

Hierarchical Models for Heterogeneous Preferences

(KPT Ch. 7)

AAEC 6564

Instructor: KLAUS MOELTNER

Matlab scripts: `mod6_data, mod6_HNRM_v1, mod6_HNRM_v2, mod6_HNRM_outage, mod6_HNRM_outage_chib, mod6_HNRM_outage_nocovs, mod6_HNRM_outage_chib, mod6_HNRM_predict`

Matlab functions: `gs_HNRM_v1, gs_HNRM_v2, gs_HNRM_chib, gs_HNRM_nocovs, gs_HNRM_nocovs_chib`

“Hierarchical modeling” has become quite popular in both Classical and Bayesian analysis. These models have one or more additional distributional layers between the base structure and the priors for some of the parameters. Hierarchical models (HMs) can be estimated via a straightforward extension of the models we have already discussed. In this Session we will consider a Hierarchical Regression Model (HRM) with some, but not all, parameters taking such a “hyper-distribution”. This is often called a regression model with “mixed coefficients”, or a “*mixed effects*”-model.

Hierarchical normal regression model

HMs require multiple observations per individual (or other observational unit of interest). Thus, panel data models with multiple observations over time per person are a special case of HMs. Any other grouping leading to multiple individual observations qualifies as well, such as visits to several hiking trails per individual, demand for several commodities per household, etc. For simplicity let's assume we have the same number of observations (say J) per person, for a total of N individuals, although such a “balanced” setup is not a requirement. Thus the total sample size is $n = NJ$. The model at the individual level is then given as

$$\mathbf{y}_i = \mathbf{X}_{i1} \boldsymbol{\beta}_f + \mathbf{X}_{i1} \boldsymbol{\beta}_{ri} + \boldsymbol{\varepsilon}_i \quad \boldsymbol{\varepsilon}_i \sim n(0, \sigma^2 \mathbf{I}_J) \quad (1)$$

$$\boldsymbol{\beta}_{ri} \sim n(\boldsymbol{\beta}_r, \boldsymbol{\Sigma})$$

where $\boldsymbol{\beta}_f$ is a vector of “fixed coefficients” shared by all individuals (as would be the case in the generic linear regression model), $\boldsymbol{\beta}_{ri}$ is an individual-specific coefficient vector (also called “random coefficients”), and the regression error $\boldsymbol{\varepsilon}_i$ has the usual properties of the basic regression model. The new element in this model is the *hyper-distribution* for the random vector $\boldsymbol{\beta}_{ri}$. Specifically, each individual “draws” her $\boldsymbol{\beta}_{ri}$ from another normal density with “grand mean” $\boldsymbol{\beta}_r$ and VCOV $\boldsymbol{\Sigma}$. The random coefficients are often interpreted as allowing for “heterogeneity in tastes” over individuals. A classical example would be the effect of characteristics of consumer goods (such as cars) on potential buyers. You might care a lot about a great sound system, but I may not care as much. Such, if your regression measures the maximum price a person would pay for a car, and “sound system quality” is one of the

regressors, it would make sense to ex ante allow for different marginal effects of this regressor over individuals.

Theoretically, one could aim at estimating all N vectors β_{ri} , but that would require N additional degrees of freedom, which may be infeasible in many applications. A more compact approach is to draw β_{ri} from a shared hyper-distribution as shown in (1). That is the essence of hierarchical modeling.

It is often convenient to re-write (1) in terms of deviations of random coefficients from the grand mean, i.e.

$$\begin{aligned} \mathbf{y}_i &= \mathbf{X}_{fi}\beta_f + \mathbf{X}_{ri}\beta_r + \mathbf{X}_{ri}\alpha_{ri} + \varepsilon_i & \text{with} & \quad \alpha_{ri} = \beta_{ri} - \beta_r \\ \alpha_{ri} &\sim n(\mathbf{0}, \Sigma) \end{aligned} \quad (2)$$

This makes it easier to recognize the modeling distribution of \mathbf{y}_i , *unconditional* on β_{ri} as

$$\mathbf{y}_i \sim n(\boldsymbol{\mu}_i, \mathbf{V}_i), \quad \text{with} \quad \boldsymbol{\mu}_i = \mathbf{X}_{fi}\beta_f + \mathbf{X}_{ri}\beta_r, \quad \mathbf{V}_i = \mathbf{X}_{ri}\Sigma\mathbf{X}'_{ri} + \sigma^2\mathbf{I}_J, \quad (3)$$

which will be convenient later in our estimation. Note that we know the exact form of this unconditional density ONLY because we specified normal distributions for both our main regression model and the hierarchical elements. For most other combination of base density plus hierarchical density (or densities), we would have to work with explicit integration.

The full model for the entire sample of $N \times J$ observations can be written as

$$\begin{aligned} \mathbf{y} &= \mathbf{X}_f\beta_f + \mathbf{X}_r\tilde{\beta}_r + \boldsymbol{\varepsilon} & \text{where} \\ \mathbf{X}_f &= \begin{bmatrix} \mathbf{X}_{f1} \\ \mathbf{X}_{f2} \\ \vdots \\ \mathbf{X}_{fN} \end{bmatrix} & \mathbf{X}_r &= \begin{bmatrix} \mathbf{X}_{r1} & 0 & \cdots & 0 \\ 0 & \mathbf{X}_{r2} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & \mathbf{X}_{rN} \end{bmatrix} & \tilde{\beta}_r &= \begin{bmatrix} \beta_{r1} \\ \beta_{r2} \\ \vdots \\ \beta_{rN} \end{bmatrix} & \boldsymbol{\varepsilon} &= \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{bmatrix} \sim n(0, \sigma^2\mathbf{I}_{N \times J}) \end{aligned} \quad (4)$$

$(N \times J) \times k_f$ $(N \times J) \times (N \times k_r)$ $(N \times k_r) \times 1$

Our goal is to derive the posterior distribution for β_f , β_r , Σ , and σ^2 . We are NOT generally interested in realizations of the individual coefficient vectors β_{ri} , although they will help us draw the main parameters of interest in our Gibbs Sampler. I recommend always collecting these draws & saving them for later, especially if you're planning to use the Chib method to evaluate the marginal likelihood.

The likelihood function for this model, unconditional on the random effects, can be written as

$$\begin{aligned}
p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2) = & \\
\prod_{i=1}^N \left\{ \int_{\boldsymbol{\beta}_{ri}} (2\pi)^{-J/2} \sigma^{-J/2} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))\right) \right\} & \quad (5) \\
\text{with } p(\boldsymbol{\beta}_{ri}) = (2\pi)^{-k_r/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2} (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)\right). &
\end{aligned}$$

In HMs, the *LHF* will usually include an integral over the random elements. This form would make it difficult to break the posterior into “well-behaved” conditionals for all parameters of interest for most models. In our case, we have a base normal density plus a hierarchical normal, which allows us to write the *LHF* without an explicit integration (or, in other words, solve the integral in (5) analytically). Using our result from (3), we can write

$$\begin{aligned}
p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2) = \prod_{i=1}^N (2\pi)^{-J/2} |\mathbf{V}_i|^{-1/2} \exp\left(-\frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu}_i)' \mathbf{V}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i)\right) & \\
(2\pi)^{-NJ/2} |\mathbf{V}_i|^{-N/2} \exp\left(-\frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - \boldsymbol{\mu}_i)' \mathbf{V}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i)\right) & \quad (6) \\
\text{with } \mathbf{V}_i = \mathbf{X}_{ri} \boldsymbol{\Sigma} \mathbf{X}_{ri}' + \sigma^2 \mathbf{I}_J, \quad \boldsymbol{\mu}_i = \mathbf{X}_{fi} \boldsymbol{\beta}_f + \mathbf{X}_{ri} \boldsymbol{\beta}_r. &
\end{aligned}$$

The priors

We use our usual combination of normal, inverse gamma, and *IW* for prior densities. Specifically:

$$\begin{aligned}
p(\boldsymbol{\beta}_f) &= n(\boldsymbol{\mu}_{0f}, \mathbf{V}_{0f}) \\
p(\boldsymbol{\beta}_r) &= n(\boldsymbol{\mu}_{0r}, \mathbf{V}_{0r}) \\
p(\boldsymbol{\Sigma}) &= IW(\nu_0, S_0) \\
p(\sigma^2) &= ig(\eta_0, \tau_0)
\end{aligned} \quad (7)$$

Note that I have switched the notation for the shape parameter in the inverse gamma from ν_0 to η_0 to distinguish it from the degrees-of-freedom parameter in the Inverse Wishart. You have already encountered the explicit forms of these densities in previous sections.

The posterior and Data Augmentation

Generically, we can write the posterior kernel as

$$p(\boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto p(\boldsymbol{\beta}_f) p(\boldsymbol{\beta}_r) p(\boldsymbol{\Sigma}) p(\sigma^2) p(\mathbf{y} | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2, \mathbf{X}) \quad (8)$$

For some of our parameters, this form would make it very awkward to derive “well-behaved” conditionals. Specifically, for our application it would be convenient to work with the individual random vectors $\boldsymbol{\beta}_{ri}$. To accomplish this without running into the integral in (5) we’ll employ “*data augmentation*” as discussed in module 5.

For this application we augment our joint posterior with the n random coefficient vectors $\boldsymbol{\beta}_{ri}$, i.e.

$$\begin{aligned}
 & p(\boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ri} (i=1 \cdots N), \boldsymbol{\Sigma}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto \\
 & p(\boldsymbol{\beta}_f) p(\boldsymbol{\beta}_r, \boldsymbol{\beta}_{ri} (i=1 \cdots N), \boldsymbol{\Sigma}) p(\sigma^2) p(\mathbf{y} | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ri} (i=1 \cdots N), \boldsymbol{\Sigma}, \sigma^2, \mathbf{X}) = \\
 & p(\boldsymbol{\beta}_f) \left(\prod_{i=1}^n p(\boldsymbol{\beta}_{ri} | \boldsymbol{\beta}_r, \boldsymbol{\Sigma}) \right) p(\boldsymbol{\beta}_r) p(\boldsymbol{\Sigma}) p(\sigma^2) p(\mathbf{y} | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ri} (i=1 \cdots N), \boldsymbol{\Sigma}, \sigma^2, \mathbf{X})
 \end{aligned} \tag{9}$$

The goal is then to integrate out the effect of the augmented data, i.e.

$$p(\boldsymbol{\theta} | \mathbf{y}) = \int_{\boldsymbol{\beta}_{ri}, i=1 \cdots N} p(\boldsymbol{\theta} | \mathbf{y}, \boldsymbol{\beta}_{ri}) p(\boldsymbol{\beta}_{ri} | \mathbf{y}) d\boldsymbol{\beta}_{ri} \tag{10}$$

where $\boldsymbol{\theta}$ includes all actual parameters. This can be accomplished in straightforward manner by making draws of $\boldsymbol{\beta}_{ri}$, conditional on actual data and all other parameters, a part of the Gibbs Sampler.

It is customary to think of the last part on the right hand side of (8) as the *conditional LHF*, which can be written as

$$\begin{aligned}
 & p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}_f, \boldsymbol{\beta}_{ri} (i=1 \cdots n), \sigma^2) = \\
 & (2\pi)^{-NJ/2} \sigma^{-NJ/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - (\mathbf{X}_{ri}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{ri}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))\right)
 \end{aligned} \tag{11}$$

We will drop the $(i=1 \cdots n)$ part henceforth for notational convenience. It will be clear from the context if we condition on all n vectors of $\boldsymbol{\beta}_{ri}$ or a specific $\boldsymbol{\beta}_{ri}$ associated with a single individual.

Also, with data augmentation, we work with an additional “prior”, i.e. the hierarchical prior

$\left(\prod_{i=1}^N p(\boldsymbol{\beta}_{ri} | \boldsymbol{\beta}_r, \boldsymbol{\Sigma}) \right)$. This hierarchical prior is normally listed along with the other priors. The priors for the moments of the hierarchical density (i.e. the priors for $\boldsymbol{\beta}_r$ and $\boldsymbol{\Sigma}$ in our case) are often called “hyper-priors”.

The full posterior kernel for the *augmented model* can then be explicitly written as:

$$\begin{aligned}
& p(\boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ri} (i=1 \dots N), \boldsymbol{\Sigma}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto \\
& \exp\left(-\frac{1}{2}(\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})' \mathbf{V}_{f0}^{-1} (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})\right) * \exp\left(-\frac{1}{2}(\boldsymbol{\beta}_r - \boldsymbol{\mu}_{r0})' \mathbf{V}_{r0}^{-1} (\boldsymbol{\beta}_r - \boldsymbol{\mu}_{r0})\right) * \\
& (\sigma^2)^{-(v_0+1)} \exp\left(-\frac{\tau_0}{\sigma^2}\right) * |\boldsymbol{\Sigma}|^{-(v_0+M+1)/2} \exp\left(-\frac{1}{2}tr(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) * \\
& |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2} \sum_{i=1}^N (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)\right) * \\
& \sigma^{-NJ/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))\right)
\end{aligned} \tag{12}$$

The Gibbs Sampler

We will first examine the GS for the fully augmented model. We start by drawing the vector of fixed coefficients $\boldsymbol{\beta}_f$. The conditional posterior kernel is given as:

$$\begin{aligned}
& p(\boldsymbol{\beta}_f | \boldsymbol{\beta}_{ri}, \sigma^2, \mathbf{y}, \mathbf{X}) \propto \\
& \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri})) + (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})' \mathbf{V}_{f0}^{-1} (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})\right)
\end{aligned} \tag{13}$$

Letting $\tilde{\mathbf{y}}_i = \mathbf{y}_i - \mathbf{X}'_{ri}\boldsymbol{\beta}_{ri}$ we can write this as

$$\begin{aligned}
& p(\boldsymbol{\beta}_f | \boldsymbol{\beta}_{ri}, \sigma^2, \mathbf{y}, \mathbf{X}) \propto \\
& \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (\tilde{\mathbf{y}}_i - \mathbf{X}_{fi}\boldsymbol{\beta}_f)' (\tilde{\mathbf{y}}_i - \mathbf{X}_{fi}\boldsymbol{\beta}_f) + (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})' \mathbf{V}_{f0}^{-1} (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})\right)
\end{aligned} \tag{14}$$

This brings us back to the exact same form of this conditional posterior kernel for the basic regression model with independent priors. We can use the results presented there to derive the conditional posterior as:

$$\boldsymbol{\beta}_f | \boldsymbol{\beta}_{ri}, \sigma^2, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_1, \mathbf{V}_1) \quad \text{with} \quad \mathbf{V}_1 = \left(\mathbf{V}_{f0}^{-1} + \frac{1}{\sigma^2} \mathbf{X}'_f \mathbf{X}_f\right)^{-1} \quad \text{and} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left(\mathbf{V}_{f0}^{-1} \boldsymbol{\mu}_{f0} + \frac{1}{\sigma^2} \mathbf{X}'_f \tilde{\mathbf{y}}\right) \tag{15}$$

Next, it is convenient to draw the individual random vectors $\boldsymbol{\beta}_{ir}$. This requires working with our data-augmented framework. We aim to draw $\boldsymbol{\beta}_{ir}$ from $p(\boldsymbol{\beta}_{ri} | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}_i, \mathbf{X}_i)$. Note that only data corresponding to individual “ i ” are relevant. The likelihood component for individual i , conditional on $\boldsymbol{\beta}_{ri}$, is given by

$$p(\mathbf{y}_i | \mathbf{X}_i, \boldsymbol{\beta}_f, \boldsymbol{\beta}_{ri}, \sigma^2) = (2\pi)^{-J/2} \sigma^{-J/2} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y}_i - (\mathbf{X}_{ri}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{ri}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))\right) \quad (16)$$

which is basically a “mini” regression model with J observations. The “prior” for $\boldsymbol{\beta}_{ri}$ is simply its hierarchical distribution, i.e. the last equation in (5). In other words, we need the i^{th} component of the conditional *LHF* in (11) plus the hierarchical prior of $\boldsymbol{\beta}_{ri}$ for this step.

Using a similar “trick” as before and defining $\tilde{\mathbf{y}}_i = \mathbf{y}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_f$, we can write the conditional kernel for drawing $\boldsymbol{\beta}_{ri}$ as

$$p(\boldsymbol{\beta}_{ri} | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}_i, \mathbf{X}_i) \propto \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2}(\tilde{\mathbf{y}}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_{ri})' (\tilde{\mathbf{y}}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}) + (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)\right)\right) \quad (17)$$

This looks exactly like the kernel for the basic regression model, except with prior mean $\boldsymbol{\mu}_0$ replaced by hierarchical mean $\boldsymbol{\beta}_r$ and prior variance \mathbf{V}_0 replaced by the hierarchical variance $\boldsymbol{\Sigma}$. We can thus immediately write the conditional density as

$$\begin{aligned} p(\boldsymbol{\beta}_{ri} | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}_i, \mathbf{X}_i) &\sim n(\boldsymbol{\mu}_{ri1}, \mathbf{V}_{ri1}) \text{ with} \\ \mathbf{V}_{ri1} &= \left(\boldsymbol{\Sigma}^{-1} + \frac{1}{\sigma^2} \mathbf{X}'_{ri} \mathbf{X}_{ri}\right)^{-1} \\ \boldsymbol{\mu}_{ri1} &= \mathbf{V}_{ri1} \left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}_r + \frac{1}{\sigma^2} \mathbf{X}'_{ri} \tilde{\mathbf{y}}_i\right) = \mathbf{V}_{ri1} \left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}_r + \frac{1}{\sigma^2} \mathbf{X}'_{ri} (\mathbf{y}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_f)\right) \end{aligned} \quad (18)$$

Keep in mind that you will have to repeat this N times, for each of the N $\boldsymbol{\beta}_{ri}$'s.

Conditional on knowing all $\boldsymbol{\beta}_{ri}$'s, drawing the vector of random coefficient means, $\boldsymbol{\beta}_r$ becomes especially simple in data-augmented setting. Upon examination of the full posterior kernel we realize that $\boldsymbol{\beta}_r$ only appears in its own prior and the hyper-distribution for $\boldsymbol{\beta}_{ri}$. Thus it's conditional posterior kernel is independent of the likelihood and data. We have:

$$\begin{aligned} p(\boldsymbol{\beta}_r | \boldsymbol{\beta}_{ri}, \boldsymbol{\Sigma}) &\propto \\ \exp\left(-\frac{1}{2}\left((\boldsymbol{\beta}_r - \boldsymbol{\mu}_{r0})' \mathbf{V}_{r0}^{-1} (\boldsymbol{\beta}_r - \boldsymbol{\mu}_{r0}) + \sum_{i=1}^N (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)\right)\right) \end{aligned} \quad (19)$$

Using transformations analogous to the SUR model we can derive the conditional posterior for $\boldsymbol{\beta}_r$ as

$$\begin{aligned} p(\boldsymbol{\beta}_r | \boldsymbol{\beta}_{ri}, \boldsymbol{\Sigma}) &= mvn(\boldsymbol{\mu}_1, \mathbf{V}_1) \quad \text{where} \\ \mathbf{V}_1 &= \left(\mathbf{V}_{r0}^{-1} + n \cdot \boldsymbol{\Sigma}^{-1}\right)^{-1} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left(\mathbf{V}_{r0}^{-1} \boldsymbol{\mu}_{r0} + \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \boldsymbol{\beta}_{ri}\right) \end{aligned} \quad (20)$$

Now on to the variance terms. As mentioned before, it will be convenient to draw σ^2 conditional on the $\boldsymbol{\beta}_{ri}$ vectors. The relevant components of the full posterior kernel yield:

$$\begin{aligned}
p(\sigma^2 | \boldsymbol{\beta}_f, \boldsymbol{\beta}_{ri}, \mathbf{y}, \mathbf{X}) &\propto \\
(\sigma^2)^{-(v_0+1)} \exp\left(-\frac{\tau_0}{\sigma^2}\right) \sigma^{-NJ/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))\right) &= \quad (21) \\
(\sigma^2)^{\frac{NJ-2v_0-2}{2}} \exp\left(-\frac{1}{2\sigma^2} \left(2\tau_0 + \sum_{i=1}^N (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))\right)\right) &
\end{aligned}$$

This lets us immediately use the results from the basic regression model with independent priors, i.e.

$$\begin{aligned}
\sigma^2 | \boldsymbol{\beta}_f, \boldsymbol{\beta}_{ri}, \mathbf{y}, \mathbf{X} &\sim ig(\eta_1, \tau_1) \quad \text{with} \\
\eta_1 = \frac{2\eta_0 + n}{2} \quad \text{and} \quad \tau_1 &= \frac{2\eta_0 + \sum_{i=1}^n (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))' (\mathbf{y}_i - (\mathbf{X}_{fi}\boldsymbol{\beta}_f + \mathbf{X}_{ri}\boldsymbol{\beta}_{ri}))}{2} \quad (22)
\end{aligned}$$

The final step in our GS is the draw of the hierarchical variance $\boldsymbol{\Sigma}$. We use again our data-augmented framework and our analogy to the SUR model. Using similar arguments as for the draw of the hierarchical mean $\boldsymbol{\beta}_r$, we get

$$\begin{aligned}
p(\boldsymbol{\Sigma} | \boldsymbol{\beta}_{ri}, \boldsymbol{\beta}_r) &= IW(v_1, S_1) \quad \text{where} \\
v_1 = v_0 + n \quad S_1 &= S_0 + \sum_{i=1}^n (\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)(\boldsymbol{\beta}_{ri} - \boldsymbol{\beta}_r)' \quad (23)
\end{aligned}$$

An implementation of this framework using artificial data, generated in script `mod6_data`, is given in Matlab script `mod6_HNRM_v1`. The GS is provided in function `gs_HNRM_v1`.

An improved Gibbs Sampler

While the Sampler described above is easy to implement, it generally leads to considerable efficiency losses (i.e. high autocorrelation in the GS, slow convergence – see Chib & Carlin, 1999, for details), especially for draws of $\boldsymbol{\beta}_f$ and $\boldsymbol{\beta}_r$ (this is evident from the results of `HNRM_v1`). Thus, for these two parameters we will stick with our non-augmented framework implicitly given in (8). This also implies that we will work with the unconditional *LHF* in (6).

Another way to look at this is from a paradigm of "efficient blocking", briefly discussed in Module 2. For the fully augmented model, the blocking is as follows:

1. Draw $\boldsymbol{\beta}_f$ from $p(\boldsymbol{\beta}_f | \boldsymbol{\beta}_{ri}, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}, \mathbf{X})$

2. Draw β_{ri} from $p(\beta_{ri} | \beta_f, \beta_r, \Sigma, \sigma^2, y_i, X_i) \quad \forall i$
3. Draw β_r from $p(\beta_r | \{\beta_{ri}\}_{i=1}^N, \Sigma)$
4. Draw σ^2 from $p(\sigma^2 | \beta_f, \beta_{ri}, y, X)$
5. Draw Σ from $p(\Sigma | \beta_{ri}, \beta_r)$

Thus, even though β_{ri} is the only posterior element that is drawn conditional on *all other* parameters, we have five separate blocks. This is because for the remaining blocks, all parameters involved are mutually conditional with respect to at least one other conditioning element. To be specific, β_r is conditioned on Σ , and vice versa, and β_f is conditioned on σ^2 and vice versa.

We have learned that, in general, a more efficient GS results if we use as few blocks as possible. Here, we can exploit the fact that we know the form of the partially conditional densities $p(\beta_f | \beta_r, \Sigma, \sigma^2, y, X)$, and $p(\beta_r | \beta_f, \sigma^2, \Sigma, y, X)$. In other words, we can draw β_f , β_r and β_{ri} in two blocks as follows (see Chib and Carlin, 1999):

1. Draw β_f and β_{ri} from $p(\beta_f, \{\beta_{ri}\}_{i=1}^N | \beta_r, \sigma^2, \Sigma, y, X)$ via:
 - a. Draw β_{ri} from $p(\{\beta_{ri}\}_{i=1}^N | \beta_f, \beta_r, \sigma^2, \Sigma, y_i, X_i)$
 - b. Draw β_f from $p(\beta_f | \sigma^2, \beta_r, \Sigma, y, X)$
2. Draw β_r from $p(\beta_r | \beta_f, \sigma^2, \Sigma, y, X)$

The variance terms are drawn as in the fully augmented model. This generally yields a more efficient posterior sampler.

The full posterior kernel for the *non*-augmented model can then be explicitly written as:

$$\begin{aligned}
& p(\beta_f, \beta_r, \Sigma, \sigma^2 | y, X) \propto \\
& \exp\left(-\frac{1}{2}(\beta_f - \mu_{f0})' \mathbf{V}_{f0}^{-1} (\beta_f - \mu_{f0})\right) * \exp\left(-\frac{1}{2}(\beta_r - \mu_{r0})' \mathbf{V}_{r0}^{-1} (\beta_r - \mu_{r0})\right) * \\
& (\sigma^2)^{-(v_0+1)} \exp\left(-\frac{\tau_0}{\sigma^2}\right) * |\Sigma|^{-(v_0+M+1)/2} \exp\left(-\frac{1}{2}tr(\mathbf{S}_0 \cdot \Sigma^{-1})\right) * \\
& \left| \left(\mathbf{X}_{ri} \Sigma \mathbf{X}_{ri}' + \sigma^2 \mathbf{I}_J \right) \right|^{-N/2} \exp\left(-\frac{1}{2} \sum_{i=1}^N (y_i - \mu_i)' \left(\mathbf{X}_{ri} \Sigma \mathbf{X}_{ri}' + \sigma^2 \mathbf{I}_J \right)^{-1} (y_i - \mu_i) \right)
\end{aligned} \tag{24}$$

In contrast to the fully augmented model in (12) the hierarchical prior for β_{ri} and the conditional likelihood, i.e. the last two lines in (12) have now been replaced by the unconditional likelihood. Note again that this option of switching between the data- augmented model and the original model for a sub-

set of parameters only exists because we happen to know the analytical form of the non-augmented likelihood for the normal-normal case. In other hierarchical settings that do not build on a normal likelihood function and / or do not use a normal hyper-distribution for random coefficients we would have to stick with the fully augmented model.

Letting again $\tilde{\mathbf{y}}_i = \mathbf{y}_i - \mathbf{X}'_{ri}\boldsymbol{\beta}_r$ we can write the conditional posterior kernel for $\boldsymbol{\beta}_f$ as

$$p(\boldsymbol{\beta}_f | \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2}\left(\sum_{i=1}^n (\tilde{\mathbf{y}}_i - \mathbf{X}_{fi}\boldsymbol{\beta}_f)' \mathbf{V}_i^{-1} (\tilde{\mathbf{y}}_i - \mathbf{X}_{fi}\boldsymbol{\beta}_f) + (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})' \mathbf{V}_{f0}^{-1} (\boldsymbol{\beta}_f - \boldsymbol{\mu}_{f0})\right)\right) \quad (25)$$

with $\mathbf{V}_i = \mathbf{X}_{ri}\boldsymbol{\Sigma}\mathbf{X}'_{ri} + \sigma^2\mathbf{I}_J$

Using similar transformations as for the SUR model, we can derive the conditional kernel of $\boldsymbol{\beta}_f$ as

$$\boldsymbol{\beta}_f | \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_{f1}, \mathbf{V}_{f1}) \quad \text{with}$$

$$\mathbf{V}_{f1} = \left(\mathbf{V}_{f0}^{-1} + \sum_{i=1}^n \mathbf{X}'_{fi}\mathbf{V}_i^{-1}\mathbf{X}_{fi}\right)^{-1} = \left(\mathbf{V}_{f0}^{-1} + \sum_{i=1}^n \mathbf{X}'_{fi}(\mathbf{X}_{ri}\boldsymbol{\Sigma}\mathbf{X}'_{ri} + \sigma^2\mathbf{I}_J)^{-1}\mathbf{X}_{fi}\right)^{-1} \quad (26)$$

$$\boldsymbol{\mu}_{f1} = \mathbf{V}_{f1} \left(\mathbf{V}_{f0}^{-1}\boldsymbol{\mu}_{f0} + \sum_{i=1}^n \mathbf{X}'_{fi}\mathbf{V}_i^{-1}\tilde{\mathbf{y}}_i\right) = \left(\mathbf{V}_{f0}^{-1}\boldsymbol{\mu}_{f0} + \sum_{i=1}^n \mathbf{X}'_{fi}(\mathbf{X}_{ri}\boldsymbol{\Sigma}\mathbf{X}'_{ri} + \sigma^2\mathbf{I}_J)^{-1}(\mathbf{y}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_r)\right)$$

Next, we draw the hierarchical mean vector $\boldsymbol{\beta}_r$. As for $\boldsymbol{\beta}_f$, we can gain efficiency by NOT conditioning on $\boldsymbol{\beta}_{ri}$. Letting $\tilde{\mathbf{y}}_i = \mathbf{y}_i - \mathbf{X}'_{fi}\boldsymbol{\beta}_f$ we can write the conditional posterior kernel as

$$p(\boldsymbol{\beta}_r | \boldsymbol{\beta}_f, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2}\left(\sum_{i=1}^n (\tilde{\mathbf{y}}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_r)' \mathbf{V}_i^{-1} (\tilde{\mathbf{y}}_i - \mathbf{X}_{ri}\boldsymbol{\beta}_r) + (\boldsymbol{\beta}_r - \boldsymbol{\mu}_{r0})' \mathbf{V}_{r0}^{-1} (\boldsymbol{\beta}_r - \boldsymbol{\mu}_{r0})\right)\right) \quad (27)$$

Thus, in perfect analogy to $\boldsymbol{\beta}_f$, we can write

$$\boldsymbol{\beta}_r | \boldsymbol{\beta}_f, \boldsymbol{\Sigma}, \sigma^2, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_{r1}, \mathbf{V}_{r1}) \quad \text{with}$$

$$\mathbf{V}_{r1} = \left(\mathbf{V}_{r0}^{-1} + \sum_{i=1}^n \mathbf{X}'_{ri}\mathbf{V}_i^{-1}\mathbf{X}_{ri}\right)^{-1} = \left(\mathbf{V}_{r0}^{-1} + \sum_{i=1}^n \mathbf{X}'_{ri}(\mathbf{X}_{ri}\boldsymbol{\Sigma}\mathbf{X}'_{ri} + \sigma^2\mathbf{I}_J)^{-1}\mathbf{X}_{ri}\right)^{-1} \quad (28)$$

$$\boldsymbol{\mu}_{r1} = \mathbf{V}_{r1} \left(\mathbf{V}_{r0}^{-1}\boldsymbol{\mu}_{r0} + \sum_{i=1}^n \mathbf{X}'_{ri}\mathbf{V}_i^{-1}\tilde{\mathbf{y}}_i\right) = \left(\mathbf{V}_{r0}^{-1}\boldsymbol{\mu}_{r0} + \sum_{i=1}^n \mathbf{X}'_{ri}(\mathbf{X}_{ri}\boldsymbol{\Sigma}\mathbf{X}'_{ri} + \sigma^2\mathbf{I}_J)^{-1}(\mathbf{y}_i - \mathbf{X}_{fi}\boldsymbol{\beta}_f)\right)$$

The next steps are conditional draws of $\boldsymbol{\beta}_{ri}$, σ^2 , and $\boldsymbol{\Sigma}$. For these we revert back to the fully augmented framework above.

An implementation of this semi-augmented framework using the same artificial data as above is given in Matlab script `mod6_HNRM_v2`. The GS is provided in function `gs_HNRM_v2`.

HNRM application

Matlab scripts: `mod6_HNRM_outage`

Let's re-visit the commercial outage data, using a different regression model. Specifically, we draw a set of 50 size 2 and 3 firms, and regress their outage costs (in \$1000) against the scenario attributes "weekday", "daytime", and outage length (in minutes). The data file is `outage_50firms`.

(We'll ignore for now that there is considerable censoring in our data since 28% of cost reports are zero.)

Since it is very likely that not all firms are affected equally by these attributes, we will specify them as random coefficients with hierarchical priors.

1. Run script `mod6_HNRM_outage`. Looking at the GS diagnostics it is obvious that the covariances in Σ are essentially zero and / or poorly identified. This suggests running a model without covariance terms for the VCOV of the random coefficients (but still with variances).
2. The No-covariance model is implemented in script `mod6_HNRM_outage_nocovs`. It is identical to the first script, except we now specify separate (but identical) priors for the diagonal elements of Σ . We then use a slightly modified GS (called `gs_HNRM_nocovs`) that replaces draws of Σ from the *IW* with *kr* draws of variances from the *ig*. However, for efficiency purposes, we draw them in one block, i.e. we do not condition one variance on draws of the others.

The diagnostics for this version are clearly improved.

A formal comparison of the two models (full Σ and diagonal Σ) requires the computation or approximation of the marginal likelihood for each case. Here we will use Chib's (1995) method to simulate the marginal likelihood, adjusted for data augmented models as discussed in modules 4 and 5.

We first initiate the Chib routine for the unconstrained model (full Σ) in script `mod6_HNRM_outage_chib`, which calls the function `gs_HNRM_chib`. The output from this script are the marginal likelihood (in log form), as well as the three separate components that feed into it, i.e. the prior evaluated at the posterior mean, the sample likelihood evaluated at the posterior mean, and the posterior itself evaluated at the posterior mean (all in log form). We then repeat this procedure for the constrained model via script `mod6_HNRM_outage_nocovs_chib` and function `gs_HNRM_nocovs_chib`.

As you can see from the log files the log-Bayes Factor for the unconstrained vs. the constrained model is approximately 7.3, given our prior settings. According to the Table for Bayes Factor interpretation provided in Kass and Raftery (1995), this assigns *strong support* to the full model, despite the negligible magnitudes for its covariance terms (a log-BF of 6 or higher indicates strong support, 12 or higher = decisive support). One thing that "hurts" the restricted model is its larger prior space for the variances –

we now need three separate ig priors compared to a single IW prior. As you can see from the Chib results, this imposes a considerable penalty on the no-covariance specification.

Generating Posterior Predictive Densities with hierarchical models

Option 1: Using the fully augmented model:

When working with HMs, PPDs often require an additional simulation step to marginalize of the hierarchical elements. For our fully augmented HNRM, for example, the PPD for a hypothetical observation with attributes \mathbf{x}_p takes the following form:

$$\begin{aligned}
 p(\hat{y}_p | \mathbf{y}) &= \int_{\boldsymbol{\theta}} \left(\int_{\boldsymbol{\beta}_{ri}} p(\hat{y}_p, \boldsymbol{\theta}, \boldsymbol{\beta}_{ri} | \mathbf{y}) d\boldsymbol{\beta}_{ri} \right) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} \left(\int_{\boldsymbol{\beta}_{ri}} p(\hat{y}_p | \boldsymbol{\theta}, \boldsymbol{\beta}_{ri}, \mathbf{y}) p(\boldsymbol{\theta}, \boldsymbol{\beta}_{ri} | \mathbf{y}) d\boldsymbol{\beta}_{ri} \right) d\boldsymbol{\theta} = \\
 &\int_{\boldsymbol{\theta}} \left(\int_{\boldsymbol{\beta}_{ri}} p(\hat{y}_p | \boldsymbol{\theta}, \boldsymbol{\beta}_{ri}) p(\boldsymbol{\beta}_{ri} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\beta}_{ri} \right) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} \left(\int_{\boldsymbol{\beta}_{ri}} p(\hat{y}_p | \boldsymbol{\theta}, \boldsymbol{\beta}_{ri}) p(\boldsymbol{\beta}_{ri} | \boldsymbol{\theta}) d\boldsymbol{\beta}_{ri} \right) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}
 \end{aligned} \tag{29}$$

where $\boldsymbol{\theta}$ captures all actual parameters. This suggests the following implementation:

1. For a given set of parameters $\boldsymbol{\theta}$ from the original GS, draw at least 1 $\boldsymbol{\beta}_{ri}$ from its hierarchical prior $p(\boldsymbol{\beta}_{ri} | \boldsymbol{\beta}_r, \boldsymbol{\Sigma})$ (more draws will smoothen the PPD).
2. Draw a value of the predictive construct \hat{y}_p from the conditional likelihood $p(\hat{y}_p | \boldsymbol{\theta}, \boldsymbol{\beta}_{ri})$.
3. Repeat for all R draws of $\boldsymbol{\theta}$.

Option 2: Using the partially augmented model:

For our very special case of the normal hierarchical regression model, we can equivalently use:

$$\begin{aligned}
 p(\hat{y}_p | \mathbf{y}) &= \int_{\boldsymbol{\theta}} p(\hat{y}_p, \boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} p(\hat{y}_p | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} = \\
 &\int_{\boldsymbol{\theta}} p(\hat{y}_p | \boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2) p(\boldsymbol{\beta}_f, \boldsymbol{\beta}_r, \boldsymbol{\Sigma}, \sigma^2 | \mathbf{y}) d\boldsymbol{\theta}
 \end{aligned}$$

with the following implementation:

1. For a given set of parameters $\boldsymbol{\theta}$ from the original GS, draw one or more values of the predictive construct \hat{y}_p from the un-augmented (original) likelihood $p(\hat{y}_p | \boldsymbol{\theta})$.
2. Repeat for all R draws of $\boldsymbol{\theta}$.

Script `mod6_HNRM_outage_predict` shows both versions.

References:

- Chib, Siddhartha and Bradley P. Carlin. 1999. "On MCMC sampling in hierarchical longitudinal models." *Statistics and Computing* 9, 17-26.
- Kass, R. E. and A. E. Raftery. 1995. "Bayes Factors." *Journal of the American Statistical Association* 90, 773-795.