

# CHAPTER 3 BAYESIAN METHODS FOR ESTIMATING STRUCTURAL EQUATION MODELS

Notations and concepts related to the Bayesian approach:

- $M$  — an arbitrary SEM with a vector of unknown parameters  $\theta$ .
- $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$  — the observed data set
- $p(\theta|M)$  (or  $p(\theta)$ ) — the prior density function of  $\theta$ .
- $p(\mathbf{Y}, \theta|M)$  — the probability density function of the joint distribution of  $\mathbf{Y}$  and  $\theta$  under  $M$ .
- $p(\theta|\mathbf{Y}, M)$  — the density function of the posterior distribution of  $\theta$  under  $M$ . This function fully describes the behavior of  $\theta$  under the given data  $\mathbf{Y}$ .
- $p(\mathbf{Y}|\theta, M)$  — the likelihood function under  $M$ .

Based on a well-known identity in probability, we have

$$p(\mathbf{Y}, \theta | M) = p(\mathbf{Y} | \theta, M) p(\theta) = p(\theta | \mathbf{Y}, M) p(\mathbf{Y} | M).$$

As  $p(\mathbf{Y} | M)$  does not depend on  $\theta$ , and can be regarded as a constant with fixed  $\mathbf{Y}$ , we have

$$p(\theta | \mathbf{Y}, M) \propto p(\mathbf{Y} | \theta, M) p(\theta), \quad \text{or} \\ \log p(\theta | \mathbf{Y}, M) = \log p(\mathbf{Y} | \theta, M) + \log p(\theta) + \text{constant}.$$

### Key steps in the Bayesian estimation:

1. Specify prior distributions.
2. Derive posterior distributions.
3. Sample from posterior distributions using MCMC algorithms.
4. Conduct sensitivity analysis.

# 1. Specify prior distributions.

The most commonly used prior distributions are informative priors (i.e.,  $p(\theta)$  and  $p(\theta|y)$  have the same form).

We consider the following model:

$$\mathbf{y}_i = \mathbf{\Lambda} \boldsymbol{\omega}_i + \boldsymbol{\epsilon}_i,$$

$$\boldsymbol{\eta}_i = \mathbf{B} \mathbf{d}_i + \mathbf{\Pi} \boldsymbol{\eta}_i + \mathbf{\Gamma} \boldsymbol{\xi}_i + \boldsymbol{\delta}_i = \mathbf{\Lambda}_\omega \mathbf{v}_i + \boldsymbol{\delta}_i,$$

where  $\boldsymbol{\omega}_i = (\boldsymbol{\eta}_i^T, \boldsymbol{\xi}_i^T)^T$ ,  $\mathbf{\Lambda}$ ,  $\mathbf{B}$ ,  $\mathbf{\Pi}$ ,  $\mathbf{\Gamma}$  are parameter matrices of unknown regression coefficients,  $\mathbf{\Lambda}_\omega = (\mathbf{B}, \mathbf{\Pi}, \mathbf{\Gamma})$ , and  $\mathbf{v}_i = (\mathbf{d}_i^T, \boldsymbol{\eta}_i^T, \boldsymbol{\xi}_i^T)^T$ . The distributions of  $\boldsymbol{\xi}_i$ ,  $\boldsymbol{\epsilon}_i$ , and  $\boldsymbol{\delta}_i$  are  $N[\mathbf{0}, \boldsymbol{\Phi}]$ ,  $N[\mathbf{0}, \boldsymbol{\Psi}_\epsilon]$ , and  $N[\mathbf{0}, \boldsymbol{\Psi}_\delta]$ .

Let  $\Lambda_k^T$  be the  $k$ th row of  $\Lambda$ ,  $\psi_{\epsilon k}$  be the  $k$ th diagonal element of  $\Psi_\epsilon$ ,  $\Lambda_{\omega k}^T$  be the  $k$ th row of  $\Lambda_\omega$ , and  $\psi_{\delta k}$  be the  $k$ th diagonal element of  $\Psi_\delta$ . Then, the conjugate prior distributions are as follows:

$$\psi_{\epsilon k} \stackrel{D}{=} IG[\alpha_{0\epsilon k}, \beta_{0\epsilon k}] \text{ or } \psi_{\epsilon k}^{-1} \stackrel{D}{=} Gamma[\alpha_{0\epsilon k}, \beta_{0\epsilon k}],$$

$$[\Lambda_k | \psi_{\epsilon k}] \stackrel{D}{=} N[\Lambda_{0k}, \psi_{\epsilon k} \mathbf{H}_{0yk}],$$

$$\Phi \stackrel{D}{=} IW_{q_2}[\mathbf{R}_0^{-1}, \rho_0], \text{ or equivalently } \Phi^{-1} \stackrel{D}{=} W_{q_2}[\mathbf{R}_0, \rho_0],$$

$$\psi_{\delta k} \stackrel{D}{=} IG[\alpha_{0\delta k}, \beta_{0\delta k}] \text{ or } \psi_{\delta k}^{-1} \stackrel{D}{=} Gamma[\alpha_{0\delta k}, \beta_{0\delta k}],$$

$$[\Lambda_{\omega k} | \psi_{\delta k}] \stackrel{D}{=} N[\Lambda_{0\omega k}, \psi_{\delta k} \mathbf{H}_{0\omega k}].$$

## Appendix 1: The Gamma, Inverted Gamma, Wishart, and Inverted Wishart distributions:

1. *Gamma distribution*:  $\theta \stackrel{D}{=} \text{Gamma}[\alpha, \beta]$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta},$$

$$E(\theta) = \alpha/\beta, \quad \text{Var}(\theta) = \alpha/\beta^2.$$

2. *Inverted Gamma distribution*:  $\theta \stackrel{D}{=} \text{Inverted Gamma}[\alpha, \beta]$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta},$$

$$E(\theta) = \frac{\beta}{\alpha - 1}, \quad \text{Var}(\theta) = \frac{\beta^2}{(\alpha - 1)^2(\alpha - 1)}.$$

3. *Relation between Gamma and inverted Gamma distributions*

If  $\theta \stackrel{D}{=} \text{Inverted Gamma}[\alpha, \beta]$ , then  $\theta^{-1} \stackrel{D}{=} \text{Gamma}[\alpha, \beta]$ .

4. *Wishart distribution*:  $\mathbf{W} \stackrel{D}{=} \text{Wishart}_q[\mathbf{R}_0, \rho_0]$   $\mathbf{R}_0 \in \mathbb{R}^{q \times q}$

$$p(\mathbf{W}) = \left[ 2^{\rho_0 q/2} \pi^{q(q-1)/4} \prod_{i=1}^q \Gamma\left(\frac{\rho_0 + 1 - i}{2}\right) \right]^{-1} \\ \times |\mathbf{R}_0|^{-\rho_0/2} \times |\mathbf{W}|^{(\rho_0 - q - 1)/2} \times \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{R}_0^{-1} \mathbf{W}) \right\},$$

$$E(\mathbf{W}) = \rho_0 \mathbf{R}_0.$$

5. *Inverted Wishart distribution*:  $\mathbf{W} \stackrel{D}{=} IW_q[\mathbf{R}_0^{-1}, \rho_0]$

$$p(\mathbf{W}) = \left[ 2^{\rho_0 q/2} \pi^{q(q-1)/4} \prod_{i=1}^q \Gamma\left(\frac{\rho_0 + 1 - i}{2}\right) \right]^{-1} \\ \times |\mathbf{R}_0|^{-\rho_0/2} \times |\mathbf{W}|^{-(\rho_0 + q + 1)/2} \times \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{R}_0^{-1} \mathbf{W}^{-1}) \right\},$$

$$E(\mathbf{W}) = \frac{\mathbf{R}_0^{-1}}{\rho_0 - q - 1}.$$

6. *Relation between Wishart and inverted Wishart distributions*

$$\text{If } \mathbf{W} \stackrel{D}{=} IW[\mathbf{R}_0^{-1}, \rho_0], \text{ then } \mathbf{W}^{-1} \stackrel{D}{=} W[\mathbf{R}_0, \rho_0].$$

The choice of hyperparameters:

1. If we have good prior information about  $\Lambda$  and  $\Lambda_\omega$ , then choose  $\Lambda_{0k}$  and  $\Lambda_{0\omega k}$  close to the true values of  $\Lambda_k$  and  $\Lambda_{\omega k}$ , along with small  $\mathbf{H}_{0yk}$  and  $\mathbf{H}_{0\omega k}$  (e.g.,  $0.5\mathbf{I}$  or  $\mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix).
2. If we have no good information about  $\Lambda$  and  $\Lambda_\omega$ , then choose  $\Lambda_{0k}$  and  $\Lambda_{0\omega k}$  at ad hoc values (e.g., zero elements), along with large  $\mathbf{H}_{0yk}$  and  $\mathbf{H}_{0\omega k}$  (e.g.,  $10\mathbf{I}$  or  $100\mathbf{I}$ ).
3. For  $\alpha_{0\epsilon k}$  and  $\beta_{0\epsilon k}$  (similarly  $\alpha_{0\delta k}$  and  $\beta_{0\delta k}$ ), we may assign their values according to the rationale: If  $\psi_{\epsilon k} \stackrel{D}{=} \text{Inverted Gamma}(\alpha_{0\epsilon k}, \beta_{0\epsilon k})$ , then

$$E(\psi_{\epsilon k}) = \beta_{0\epsilon k} / (\alpha_{0\epsilon k} - 1),$$

$$\text{Var}(\psi_{\epsilon k}) = \beta_{0\epsilon k}^2 / \{(\alpha_{0\epsilon k} - 1)^2(\alpha_{0\epsilon k} - 2)\}.$$

4. For  $\mathbf{R}_0$  and  $\rho_0$ , since  $E(\Phi) = \mathbf{R}_0^{-1} / (\rho_0 - q_2 - 1)$ , we may choose  $\mathbf{R}_0^{-1}$  and  $\rho_0$  such that

$$\mathbf{R}_0^{-1} = (\rho_0 - q_2 - 1)\Phi_0.$$



## 2. Derive posterior distributions.

(1) *Conditional Distribution*  $[\Omega|\mathbf{Y}, \theta]$ 

For  $i = 1, \dots, n$ , as  $\omega_i$  are conditionally independent and  $\mathbf{y}_i$  are also conditionally independent, then

$$\begin{cases} \mathbf{y} = \Lambda \boldsymbol{\omega} + \boldsymbol{\varepsilon} \\ \boldsymbol{\eta} = \mathbf{B} \boldsymbol{\delta} + \boldsymbol{\Pi} \boldsymbol{\eta} + \boldsymbol{\Gamma} \boldsymbol{\xi} + \boldsymbol{\delta} \\ \omega_i = \begin{bmatrix} \eta_i \\ \xi_i \end{bmatrix} \\ \text{Cov}(\omega_i) = \begin{bmatrix} \text{Var} \eta_i & \text{Cov}(\eta_i, \xi_i) \\ \text{Cov}(\xi_i, \eta_i) & \text{Var} \xi_i \end{bmatrix} \end{cases}$$

$$p(\Omega|\mathbf{Y}, \theta) = \prod_{i=1}^n p(\omega_i|\mathbf{y}_i, \theta).$$

It implies that the conditional distributions of  $\omega_i$  given  $(\mathbf{y}_i, \theta)$  are mutually independent for different  $i$ , and

$$p(\omega_i|\mathbf{y}_i, \theta) = \frac{p(\omega_i, \mathbf{y}_i|\theta)}{p(\mathbf{y}_i|\theta)} \propto p(\omega_i|\theta)p(\mathbf{y}_i|\omega_i, \theta).$$

Let  $\boldsymbol{\Pi}_0 = \mathbf{I} - \boldsymbol{\Pi}$  and the covariance matrix of  $\omega_i$  be

$$\boldsymbol{\Sigma}_{\omega} = \begin{bmatrix} \boldsymbol{\Pi}_0^{-1}(\boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}^T + \boldsymbol{\Psi}_{\delta})\boldsymbol{\Pi}_0^{-T} & \boldsymbol{\Pi}_0^{-1}\boldsymbol{\Gamma}\boldsymbol{\Phi} \\ \boldsymbol{\Phi}\boldsymbol{\Gamma}^T\boldsymbol{\Pi}_0^{-T} & \boldsymbol{\Phi} \end{bmatrix}.$$

It can be shown that

$$[\omega_i | \theta] \stackrel{D}{=} N \left[ \begin{pmatrix} \Pi_0^{-1} \mathbf{B} \mathbf{d}_i \\ \mathbf{0} \end{pmatrix}, \Sigma_\omega \right],$$

and

$$[\mathbf{y}_i | \omega_i, \theta] \stackrel{D}{=} N[\Lambda \omega_i, \Psi_\epsilon].$$

Let  $\mu_\omega = \begin{pmatrix} \Pi_0^{-1} \mathbf{B} \mathbf{d}_i \\ \mathbf{0} \end{pmatrix}$ . Then, we have

$$\begin{aligned} p(\omega_i | \mathbf{y}_i, \theta) &\propto p(\mathbf{y}_i | \omega_i, \theta) p(\omega_i | \theta) \\ &\propto \exp \left\{ -\frac{1}{2} (\mathbf{y}_i - \Lambda \omega_i)^T \Psi_\epsilon^{-1} (\mathbf{y}_i - \Lambda \omega_i) \right\} \times \\ &\quad \exp \left\{ -\frac{1}{2} (\omega_i - \mu_\omega)^T \Sigma_\omega^{-1} (\omega_i - \mu_\omega) \right\} \\ &\propto \exp \left\{ -\frac{1}{2} [\mathbf{y}_i^T \Psi_\epsilon^{-1} \mathbf{y}_i - 2 \omega_i^T \Lambda^T \Psi_\epsilon^{-1} \mathbf{y}_i + \omega_i^T (\Lambda^T \Psi_\epsilon^{-1} \Lambda) \omega_i \right. \\ &\quad \left. + \omega_i^T \Sigma_\omega^{-1} \omega_i - 2 \omega_i^T \Sigma_\omega^{-1} \mu_\omega] \right\} \end{aligned}$$

$$\begin{aligned}
&\propto \exp \left\{ -\frac{1}{2} \left[ -2\boldsymbol{\omega}_i^T (\boldsymbol{\Lambda}^T \boldsymbol{\Psi}_\epsilon^{-1} \mathbf{y}_i + \boldsymbol{\Sigma}_\omega^{-1} \boldsymbol{\mu}_\omega) + \right. \right. \\
&\quad \left. \left. \boldsymbol{\omega}_i^T (\boldsymbol{\Sigma}_\omega^{-1} + \boldsymbol{\Lambda}^T \boldsymbol{\Psi}_\epsilon^{-1} \boldsymbol{\Lambda}) \boldsymbol{\omega}_i \right] \right\} \\
&\propto \exp \left\{ -\frac{1}{2} \left[ \boldsymbol{\omega}_i - \boldsymbol{\Sigma}^{*-1} (\boldsymbol{\Lambda}^T \boldsymbol{\Psi}_\epsilon^{-1} \mathbf{y}_i + \boldsymbol{\Sigma}_\omega^{-1} \boldsymbol{\mu}_\omega) \right]^T \boldsymbol{\Sigma}^* \right. \\
&\quad \left. \left[ \boldsymbol{\omega}_i - \boldsymbol{\Sigma}^{*-1} (\boldsymbol{\Lambda}^T \boldsymbol{\Psi}_\epsilon^{-1} \mathbf{y}_i + \boldsymbol{\Sigma}_\omega^{-1} \boldsymbol{\mu}_\omega) \right] \right\}.
\end{aligned}$$

Thus,

$$[\boldsymbol{\omega}_i | \mathbf{y}_i, \boldsymbol{\theta}] \stackrel{D}{=} N \left[ \boldsymbol{\Sigma}^{*-1} \boldsymbol{\Lambda}^T \boldsymbol{\Psi}_\epsilon^{-1} \mathbf{y}_i + \boldsymbol{\Sigma}^{*-1} \boldsymbol{\Sigma}_\omega^{-1} \begin{pmatrix} \boldsymbol{\Pi}_0^{-1} \mathbf{B} \mathbf{d}_i \\ \mathbf{0} \end{pmatrix}, \boldsymbol{\Sigma}^{*-1} \right],$$

where  $\boldsymbol{\Sigma}^* = \boldsymbol{\Sigma}_\omega^{-1} + \boldsymbol{\Lambda}^T \boldsymbol{\Psi}_\epsilon^{-1} \boldsymbol{\Lambda}$ . The conditional distribution  $[\boldsymbol{\omega}_i | \mathbf{y}_i, \boldsymbol{\theta}]$  is a normal distribution.

(2) *Conditional Distribution*  $[\theta|\mathbf{Y}, \Omega]$ 

Note that

$$p(\theta|\mathbf{Y}, \Omega) \propto p(\theta)p(\mathbf{Y}, \Omega|\theta).$$

Let  $\theta_y$  be the unknown parameters in  $\Lambda$  and  $\Psi_\epsilon$  associated with the measurement equation, and  $\theta_\omega$  be the unknown parameters in  $\mathbf{B}$ ,  $\Pi$ ,  $\Gamma$ ,  $\Phi$ , and  $\Psi_\delta$  associated with the structural equation.

Assumed that the prior distribution of  $\theta_y$  is independent of the prior distribution of  $\theta_\omega$ , that is,  $p(\theta) = p(\theta_y)p(\theta_\omega)$ . Note also that

$$p(\mathbf{Y}|\Omega, \theta) = p(\mathbf{Y}|\Omega, \theta_y), \quad p(\Omega|\theta) = p(\Omega|\theta_\omega)$$

and

$$p(\theta_y|\mathbf{Y}, \Omega) \propto p(\mathbf{Y}|\Omega, \theta_y)p(\theta_y), \quad p(\theta_\omega|\mathbf{Y}, \Omega) \propto p(\Omega|\theta_\omega)p(\theta_\omega),$$

Hence, these conditional densities can be treated separately.

(a)  $p(\theta_y|\mathbf{Y}, \boldsymbol{\Omega})$  or  $p(\boldsymbol{\Lambda}, \psi_{\epsilon k}|\mathbf{Y}, \boldsymbol{\Omega})$ :

For notation simplicity, we temporarily assume that all the elements of  $\boldsymbol{\Lambda}_k$  are unknown. Let  $\nu_k = \psi_{\epsilon k}^{-1}$ , we have

$$p(\nu_k) \propto \nu_k^{\alpha_{0\epsilon k}-1} \exp(-\beta_{0\epsilon k} \nu_k),$$

$$p(\boldsymbol{\Lambda}_k|\nu_k) \propto \nu_k^{q/2} \exp\left\{-\frac{1}{2}(\boldsymbol{\Lambda}_k - \boldsymbol{\Lambda}_{k0})^T \mathbf{H}_{0yk}^{-1}(\boldsymbol{\Lambda}_k - \boldsymbol{\Lambda}_{k0})\nu_k\right\},$$

and the likelihood of  $\mathbf{Y}$  is

$$p(\mathbf{Y}|\theta_y, \boldsymbol{\Omega}) \propto |\boldsymbol{\Psi}_\epsilon|^{-n/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\Lambda} \boldsymbol{\omega}_i)^T \boldsymbol{\Psi}_\epsilon^{-1} (\mathbf{y}_i - \boldsymbol{\Lambda} \boldsymbol{\omega}_i)\right\}.$$

Let  $\mathbf{Y}_k^T$  be the  $k$ th row of  $\mathbf{Y}$ ,  $y_{ik}$  be the  $i$ th component of  $\mathbf{Y}_k^T$ ,  $\mathbf{A}_k^* = (\Omega\Omega^T)^{-1}\Omega\mathbf{Y}_k$ , and  $b_k = \mathbf{Y}_k^T\mathbf{Y}_k - \mathbf{Y}_k^T\Omega^T(\Omega\Omega^T)^{-1}\Omega\mathbf{Y}_k = \mathbf{Y}_k^T\mathbf{Y}_k - \mathbf{A}_k^{*T}(\Omega\Omega^T)\mathbf{A}_k^*$ . The exponential term in  $p(\mathbf{Y}|\theta_y, \Omega)$  can be expressed as

$$\begin{aligned}
& -\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \Lambda\omega_i)^T \Psi_\epsilon^{-1} (\mathbf{y}_i - \Lambda\omega_i) = -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^p \nu_k (y_{ik} - \Lambda_k^T \omega_i)^2 \\
& = -\frac{1}{2} \sum_{k=1}^p \left\{ \nu_k \left[ \sum_{i=1}^n y_{ik}^2 - 2\Lambda_k^T \sum_{i=1}^n y_{ik} \omega_i + \text{tr}(\Lambda_k \Lambda_k^T \sum_{i=1}^n \omega_i \omega_i^T) \right] \right\} \\
& = -\frac{1}{2} \sum_{k=1}^p \left\{ \nu_k [\mathbf{Y}_k^T \mathbf{Y}_k - 2\Lambda_k^T \Omega \mathbf{Y}_k + \Lambda_k^T (\Omega\Omega^T) \Lambda_k] \right\} \\
& = -\frac{1}{2} \sum_{k=1}^p \left\{ \nu_k [\mathbf{Y}_k^T \mathbf{Y}_k - \mathbf{Y}_k^T \Omega^T (\Omega\Omega^T)^{-1} \Omega \mathbf{Y}_k] + \right. \\
& \quad \left. \nu_k [\Lambda_k - (\Omega\Omega^T)^{-1} \Omega \mathbf{Y}_k]^T (\Omega\Omega^T) [\Lambda_k - (\Omega\Omega^T)^{-1} \Omega \mathbf{Y}_k] \right\} \\
& = -\frac{1}{2} \sum_{k=1}^p \left\{ \nu_k [b_k + (\Lambda_k - \mathbf{A}_k^*)^T (\Omega\Omega^T) (\Lambda_k - \mathbf{A}_k^*)] \right\}.
\end{aligned}$$

Thus,

$$\begin{aligned}
 p(\theta_y | \mathbf{Y}, \Omega) &= p(\Psi_\epsilon) p(\Lambda | \Psi_\epsilon) p(\mathbf{Y} | \theta_y, \Omega) \\
 &\propto \prod_{k=1}^p \left\{ \nu_k^{n/2+q/2+\alpha_{0\epsilon k}-1} \exp \left[ -\frac{1}{2} \nu_k \{ (\Lambda_k - \mathbf{A}_k^*)^T (\Omega \Omega^T) (\Lambda_k - \mathbf{A}_k^*) \right. \right. \\
 &\quad \left. \left. + (\Lambda_k - \Lambda_{0k})^T \mathbf{H}_{0yk}^{-1} (\Lambda_k - \Lambda_{0k}) \} - \nu_k (\beta_{0\epsilon k} + b_k/2) \right] \right\} \\
 &= \prod_{k=1}^p p(\Lambda_k, \nu_k | \mathbf{Y}, \Omega).
 \end{aligned}$$

To derive  $p(\theta_y | \mathbf{Y}, \Omega)$ , it is sufficient to derive  $p(\Lambda_k, \nu_k | \mathbf{Y}, \Omega)$ . Let  $\mathbf{A}_k = (\mathbf{H}_{0yk}^{-1} + \Omega \Omega^T)^{-1}$  and  $\mathbf{a}_k = \mathbf{A}_k (\mathbf{H}_{0yk}^{-1} \Lambda_{0k} + \Omega \mathbf{Y}_k)$ , it follows that

$$\begin{aligned}
 &(\Lambda_k - \mathbf{A}_k^*)^T (\Omega \Omega^T) (\Lambda_k - \mathbf{A}_k^*) + (\Lambda_k - \Lambda_{0k})^T \mathbf{H}_{0yk}^{-1} (\Lambda_k - \Lambda_{0k}) \\
 &= (\Lambda_k - \mathbf{a}_k)^T \mathbf{A}_k^{-1} (\Lambda_k - \mathbf{a}_k) - \mathbf{a}_k^T \mathbf{A}_k^{-1} \mathbf{a}_k \\
 &\quad + \mathbf{A}_k^{*T} \Omega \Omega^T \mathbf{A}_k^* + \Lambda_{0k}^T \mathbf{H}_{0yk}^{-1} \Lambda_{0k}.
 \end{aligned}$$

Hence,

$$p(\mathbf{\Lambda}_k, \nu_k | \mathbf{Y}, \mathbf{\Omega}) = p(\nu_k | \mathbf{Y}, \mathbf{\Omega}) p(\mathbf{\Lambda}_k | \mathbf{Y}, \mathbf{\Omega}, \nu_k) \\ \propto \left\{ \nu_k^{n/2 + \alpha_{0\epsilon k} - 1} \exp(-\beta_{\epsilon k} \nu_k) \right\} \left\{ \nu_k^{q/2} \exp \left[ -\frac{1}{2} (\mathbf{\Lambda}_k - \mathbf{a}_k)^T \mathbf{A}_k^{-1} (\mathbf{\Lambda}_k - \mathbf{a}_k) \nu_k \right] \right\}$$

where  $\beta_{\epsilon k} = \beta_{0\epsilon k} + \frac{1}{2} (\mathbf{Y}_k^T \mathbf{Y}_k - \mathbf{a}_k^T \mathbf{A}_k^{-1} \mathbf{a}_k + \mathbf{\Lambda}_{0k}^T \mathbf{H}_{0yk}^{-1} \mathbf{\Lambda}_{0k})$ .

Finally, we can conclude that  $p(\mathbf{\Lambda}_k, \nu_k | \mathbf{Y}, \mathbf{\Omega})$  is the following normal-gamma distribution:

$$[\nu_k | \mathbf{Y}, \mathbf{\Omega}] \stackrel{D}{=} \text{Gamma}[n/2 + \alpha_{0\epsilon k}, \beta_{\epsilon k}], \quad \text{and} \\ [\mathbf{\Lambda}_k | \mathbf{Y}, \mathbf{\Omega}, \nu_k] \stackrel{D}{=} N[\mathbf{a}_k, \nu_k^{-1} \mathbf{A}_k],$$



If some elements of  $\Lambda_k$  are fixed, we identify the positions of the fixed elements via an index matrix  $\mathbf{L}$  with the following elements:

$$l_{kj} = \begin{cases} 0, & \text{if } \lambda_{kj} \text{ is fixed,} \\ 1, & \text{if } \lambda_{kj} \text{ is free;} \end{cases} \quad \text{for } j = 1, \dots, q \text{ and } k = 1, \dots, p.$$

Let  $\Lambda_k^*$  be a vector of unknown parameters in  $\Lambda_k$ ,  $\mathbf{Y}_k$  be the submatrix of  $\mathbf{Y}$  such that all the rows corresponding to  $l_{kj} = 0$  are deleted; and let  $\mathbf{Y}_k^{*T} = (y_{1k}^*, \dots, y_{nk}^*)$  with

$$y_{ik}^* = y_{ik} - \sum_{j=1}^q \lambda_{kj} y_{ij} (1 - l_{kj}).$$

where  $y_{ij}$  is the  $j$ th element of  $\mathbf{y}_i$ . Then,

$$[\nu_k | \mathbf{Y}, \Omega] \stackrel{D}{=} \text{Gamma}[n/2 + \alpha_{0\epsilon k}, \beta_{\epsilon k}], \quad [\Lambda_k^* | \mathbf{Y}, \Omega, \nu_k] \stackrel{D}{=} N[\mathbf{a}_k, \nu_k \mathbf{A}_k],$$

where  $\mathbf{A}_k = (\mathbf{H}_{0yk}^{-1} + \mathbf{Y}_k \mathbf{Y}_k^T)^{-1}$ ,  $\mathbf{a}_k = \mathbf{A}_k (\mathbf{H}_{0yk}^{-1} \Lambda_{0yk} + \Omega \mathbf{Y}_k^*)$ , and  $\beta_{\epsilon k} = \beta_{0\epsilon k} + \frac{1}{2} (\mathbf{Y}_k^{*T} \mathbf{Y}_k^* - \mathbf{a}_k^T \mathbf{A}_k^{-1} \mathbf{a}_k + \Lambda_{0k}^T \mathbf{H}_{0yk}^{-1} \Lambda_{0k})$ .

(b)  $p(\theta_\omega | \mathbf{Y}, \Omega)$

Note that

$$p(\theta_\omega | \mathbf{Y}, \Omega) \propto p(\Omega | \theta_\omega) p(\theta_\omega).$$

Let  $\Omega_1 = (\eta_1, \dots, \eta_n)$  and  $\Omega_2 = (\xi_1, \dots, \xi_n)$ . Since the distribution of  $\xi_i$  only involves  $\Phi$ ,  $p(\Omega_2 | \theta_\omega) = p(\Omega_2 | \Phi)$ . Under the assumption that the prior distribution of  $\Phi$  is independent of the prior distributions of  $\mathbf{B}, \Pi, \Gamma$ , and  $\Psi_\delta$ , we have

$$p(\Omega | \theta_\omega) p(\theta_\omega) = [p(\Omega_1 | \Omega_2, \mathbf{B}, \Pi, \Gamma, \Psi_\delta) p(\mathbf{B}, \Pi, \Gamma, \Psi_\delta)] [p(\Omega_2 | \Phi) p(\Phi)].$$

Hence, the conditional distributions of  $(\mathbf{B}, \Pi, \Gamma, \Psi_\delta)$  and  $\Phi$  can be treated separately.

Consider a conjugate type prior distribution for  $\Phi$  with  $\Phi \stackrel{D}{=} IW_{q_2}[\mathbf{R}_0^{-1}, \rho_0]$  or  $\Phi^{-1} \stackrel{D}{=} W_{q_2}[\mathbf{R}_0, \rho_0]$ , with hyperparameters  $\rho_0$  and  $\mathbf{R}_0^{-1}$  or  $\mathbf{R}_0$ .

To derive  $p(\Phi|\Omega_2)$ , we first note that it is proportional to  $p(\Phi)p(\Omega_2|\Phi)$ . As  $\xi_i$  are independent, we have

$$p(\Phi|\Omega_2) \propto p(\Phi) \prod_{i=1}^n p(\xi_i|\theta).$$

Moreover, since the distribution of  $\xi_i$  given  $\Phi$  is  $N(\mathbf{0}, \Phi)$ , we have

$$\begin{aligned} p(\Phi|\Omega_2) &\propto \left[ |\Phi|^{-(\rho_0+q_2+1)/2} \exp \left\{ -\frac{1}{2} \text{tr}[\mathbf{R}_0^{-1}\Phi^{-1}] \right\} \right] \times \\ &\quad \left[ |\Phi|^{-n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \xi_i^T \Phi^{-1} \xi_i \right\} \right] \\ &= |\Phi|^{-(n+\rho_0+q_2+1)/2} \exp \left\{ -\frac{1}{2} \text{tr}[\Phi^{-1}(\Omega_2\Omega_2^T + \mathbf{R}_0^{-1})] \right\}. \end{aligned}$$

*invert Wishart dist.*

Since the right hand side is proportional to the density function of an inverted Wishart distribution (Zellner, 1971), it follows that the conditional distribution of  $\Phi$  given  $\Omega_2$  is given by

$$[\Phi|\Omega_2] \stackrel{D}{=} IW_{q_2}[(\Omega_2\Omega_2^T + \mathbf{R}_0^{-1}), n + \rho_0]$$

Recall that

$$\boldsymbol{\eta}_i = \boldsymbol{\Lambda}_\omega \mathbf{v}_i + \boldsymbol{\delta}_i,$$

where  $\boldsymbol{\Lambda}_\omega = (\mathbf{B}, \boldsymbol{\Pi}, \boldsymbol{\Gamma})$  with general elements  $\lambda_{\omega kj}$  for  $k = 1, \dots, q_1$ , and  $\mathbf{v}_i = (\mathbf{d}_i^T, \boldsymbol{\eta}_i^T, \boldsymbol{\xi}_i^T)^T = (\mathbf{d}_i^T, \boldsymbol{\omega}_i^T)^T$  be an  $(r_2 + q_1 + q_2) \times 1$  vector. The model  $\boldsymbol{\eta}_i = \boldsymbol{\Lambda}_\omega \mathbf{v}_i + \boldsymbol{\delta}_i$  is similar to  $\mathbf{y}_i = \boldsymbol{\Lambda} \boldsymbol{\omega}_i + \boldsymbol{\epsilon}_i$  considered before. Hence, the derivation for the conditional distributions corresponding to  $\boldsymbol{\theta}_\omega$  is similar to that corresponding to  $\boldsymbol{\theta}_y$ .

Let  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$ ,  $\mathbf{L}_\omega$  be the index matrix with general elements  $l_{\omega kj}$  that similarly defined as  $\mathbf{L}$  to indicate the fixed known parameters in  $\boldsymbol{\Lambda}_\omega$ ;  $\psi_{\delta k}$  be the  $k$ th diagonal element of  $\boldsymbol{\Psi}_\delta$  and  $\boldsymbol{\Lambda}_{\omega k}^T$  be the row vector that contains the unknown parameters in the  $k$ th row of  $\boldsymbol{\Lambda}_\omega$ . Let  $\mathbf{V}_k$  be the submatrix of  $\mathbf{V}$  such that all the rows corresponding to  $l_{\omega kj} = 0$  are deleted; and let  $\boldsymbol{\Xi}_k^T = (\eta_{1k}^*, \dots, \eta_{nk}^*)$  where

$$\eta_{ik}^* = \eta_{ik} - \sum_{j=1}^{r_2+q} \lambda_{\omega kj} v_{ij} (1 - l_{\omega kj}).$$

It can be shown that

$$[\psi_{\delta k}^{-1}|\Omega] \stackrel{D}{=} \text{Gamma}[n/2 + \alpha_{0\delta k}, \beta_{\delta k}],$$

$$[\Lambda_{\omega k}|\Omega, \psi_{\delta k}^{-1}] \stackrel{D}{=} N[\mathbf{a}_{\omega k}, \psi_{\delta k} \mathbf{A}_{\omega k}],$$

where  $\mathbf{A}_{\omega k} = (\mathbf{H}_{0\omega k}^{-1} + \mathbf{V}_k \mathbf{V}_k^T)^{-1}$ ,  $\mathbf{a}_{\omega k} = \mathbf{A}_{\omega k}(\mathbf{H}_{0\omega k}^{-1} \Lambda_{0\omega k} + \mathbf{V}_k \Xi_k)$ , and

$$\beta_{\delta k} = \beta_{0\delta k} + \frac{1}{2}(\Xi_k^T \Xi_k - \mathbf{a}_{\omega k}^T \mathbf{A}_{\omega k}^{-1} \mathbf{a}_{\omega k} + \Lambda_{0\omega k}^T \mathbf{H}_{0\omega k}^{-1} \Lambda_{0\omega k}).$$

Again, the conditional distribution

$$[\theta_{\omega}|\Omega] = \prod_{k=1}^{q_1} [\psi_{\delta k}^{-1}|\Omega][\Lambda_{\omega k}|\Omega, \psi_{\delta k}^{-1}]$$

is a normal-Gamma distribution.

### 3. Sample from posterior distribution using MCMC methods.

#### (1) Data augmentation:

Instead of working on the intractable posterior density  $p(\theta|\mathbf{Y})$ , we work on  $p(\theta, \Omega|\mathbf{Y})$ , where  $\Omega$  is the set of latent variables in the model.

*Given the latent  $\Omega$  to form complete data  $(\mathbf{Y}, \Omega)$*

#### (2) The Gibbs sampler:

The Gibbs sampler (Geman and Geman, 1984) is used to simulate observations from  $p(\theta, \Omega|\mathbf{Y})$  by drawing observations iteratively from their full conditional densities  $p(\theta|\Omega, \mathbf{Y})$  and  $p(\Omega|\theta, \mathbf{Y})$ .

#### (3) Check convergence:

The burn-in iterations for achieving convergence of the Gibbs sampler can be determined by plots of the simulated sequences of the individual parameters. At convergence, parallel sequences generated with different starting values should mix well together. **EPSR** values can also be used to check convergence of MCMC algorithms.

The 'estimated potential scale reduction (EPSR)' value is used to monitor the convergence of MCMC algorithm. As suggested by Gelmen (1996), convergence is achieved when the EPSR values are all less than 1.2. The computation of the EPSR values is based on several simulation sequences generated independently from different starting points. The EPSR approach monitors each parameter of interest separately.

Let  $n$  be the length of the each sequence. For each parameter estimate, say  $\theta$ , let  $\theta_{jk}$  ( $j = 1, \dots, n; k = 1, \dots, K$ ) be the observations from  $K$  parallel sequences of length  $n$ . The between- and within-sequence variances are computed as

*\* Introduced in PML: Intro.*

$$B = \frac{n}{K-1} \sum_{k=1}^K (\theta_{.k} - \theta_{..})^2, \quad \theta_{.k} = n^{-1} \sum_{j=1}^n \theta_{jk}, \quad \theta_{..} = K^{-1} \sum_{k=1}^K \theta_{.k},$$

$$W = \frac{1}{K} \sum_{k=1}^K s_k^2, \quad s_k^2 = (n-1)^{-1} \sum_{j=1}^n (\theta_{jk} - \theta_{.k})^2.$$

The estimate of  $\text{var}(\theta|\mathbf{Y})$ , the marginal posterior variance of the estimate, is then obtained by a weighted average of  $B$  and  $W$  as follows:

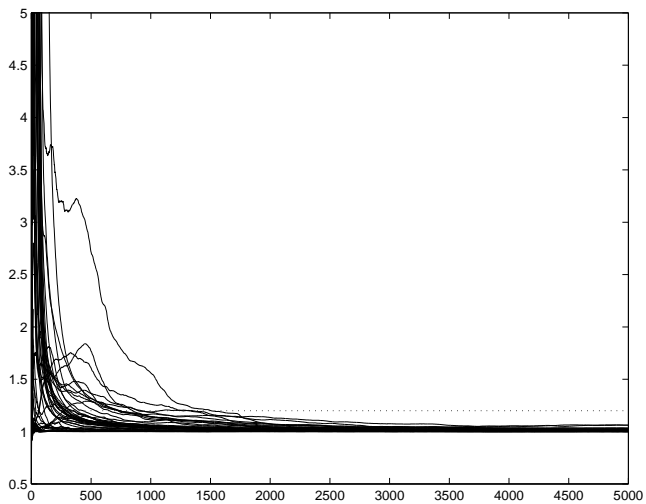
$$\widehat{\text{var}}(\theta) = \frac{n-1}{n}W + \frac{1}{n}B.$$

The EPSR is defined as

$$\hat{R}^{1/2} = [\widehat{\text{var}}(\theta)/W]^{1/2}.$$

As the algorithm converges,  $\hat{R}^{1/2}$  should be close to 1.0. In monitoring convergence, all EPSR values for all parameters are computed. Convergence is achieved if all the EPSR values are less than 1.2.





## (4) Statistical inference:

Statistical inference of the model can be conducted on the basis of a simulated sample of observations from  $p(\boldsymbol{\theta}, \boldsymbol{\Omega} | \mathbf{Y})$ , namely,  $\{(\boldsymbol{\theta}^{(t)}, \boldsymbol{\Omega}^{(t)}) : t = 1, \dots, T^*\}$ . The Bayesian estimate of  $\boldsymbol{\theta}$  and its standard error estimate can be obtained as follows:

$$\hat{\boldsymbol{\theta}} = T^{*-1} \sum_{t=1}^{T^*} \boldsymbol{\theta}^{(t)}, \quad \hat{\omega}_i = T^{*-1} \sum_{t=1}^{T^*} \omega_i^{(t)}.$$

$$\text{Var}(\hat{\boldsymbol{\theta}} | \mathbf{Y}) = (T^* - 1)^{-1} \sum_{t=1}^{T^*} (\boldsymbol{\theta}^{(t)} - \hat{\boldsymbol{\theta}})(\boldsymbol{\theta}^{(t)} - \hat{\boldsymbol{\theta}})^T,$$

Other statistical inference on  $\boldsymbol{\theta}$  can be carried out based on the simulated sample,  $\{\boldsymbol{\theta}^{(t)} : t = 1, \dots, T^*\}$ . For instance, the 2.5% and 97.5% quantiles of the sampled distribution of an individual parameter can give a 95% posterior credible interval and convey skewness in its marginal posterior density.

$\Rightarrow$  Significant or not.

#### 4. Sensitivity analysis.

A sensitivity analysis is usually necessary in the Bayesian analysis. The main objective is to check whether Bayesian results are sensitive to given prior inputs. It can be conducted by disturbing prior inputs in certain ranges.

Based on our experience, Bayesian results are

- (1) relatively robust to the prior inputs for regression type parameters;
- (2) relatively sensitive to the prior inputs for variance and covariance parameters.

 $\varphi_e \varphi_\beta$ 
 $\Lambda \Lambda_0$ 

Developing robust Bayesian methods and less sensitive prior distributions has attracted continuous attention in the statistical literature.