\tilde{\mathbf{B}}, \gamma, \eta) \log \left( \frac{\mathbf{q}(\tilde{\mathbf{B}}, \gamma, \eta)}{\mathbf{P}(\tilde{\mathbf{B}}, \gamma, \eta, \mathbf{Y}|\mathbf{X}; \theta)} \right) d\tilde{\mathbf{B}}$ .

$$\min_{\mathbf{q}} KL(\mathbf{q}||\mathbf{P}) \iff \max_{\mathbf{q}} \mathbf{L}(\mathbf{q}).$$

- ▶ Based on prior assumptions, the proposed variational Bayesian algorithm iterates the following steps. E-step: take the expectation of  $\mathbf{L}(\mathbf{q})$  with respect to  $\eta, \gamma$  and  $\tilde{\mathbf{B}}$  to update  $\mathbf{q}(\eta, \gamma, \tilde{\mathbf{B}})$ .

$$\begin{aligned} \mathbf{q}(\eta, \gamma, \tilde{\mathbf{B}}) &= \prod_{\ell} \prod_i \prod_k \left( \phi_{\ell,i,i} \mathbf{N}(\mu_{\ell,i,i}, \Sigma_{\mathbf{B}_{\ell,i,i}}) \right)^{\gamma_{\ell,i,i}} \left( (1 - \phi_{\ell,i,i}) \mathbf{N}(\mathbf{0}, \sigma_{\mathbf{B}}^2) \right)^{(1-\gamma_{\ell,i,i})} \\ &\quad \left( \phi_{\ell,i,\tilde{\mathbf{s}}_k} MN_{1 \times n(\tilde{\mathbf{s}}_k)}(\mu_{\ell,i,\tilde{\mathbf{s}}_k}, \Sigma_{\mathbf{B}_{\ell,i,\tilde{\mathbf{s}}_k}}) \right)^{\eta_{\ell,i,\tilde{\mathbf{s}}_k}} \left( (1 - \phi_{\ell,i,\tilde{\mathbf{s}}_k}) MN_{1 \times n(\tilde{\mathbf{s}}_k)}(\mathbf{0}, \sigma_{\mathbf{B}}^2 \mathbf{I}_{n(\tilde{\mathbf{s}}_k)}) \right)^{(1-\eta_{\ell,i,\tilde{\mathbf{s}}_k})}, \end{aligned}$$

- ▶ M-step: take the derivative of  $\mathbf{L}(\mathbf{q})$  with respect to  $\theta = \{\hat{\Sigma}, \hat{\sigma}_{\mathbf{B}}^2, \hat{\pi}_1, \hat{\pi}_2\}$ .

## Simulation

UG: The true models follow the universal grouping structure with  $(m, p) = (10, 5), (20, 5)$  and  $(50, 5)$ , respectively. There are 72, 145 and 355 nonzero coefficients in three cases.

SG: The segmentation structure is considered in these three simulations with  $m = 10, 20$ , and 50. Let  $\mathbf{S}_{m,g}$  denote a specific group structure with  $m$  time series for  $g$  disjoint groups. We set  $\mathbf{S}_{10,3} = \{(1, 2, 3), (4, 5, 6), (7, 8, 9, 10)\}$  for  $(m, p) = (10, 5)$ ,  $\mathbf{S}_{20,4} = \{(1, \dots, 5), (6, \dots, 10), (11, 12, 13), (14, \dots, 20)\}$  for  $(m, p) = (20, 5)$ , and  $\mathbf{S}_{50,8} = \{(1, \dots, 5), (6, \dots, 10), (11, 12, 13), (14, \dots, 20), (21, \dots, 30), (31, \dots, 35), (36, \dots, 40), (41, \dots, 50)\}$  for  $(m, p) = (50, 5)$ . There are 40, 109, and 360 nonzero coefficients, respectively.

NG: In the "no grouping" cases, there are 18, 27 and 128 nonzero coefficients for  $(m, p) = (10, 5), (20, 5)$ , and  $(m, p) = (50, 5)$  respectively.

- ▶ The simulation structures with  $T = 301$  in  $m = 20$  and 10 and  $T = 701$  in  $m = 50$  are conducted according to Chu et al. [3]. The simulation results are summarized based on 100 replications.
- ▶ For the comparison purpose, we implement VAGSA [3] and BIVAS [5].
- ▶ Initial prior  $\pi_1$  and  $\pi_2$  for m50NG are 0.5, other cases are 0.01.

Table: Numerical comparison results for VB-NAR, VAGSA and BIVAS for the cases with  $p = 10$ .

VB-NAR/VAGSA/BIVAS	TPR(%)	FPR(%)	AMS	MSPE	Ave. CPU Time (sec)
m10UG- $\mathbf{I}_0$	100/100/98	0.07/0.31/1.66	72.62/74.90/86.02	1.00/1.01/0.98	4/148/8
m10UG- $\Sigma_{10}$	100/100/97	0.06/0.28/1.79	72.51/74.59/86.59	1.00/0.83/0.94	5/138/9
m10SG- $\mathbf{I}_0$	100/100/99	0.15/0.63/0.63	41.35/46.02/45.57	1.07/1.09/1.08	12/303/8
m10SG- $\Sigma_{10}$	100/100/98	0.13/0.26/0.74	41.17/46.02/46.16	0.86/1.09/0.90	11/303/9
m10NG- $\mathbf{I}_0$	98/97/96	0.15/0.11/0.26	19.07/18.58/19.91	1.00/1.00/1.02	40/628/16
m10NG- $\Sigma_{10}$	99/99/96	0.11/0.08/0.33	18.86/18.55/20.55	0.89/0.88/0.89	36/653/18
m20UG- $\mathbf{I}_0$	100/100/94	0.03/0.18/1.04	145.79/151.89/178.23	0.98/0.99/0.98	15/236/83
m20UG- $\Sigma_{20}$	100/100/97	0.02/0.14/0.63	145.56/150.54/45.57	0.95/0.95/1.00	180/237/93
m20SG- $\mathbf{I}_0$	100/100/89	0.06/0.26/0.48	109.25/117.05/115.81	1.10/1.11/1.17	44/652/63
m20SG- $\Sigma_{20}$	100/100/99	0.06/0.21/0.79	109.22/115.20/117.20	0.90/0.90/0.99	51/672/70
m20NG- $\mathbf{I}_0$	97/98/94	0.08/0.15/0.12	29.29/32.39/30.03	0.99/0.99/0.98	291/2183/144
m20NG- $\Sigma_{20}$	97/98/92	0.06/0.13/0.14	28.55/31.43/30.49	0.86/0.86/0.95	262/2247/151
m50UG- $\mathbf{I}_0$	100/-/97	0.00/-/0.26	355.47/-/409.03	1.02/-/2.66	213/-/769
m50UG- $\Sigma_{50}$	100/-/98	0.00/-/0.27	355.57/-/412.43	0.87/-/2.28	239/-/731
m50SG- $\mathbf{I}_0$	100/-/90	0.01/-/0.14	361.58/-/394.43	1.02/-/3.50	1082/-/645
m50SG- $\Sigma_{50}$	100/-/89	0.01/-/0.16	361.15/-/395.07	0.88/-/3.09	1767/-/663
m50NG- $\mathbf{I}_0$	100/-/98	0.03/-/0.07	134.67/-/143.61	1.02/-/1.92	7860/-/1163
m50NG- $\Sigma_{50}$	92/-/98	0.03/-/0.08	123.33/-/145.00	0.89/-/1.64	58194/-/1177

## Empirical Study

- ▶ Natural gas flows for the German network with 51 nodes from Oct 1st, 2013, to Sep 30th, 2015
- ▶ The segmented grouping structures based on the types of nodes that are used for the NAR analysis.
- ▶ We set the number of lags  $p$  as 14 to incorporate social dependence up to two weeks ahead.
- ▶ The data from Oct 1st, 2013, to March 31th, 2015, is used as the in-sample training data.
- ▶ We perform a one-step ahead forecast for remaining 0.5 years.

In order to illustrate the performance of the proposed VB method, in addition to the mean absolute percentage estimator (MAPE), the normalized MSE (NRMSE) for the 51 nodes is also used.

$$\begin{aligned} \text{MAPE} &= \frac{1}{N \times \# \text{ of type element}} \sum_t \sum_{i=1}^{\# \text{ of type element}} \frac{|Y_{t+1,i} - \hat{Y}_{t+1,i}|}{|Y_{t+1,i}|}, \\ \text{NRMSE} &= \frac{1}{N \times \# \text{ of type element}} \left( \sum_t \sum_{i=1}^{\# \text{ of type element}} (Y_{t+1,i} - \hat{Y}_{t+1,i})^2 \right)^{1/2} \\ &\quad \frac{1}{N \times \# \text{ of type element} \sum_t \sum_{i=1}^{\# \text{ of type element}} Y_{t+1,i}}. \end{aligned}$$

Table: The MAPEs and NRMSEs for the fitted and forecast results classified by type.

In sample from 01/10/2013 to 31/03/2015							
Type	Number	MAPE (%)	Range (%)	SD	NRMSE	Range	SD
Municipal	34	6.28	(3.20, 14.70)	1.93	0.09	(0.04, 0.36)	0.05
Industrial	11	11.69	(0.89, 34.22)	11.74	0.12	(0.01, 0.31)	0.10
Border	1	9.36	—	—	0.11	—	—
Others	5	13.05	(3.78, 36.30)	13.26	0.14	(0.05, 0.34)	0.12
Out of sample from 01/04/2015 to 30/09/2015							
Municipal	34	7.89	(4.66, 11.85)	1.60	0.10	(0.06, 0.15)	0.02
Industrial	11	15.98	(1.70, 54.57)	19.00	0.12	(0.02, 0.30)	0.10
Border	1	11.77	—	—	0.14	—	—
Others	5	8.65	(5.52, 13.47)	3.32	0.10	(0.07, 0.12)	0.02

## Summary

- ▶ We focus on Bayesian analysis for NAR/VAR models, especially when there are many time series are taken into considered simultaneously.
- ▶ The simulation results support that the proposed VB approach not only identifies the proper active structures in NAR/VAR models but also significantly reduces the computational cost.
- ▶ Finally, German gas flow network data with 51 nodes are analyzed to illustrate the performance of the proposed method. The analytical results of the proposed VB method yield the trends of nodes, which may be useful for assisting operators with performing appropriate operations.

## References

- Carbonetto P, Stephens M. (2012), Scalable variational inference for Bayesian variable selection in regression, and its accuracy in genetic association studies. *Bayesian Anal* 7: 73–108.
- Christopher M. Bishop. (2006), *Pattern Recognition and Machine Learning*. Springer-Verlag New York.
- Chi-Hsiang Chu, Mong-Na Lo Huang, Shih-Feng Huang, and Ray-Bing Chen. Bayesian structure selection for vector autoregression model. *Journal of Forecasting*, 38:422–439, 2019.
- George, E. I. and McCulloch, R. E. (1993), Variable selection via Gibbs sampling. *J. Am. Statist. Ass.*, 88, 881–889.
- Mingxuan Cai, Mingwei Dai, Jingsi Ming, Heng Peng, Jin Liu, and Can Yang. Bivas: a scalable bayesian method for bi-level variable selection with applications. *Journal of Computational and Graphical Statistics*, 29(1):40–52, 2020.
- Song, S. and Bickel, P. J. (2011), Large vector auto regressions. Preprint.
- Titsias, M. K. and Lazaro-Gredilla, M. (2011), Spike and slab variational inference for multi-task and multiple kernel learning. In NIPS.