

1 Model and Likelihood

Consider a linear K -factor model for D assets of the form

$$y_{it} = \alpha_d + \mathbf{f}'_t \boldsymbol{\beta}_d + \varepsilon_{it}$$

where $d = 1, \dots, D$ and $t = 1, \dots, T$ and $\mathbf{f}'_t = (f_{t1}, \dots, f_{tK})$ is a $K \times 1$ vector. This is a special case of the seemingly unrelated regression (SUR) model in which the regressors are *identical* across equations. Stacking observations for a given time period across assets, define $\mathbf{y}'_t = (y_{1t}, \dots, y_{Dt})$ and analogously $\boldsymbol{\varepsilon}'_t = (\varepsilon_{t1}, \dots, \varepsilon_{tD})$. Now let $\mathbf{x}'_t = (1, \mathbf{f}'_t)$ and $\boldsymbol{\gamma}'_d = (\alpha_d, \boldsymbol{\beta}'_d)$ so we have

$$\mathbf{y}_t = X_t \boldsymbol{\gamma} + \boldsymbol{\varepsilon}_t$$

where $X_t = I_D \otimes \mathbf{x}'_t$ and $\boldsymbol{\gamma}' = (\boldsymbol{\gamma}'_1, \dots, \boldsymbol{\gamma}'_D)$. Now, suppose that

$$\boldsymbol{\varepsilon}_t | \mathbf{x}_t \sim \text{iid } \mathcal{N}_D(0, \Omega^{-1})$$

Let Y_T denote the full data sample, i.e. $\{\mathbf{y}_t, \mathbf{x}_t\}_{t=1}^T$. Then the likelihood is

$$\pi(\boldsymbol{\gamma}, \Omega^{-1} | Y_T) \propto |\Omega^{-1}|^{T/2} \exp \left[-\frac{1}{2} \sum_{t=1}^T (\mathbf{y}_t - X_t \boldsymbol{\gamma})' \Omega^{-1} (\mathbf{y}_t - X_t \boldsymbol{\gamma}) \right]$$

where we parameterize this problem in terms of the $D \times D$ *precision* matrix Ω^{-1} and the $p \times 1$ vector of regression coefficients $\boldsymbol{\gamma}$, where $p = D(K + 1)$.

2 Prior and Posterior Distribution

To complete the model we specify the following prior distribution

$$\pi(\boldsymbol{\gamma}, \Omega^{-1}) = \mathcal{N}_p(\boldsymbol{\gamma} | \boldsymbol{\gamma}_0, G_0) \mathcal{W}_D(\Omega^{-1} | \rho_0, R_0)$$

This prior is conditionally conjugate with the normal likelihood. In particular, we have $\boldsymbol{\gamma} | \Omega^{-1}, Y_T \sim \mathcal{N}_p(\bar{\boldsymbol{\gamma}}, G_T)$ where

$$\begin{aligned} G_T &= \left[G_0^{-1} + \sum_{t=1}^T X_t' \Omega^{-1} X_t \right]^{-1} \\ \bar{\boldsymbol{\gamma}} &= G_T \left[G_0^{-1} \boldsymbol{\gamma}_0 + \sum_{t=1}^T X_t' \Omega^{-1} \mathbf{y}_t \right] \end{aligned}$$

and $\Omega^{-1}|Y_T \sim \mathcal{W}_D(\rho_0 + T, R_T)$ where

$$R_T = \left[R_0^{-1} + \sum_{t=1}^T (\mathbf{y}_t - X_t \boldsymbol{\gamma}) (\mathbf{y}_t - X_t \boldsymbol{\gamma})' \right]^{-1}$$

3 MCMC

Using the full set of conditional posteriors, given in the preceding section, we can simulate from the joint posterior for this model using a Gibbs sampler:

1. Select a starting value $\Omega^{-1(0)}$ for the precision matrix.
2. Draw $\boldsymbol{\gamma}^{(1)} \sim \mathcal{N}(\bar{\boldsymbol{\gamma}}^{(1)}, G_T^{(1)})$ where

$$\begin{aligned} G_T^{(1)} &= \left[G_0^{-1} + \sum_{t=1}^T X_t' \Omega^{-1(0)} X_t \right]^{-1} \\ \bar{\boldsymbol{\gamma}}^{(1)} &= G_T^{(1)} \left[G_0^{-1} \boldsymbol{\gamma}_0 + \sum_{t=1}^T X_t' \Omega^{-1(0)} \mathbf{y}_t \right] \end{aligned}$$

3. Draw $\Omega^{-1(1)} \sim \mathcal{W}_D(\rho_T, R_T^{(1)})$ where

$$R_T^{(1)} = \left[R_0^{-1} + \sum_{t=1}^T (\mathbf{y}_t - X_t \boldsymbol{\gamma}^{(1)}) (\mathbf{y}_t - X_t \boldsymbol{\gamma}^{(1)})' \right]^{-1}$$

4. Repeat the preceding two steps a total of G times. In the g th iteration:

- (i) Draw $\boldsymbol{\gamma}^{(g)} \sim \mathcal{N}(\bar{\boldsymbol{\gamma}}^{(g)}, G_T^{(g)})$ where

$$\begin{aligned} G_T^{(g)} &= \left[G_0^{-1} + \sum_{t=1}^T X_t' \Omega^{-1(g-1)} X_t \right]^{-1} \\ \bar{\boldsymbol{\gamma}}^{(g)} &= G_T^{(g)} \left[G_0^{-1} \boldsymbol{\gamma}_0 + \sum_{t=1}^T X_t' \Omega^{-1(g-1)} \mathbf{y}_t \right] \end{aligned}$$

(ii) Draw $\Omega^{-1(g)} \sim \mathcal{W}_D \left(\rho_T, R_T^{(g)} \right)$ where

$$R_T^{(g)} = \left[R_0^{-1} + \sum_{t=1}^T (\mathbf{y}_t - X_t \boldsymbol{\gamma}^{(g)}) (\mathbf{y}_t - X_t \boldsymbol{\gamma}^{(g)})' \right]^{-1}$$

5. Discard the first B draws.

Note that in iteration g , $G_T^{(g)}$ and $\tilde{\boldsymbol{\gamma}}^{(g)}$ are calculated using $\Omega^{-1(g-1)}$ while $R_T^{(g)}$ is calculated using $\boldsymbol{\gamma}^{(0)}$. This is because we choose to initialize the sample with a starting value $\Omega^{-1(0)}$ for the precision matrix rather than for the vector of regression coefficients.

4 Numerical Details for the Gibbs Sampler

4.1 Efficient Calculation of R_T

In the second step of each iteration we compute $\left(R_0^{-1} + \sum_{t=1}^T \hat{\boldsymbol{\epsilon}}_t \hat{\boldsymbol{\epsilon}}_t' \right)^{-1}$ where $\hat{\boldsymbol{\epsilon}}_t = \mathbf{y}_t - X_t \boldsymbol{\gamma}$. Since R_0 is simply the prior scale matrix for Ω^{-1} and hence remains unchanged during the iterations, we can pre-compute it and store the result before starting the sampler. Since X_t is a sparse matrix, there is a much more efficient and compact way to compute the sum of outer products of residuals. Define:

$$\tilde{Y} = \begin{bmatrix} \mathbf{y}'_1 \\ \vdots \\ \mathbf{y}'_T \end{bmatrix}, \quad \tilde{X} = \begin{bmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_T \end{bmatrix}, \quad \Gamma = \begin{bmatrix} \gamma_1 & \cdots & \gamma_D \end{bmatrix}, \quad \hat{\boldsymbol{\epsilon}} = \begin{bmatrix} \hat{\boldsymbol{\epsilon}}'_1 \\ \vdots \\ \hat{\boldsymbol{\epsilon}}'_T \end{bmatrix}$$

so that $\hat{\boldsymbol{\epsilon}} = \tilde{Y} - \tilde{X}\Gamma$. Note that the vector of regression coefficients $\boldsymbol{\gamma}$ is the vec of the *matrix* of regression coefficients Γ . Thus, expressed in terms of dense matrix operations

$$R_T^{-1} = R_0^{-1} + \left(\tilde{Y} - \tilde{X}\Gamma \right)' \left(\tilde{Y} - \tilde{X}\Gamma \right)$$

The final step is to invert this sum (which is positive definite) to calculate R_T . Note that the Matrix Inversion Lemma (Sherman-Morrison-Woodbury Formula) does *not* simplify this calculation unless $D > T$.

4.2 Efficient Calculation of G_T

Because we parameterize our multivariate normal sampler in terms of the *precision* matrix rather than the covariance matrix, we work with the *inverse* of G_T , namely

$$G_T^{-1} = G_0^{-1} + \sum_{t=1}^T X_t' \Omega^{-1} X_t$$

Since it is simply the prior precision matrix for the vector γ of regression coefficients we can pre-compute G_0 (assuming that we elicit a prior in terms of the covariance matrix). Now, the sum over $X_t' \Omega^{-1} X_t$ can in fact be simplified using the properties of the Kronecker product.¹ Recall that $X_t = I_D \otimes \mathbf{x}_t$. Since $(A \otimes B)' = A' \otimes B'$,

$$X_t' \Omega^{-1} X_t = (I_D \otimes \mathbf{x}_t) \Omega^{-1} X_t$$

Since $\Omega^{-1} X_t = (\Omega^{-1} X_t) \otimes 1$, $\Omega^{-1} = \Omega^{-1} \otimes 1$, and $(A \otimes B)(C \otimes D) = AC \otimes BD$, provided that everything is conformable, we have

$$\begin{aligned} (I_D \otimes \mathbf{x}_t) \Omega^{-1} X_t &= (I_D \otimes \mathbf{x}_t) (\Omega^{-1} X_t \otimes 1) \\ &= \Omega^{-1} X_t \otimes \mathbf{x}_t = [\Omega^{-1} (I_D \otimes \mathbf{x}_t')] \otimes \mathbf{x}_t \\ &= [(\Omega^{-1} \otimes 1) (I_D \otimes \mathbf{x}_t')] \otimes \mathbf{x}_t \\ &= \Omega^{-1} \otimes \mathbf{x}_t' \otimes \mathbf{x}_t \end{aligned}$$

Finally, since $A \otimes (B + C) = A \otimes B + A \otimes C$,

$$\sum_{t=1}^T X_t' \Omega^{-1} X_t = \sum_{t=1}^T \Omega^{-1} \otimes \mathbf{x}_t' \otimes \mathbf{x}_t = \Omega^{-1} \otimes \left(\sum_{t=1}^T \mathbf{x}_t' \otimes \mathbf{x}_t \right)$$

This is an extremely useful simplification: because $\sum_{t=1}^T \mathbf{x}_t' \otimes \mathbf{x}_t$ involves neither Ω^{-1} nor γ , only the data, we can pre-compute this quantity. In fact, there is one final simplification that makes this quantity even simpler. By

¹See, e.g., Horn and Johnson (1994) Chapter 4.2.

writing out the definition of the Kronecker Product, we see that $\mathbf{x}_t' \otimes \mathbf{x}_t = \mathbf{x}_t \mathbf{x}_t'$ and hence

$$\sum_{t=1}^T X_t' \Omega^{-1} X_t = \Omega^{-1} \otimes \left(\sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t' \right) = \Omega^{-1} \otimes \tilde{X}' \tilde{X}$$

where $\tilde{X}' = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_T \end{bmatrix}$. Thus, we have $G_T^{-1} = G_0^{-1} + \Omega^{-1} \otimes \tilde{X}' \tilde{X}$.

4.3 Efficient Calculation of $\bar{\gamma}$

The vector $\bar{\gamma}$ is constructed from several pieces. The first is $G_0^{-1} \gamma_0$, the solution to the linear system $G_0 \mathbf{v} = \gamma_0$. Since this piece depends only on the prior, we can pre-compute it. The next piece is the sum $\sum_{t=1}^T X_t' \Omega^{-1} \mathbf{y}_t$. We noted above, X_t is sparse so there is a more efficient way to compute this quantity. Indeed, while this is far from obvious at first glance, it is possible to *factor* Ω^{-1} outside of the sum using some clever matrix operations, allowing us to drastically reduce the computational complexity of the sampler. To accomplish this simplification we combine the definition of X_t as $I_D \otimes \mathbf{x}_t'$ with two properties of the Kronecker Product, namely:

$$(A \otimes B)(C \otimes D) = AC \otimes BD$$

which holds provided that the respective matrices are conformable and

$$\text{vec}(AB) = (B' \otimes I_k) \text{vec}(A)$$

where A is $k \times \ell$ and B is $\ell \times m$. Applying the first property twice in succession followed by the second property, we find that

$$\begin{aligned} X_t' \Omega^{-1} \mathbf{y}_t &= (I_D \otimes \mathbf{x}_t')' \Omega^{-1} \mathbf{y}_t = (I_D \otimes \mathbf{x}_t) \Omega^{-1} \mathbf{y}_t \\ &= (I_D \otimes \mathbf{x}_t) (\Omega^{-1} \mathbf{y}_t \otimes 1) = I_D \Omega^{-1} \mathbf{y}_t \otimes \mathbf{x}_t 1 \\ &= \Omega^{-1} \mathbf{y}_t \otimes \mathbf{x}_t = [(\Omega^{-1} \mathbf{y}_t) \ 1] \otimes [I_{K+1} \mathbf{x}_t] \\ &= (\Omega^{-1} \mathbf{y}_t \otimes I_{K+1}) (1 \otimes \mathbf{x}_t) \\ &= ([\Omega^{-1} \mathbf{y}_t] \otimes I_{K+1}) \mathbf{x}_t = \left([\mathbf{y}_t' \Omega^{-1}]' \otimes I_{K+1} \right) \text{vec}(\mathbf{x}_t) \\ &= \text{vec}(\mathbf{x}_t \mathbf{y}_t' \Omega^{-1}) \end{aligned}$$

where we have used the fact that $\text{vec}(\mathbf{x}_t) = \mathbf{x}_t$. Finally, since we can interchange the vec summation operations,

$$\begin{aligned} \sum_{t=1}^T \text{vec}(\mathbf{x}_t \mathbf{y}_t' \Omega^{-1}) &= \text{vec} \left[\sum_{t=1}^T \mathbf{x}_t \mathbf{y}_t' \Omega^{-1} \right] = \text{vec} \left[\left(\sum_{t=1}^T \mathbf{x}_t \mathbf{y}_t' \right) \Omega^{-1} \right] \\ &= \text{vec}(\tilde{X}' \tilde{Y} \Omega^{-1}) \end{aligned}$$

where, as above,

$$\tilde{Y} = \begin{bmatrix} \mathbf{y}_1' \\ \vdots \\ \mathbf{y}_T' \end{bmatrix}, \quad \tilde{X} = \begin{bmatrix} \mathbf{x}_1' \\ \vdots \\ \mathbf{x}_T' \end{bmatrix}$$

Thus we see that

$$\bar{\gamma} = G_T \left[G_0^{-1} \gamma_0 + \text{vec}(\tilde{X}' \tilde{Y} \Omega^{-1}) \right]$$

Because it does not change between iterations, we can pre-compute the product $\tilde{X}' \tilde{Y}$. The only term that remains to be addressed is G_T . Because our normal sampler is parameterized in terms of the precision matrix rather than the covariance matrix we calculated G_T^{-1} rather than G_T above. Rather than inverting it in this step, which is a very bad idea given its size, we notice that our expression for $\bar{\gamma}$ takes the form $\mathbf{v} = A^{-1} \mathbf{b}$. Therefore,

$$\bar{\gamma} = \text{solve} \left[G_T^{-1}, G_0^{-1} \gamma_0 + \text{vec}(\tilde{X}' \tilde{Y} \Omega^{-1}) \right]$$

5 Calculating the Marginal likelihood

The marginal likelihood is available by the method of Chib (1995). From the Chib (1995) identity, we have

$$\log m(\{y_t\}) = \log \pi(\gamma^*) + \log \pi(\Omega^{-1*}) + \sum_{t=1}^n \log p(y_t | X_t \gamma^*, \Omega^*) - \log \pi(\gamma^*, \Omega^{-1*} | Y_n)$$

where the last term is calculated as

$$\pi(\Omega^{-1*} | Y_n) \times \pi(\gamma^* | Y_n, \Omega^{-1*})$$

in which the first term is estimated by averaging the full-conditional Wishart density over the draws $\{\gamma^{(g)}\}_{g=1}^G$ from the main MCMC run

$$\pi(\Omega^{-1*}|Y_n) = \frac{1}{G} \sum_{g=1}^G \mathcal{W}_d \left(\Omega^{-1*} | \rho_0 + n, \left(R_0^{-1} + \sum_{t=1}^n (y_t - X_t \gamma^{(g)}) (y_t - X_t \gamma^{(g)})' \right)^{-1} \right)$$

and the second term $\pi(\gamma^*|Y_n, \Omega^{-1*})$ is available directly as

$$\pi(\gamma^*|Y_n, \Omega^{-1*}) = \mathcal{N}_{d+p}(\gamma^* | \hat{\gamma}^*, G_n^*)$$

where

$$\begin{aligned} \hat{\gamma}^* &= G_n^* \left(G_0^{-1} \gamma_0 + \sum_{t=1}^n X_t' \Omega^{-1*} y_t \right) \\ G_n^* &= \left(G_0^{-1} + \sum_{t=1}^n X_t' \Omega^{-1*} X_t \right)^{-1} \end{aligned}$$

6 Prediction

Suppose we are interested in predicting the cross-section of returns y_{n+1} at time $(n+1)$. The Bayes prediction density of these returns, conditioned on the data Y_{n+1} and the factors f_{n+1} , is given by

$$p(y_{n+1}|Y_n, f_{n+1}) = \int_{\gamma, \Omega^{-1}} \mathcal{N}_d(y_{n+1} | X_{n+1} \gamma, \Omega) d\pi(\gamma, \Omega^{-1} | Y_n)$$

which is estimated by the ergodic Monte Carlo average

$$p(y_{n+1}|Y_n, f_{n+1}) = \frac{1}{G} \sum_{g=1}^G \mathcal{N}_d(y_{n+1} | X_{n+1} \gamma^{(g)}, \Omega^{(g)})$$

with the MCMC draws $\{\gamma^{(g)}, \Omega^{(g)}\}$ from the posterior distribution.

7 Student-t errors

Suppose now that the errors are distributed as multivariate-t

$$\varepsilon_t \sim t_{d,\nu}(0, \Omega)$$

so that

$$\begin{aligned} E(\varepsilon_t) &= 0, \quad \nu > 1 \\ Var(\varepsilon_t) &= \frac{\nu}{\nu - 2} \Omega, \quad \nu > 2 \end{aligned}$$

The analysis of this model utilizes the hierarchical representation

$$\begin{aligned} \varepsilon_t | \lambda_t &\sim N(0, \lambda_t^{-1} \Omega) \\ \lambda_t &\sim G\left(\frac{\nu}{2}, \frac{\nu}{2}\right) \end{aligned}$$

which means that conditioned on $(\nu, \{\lambda_t\})$, the results presented for the Gaussian model can be applied with minor modifications. The MCMC sampling is completed with the sampling of $(\nu, \{\lambda_t\})$.

Following Albert and Chib (1993), let us assume that the support of ν is the set of values $\{\nu_j\}_{j=1}^J$, for example, $\{4, 6, 8, 10, 12, 14, 16\}$ and that a priori

$$\Pr(\nu = \nu_j) = q_j$$

Then, simple calculations show that

$$\gamma | Y_n, \Omega^{-1}, \nu, \{\lambda_t\} \sim \mathcal{N}_{d+p}(\hat{\gamma}_\lambda, G_{n,\lambda})$$

where

$$\begin{aligned} \hat{\gamma}_\lambda &= G_{n,\lambda} \left(G_0^{-1} \gamma_0 + \sum_{t=1}^n \lambda_t X_t' \Omega^{-1} y_t \right) \\ G_{n,\lambda} &= \left(G_0^{-1} + \sum_{t=1}^n \lambda_t X_t' \Omega^{-1} X_t \right)^{-1} \end{aligned}$$

and

$$\Omega^{-1}|Y_n, \gamma, \nu, \{\lambda_t\} \sim \mathcal{W}_d \left(\rho_0 + n, \left(R_0^{-1} + \sum_{t=1}^n \lambda_t (y_t - X_t \gamma) (y_t - X_t \gamma)' \right)^{-1} \right)$$

Moreover,

$$\Pr(\nu = \nu_j | Y_n, \gamma, \Omega^{-1}) \propto q_j \prod_{t=1}^n t_{d, \nu_j}(y_t | X_t \gamma, \Omega)$$

and

$$\lambda_t | Y_n, \gamma, \nu \sim G \left(\frac{\nu + d}{2}, \frac{\nu + (y_t - X_t \gamma)' (y_t - X_t \gamma)}{2} \right)$$

One sweep of the MCMC sampling is completed by sampling these four distributions in this order.

7.1 Marginal likelihood

The Chib (1995) method can again be applied to find the log marginal likelihood as

$$\log \Pr(\nu^*) + \log \pi(\gamma^*) + \log \pi(\Omega^{-1*}) + \sum_{t=1}^n \log t_{d, \nu^*}(y_t | X_t \gamma^*, \Omega^*) - \log \pi(\nu^*, \gamma^*, \Omega^{-1*} | Y_n)$$

where ν^* is the posterior mode (which is easily computed from the sampled values), the last term is calculated as

$$\Pr(\nu^* | Y_n) \times \pi(\Omega^{-1*} | Y_n, \nu^*) \times \pi(\gamma^* | Y_n, \Omega^{-1*}, \nu^*)$$

in which the first term is obtained from the posterior frequency distribution of ν , the second term is obtained from a reduced run in which ν is fixed at ν^* and the remaining three distributions are sampled and the draws

$$\left\{ \gamma^{(g)}, \lambda_t^{(g)} \right\}_{g=1}^G$$

from this reduced MCMC run are used to calculate $\pi(\Omega^{-1*} | Y_n, \nu^*)$ as

$$\frac{1}{G} \sum_{g=1}^G \mathcal{W}_d \left(\Omega^{-1*} | \rho_0 + n, \left(R_0^{-1} + \sum_{t=1}^n \lambda_t^{(g)} (y_t - X_t \gamma^{(g)}) (y_t - X_t \gamma^{(g)})' \right)^{-1} \right)$$

and the final term $\pi(\gamma^*|Y_n, \Omega^{-1*})$ is obtained from a second reduced run in which ν is fixed at ν^* and Ω^{-1} is fixed at Ω^{-1*} and the draws

$$\{\lambda_t^{(g)}\}$$

from this reduced run are used to give

$$\pi(\gamma^*|Y_n, \Omega^{-1*}, \nu^*) = \frac{1}{G} \sum_{g=1}^G \mathcal{N}_{d+p}(\gamma^*|\hat{\gamma}_{\lambda^{(g)}}^*, G_{n, \lambda^{(g)}}^*)$$

where

$$\begin{aligned} \hat{\gamma}_{\lambda^{(g)}}^* &= G_{n, \lambda^{(g)}}^* \left(G_0^{-1} \gamma_0 + \sum_{t=1}^n \lambda_t^{(g)} X_t' \Omega^{-1*} y_t \right) \\ G_{n, \lambda^{(g)}}^* &= \left(G_0^{-1} + \sum_{t=1}^n \lambda_t^{(g)} X_t' \Omega^{-1*} X_t \right)^{-1} \end{aligned}$$