

Metropolis-Hastings Algorithm

(K Ch. 5, KPT Ch. 11)

Matlab scripts: mod7_Poisson_data, mod7_Poisson_rwc_prep, mod7_Poisson_rwc,
 mod7_Poisson_rwc_chib, mod7_Poisson_ic,
 mod7_Poisson_rwc_ACplots,
 mod7_Poisson_ic_ACplots, mod7_Poisson_fishing_rwc_prep,
 mod7_Poisson_fishin_rwc, mod7_Poisson_fishing_ic,
 mod7_HP_mudsnail, mod7_HP-mudsnail_parallel, AR1data, AR1ic

Matlab functions: gs_Poisson_rwc, gs_Poisson_rwc_chib, gs_Poisson_ic, gs_HP,
 gs_HPpar, Poisson_beta_mle, gs_mudsnail, HP_beta_mle,
 HP_gi_mle, gs_AR1_ic, rho_mode

Generally, whenever even the *conditional posterior kernel* for a parameter or set of parameters is unknown, we can no longer draw from it as would be the case in a "well-behaved" Gibbs Sampler. Instead, we need to approximate the unknown density via tools that have been specifically designed for this purpose, such as Importance Sampling, or, as is becoming increasingly popular given computational advances, a Metropolis-Hastings (MH) algorithm.

Let's assume the *unknown* conditional posterior for some parameter θ (scalar or vector) is generically given by $p(\theta | \Gamma, \mathbf{y}) = \frac{p(\theta, \mathbf{y} | \Gamma)}{p(\mathbf{y} | \Gamma)} = \frac{p(\theta)p(\mathbf{y} | \theta, \Gamma)}{p(\mathbf{y} | \Gamma)}$, where Γ are other model parameters that enter the conditional posterior. However, we do know its kernel $\tilde{p}(\theta | \Gamma, \mathbf{y}) = p(\theta)p(\mathbf{y} | \theta, \Gamma)$. Thus, the only unknown element is the normalizing constant $p(\mathbf{y} | \Gamma)$. The rationale of the MH algorithm is to use $\tilde{p}(\theta | \Gamma, \mathbf{y})$, plus a *candidate-generating density* (CGD) or *proposal density* $q(\theta)$ to obtain draws from the unknown $p(\theta | \Gamma, \mathbf{y})$. Note that $q(\theta)$ can also be a function of the data in addition to θ .

Suppose that the most recent draw of θ in the GS (which will initially be the starting draw) is θ^a . Now obtain a new "candidate draw" of θ , call it θ^b , from $q(\theta)$. Let's call this CGD generically $q(\theta)$. The new draw of θ is then accepted with probability

$$\alpha(\theta^a, \theta^b) = \min\left(\frac{p(\theta^b | \Gamma, \mathbf{y})q(\theta^a)}{p(\theta^a | \Gamma, \mathbf{y})q(\theta^b)}, 1\right) = \min\left(\frac{\tilde{p}(\theta^b | \Gamma, \mathbf{y})q(\theta^a)}{\tilde{p}(\theta^a | \Gamma, \mathbf{y})q(\theta^b)}, 1\right), \quad (1)$$

since the normalizing constant $p(\mathbf{y} | \Gamma)$ cancels out in the ratio. We usually work with logs:

$$\begin{aligned} \log(\alpha(\theta^a, \theta^b)) &= \min(\log \tilde{p}(\theta^b | \Gamma, \mathbf{y}) + \log q(\theta^a) - \log \tilde{p}(\theta^a | \Gamma, \mathbf{y}) - \log q(\theta^b), 0) = \\ &= \min((\log \tilde{p}(\theta^b | \Gamma, \mathbf{y}) - \log q(\theta^b)) - (\log \tilde{p}(\theta^a | \Gamma, \mathbf{y}) - \log q(\theta^a)), 0) \end{aligned} \quad (2)$$

In practice this is implemented by comparing the α value to a random uniform $[0,1]$ draw (or, equivalently, $\log(\alpha)$ to the log of a uniform draw). If α exceeds the random value, the new draw θ^b is accepted, otherwise θ^a remains the most current draw. As discussed in Gelman et al (Ch. 12), Koop (Ch. 5) and KPT, Ch. 11, after a sufficient number of “burn-ins” the sequence of draws of θ will converge to the desired underlying density.

If the proposal density is symmetric, i.e. $q(\theta^a) = q(\theta^b)$ the acceptance probability reduces to

$$\alpha(\theta^a, \theta^b) = \min\left(\frac{\tilde{p}(\theta^b | \Gamma, \mathbf{y})}{\tilde{p}(\theta^a | \Gamma, \mathbf{y})}, 1\right) \quad (3)$$

Also note that the generic GS is a special case of the MH with $q(\theta^j) = p(\theta^j | \Gamma, \mathbf{y})$, $j = a, b$ such that the acceptance probability is always one, i.e. every new draw of θ is accepted by default.

There are two basic types of MH algorithms – the *Random Walk Chain MH (RWMH)* and the *Independence Chain MH (IMH)*. In the *RWMH* the proposal density is a function of the current draw of θ , i.e. takes the form of $q(\theta^b | \theta^a)$. For example, the current draw θ^a might be designated the mean of the proposal density. In this case, the acceptance probability takes the form of

$$\alpha(\theta^a, \theta^b) = \min\left(\frac{\tilde{p}(\theta^b | \Gamma, \mathbf{y})q(\theta^a | \theta^b)}{\tilde{p}(\theta^a | \Gamma, \mathbf{y})q(\theta^b | \theta^a)}, 1\right) \quad (4)$$

Usually, $q(\theta^b | \theta^a)$ is chosen to be symmetric to yield the simplified expression

$$\alpha(\theta^a, \theta^b) = \min\left(\frac{\tilde{p}(\theta^b | \Gamma, \mathbf{y})}{\tilde{p}(\theta^a | \Gamma, \mathbf{y})}, 1\right) \quad (5)$$

For the *IMH* the proposal density is not a direct function of the most recent draw of θ . A popular choice is a t-density or multivariate t-density with the mean set to the MLE solution (mode) of the posterior kernel $\tilde{p}(\theta | \Gamma, \mathbf{y})$ and the variance set to the (possibly scaled) inverted negative Hessian coming out of the MLE sub-routine. Naturally, this is more complex and takes longer than the *RWMH*, but it tends to reduce auto-correlation compared to draws from the *RWMH*.

How can we influence acceptance rates? In the *IMH*, we can choose the scalar for the Inverted Hessian to boost or reduce the variance matrix in the t-density, and also the degrees of freedom for the t density. In the *RWMH*, we usually specify $q(\theta^b | \theta^a)$ to be a normal density or multivariate normal density with mean θ^a and an arbitrary variance. This variance is specified at the onset, along with the priors, then refined after a few “practice runs” to achieve optimal acceptance rates.

To assure that the MH algorithm covers the entire posterior kernel (with more weight given to areas with higher density), Gelman et al recommend an optimal acceptance rate of about 40-45% of all draws for a single-valued θ , and about 20-25% if θ is a vector. **However, these target thresholds only apply to the RWMH with a normal CGD!** For any other version of a MH algorithm IEF scores and convergence

diagnostics should be used for guidance on the optimal acceptance rate. For example, for the Poisson fishing application, boosting acceptance rates in the IMH from 35% to 75% doubled efficiency (= halved IEF scores), but in some other cases acceptance rates in the 30-50% may be sufficient.

An informal “proof” of the MH algorithm based on KPT, p. 153, is given at the end of this chapter.

Example 1: Parameterized Poisson Model

Consider a Poisson model with parameterized expectation. These models are useful in contexts such as visitation counts to recreation sites, patient visits to medical facilities, etc. The basic Poisson density has a single parameter, λ , which denotes both the expectation and variance of the distribution. Note that while a Poisson random variable is a non-negative integer, λ can be any positive real number.

The explanatory variables enter the scene as linear regressors for the logged expectation. Thus we have:

$$f(y_i | \lambda_i) = \frac{\exp(-\lambda_i) \lambda_i^{y_i}}{y_i!} \quad y_i = 0, 1, 2, \dots, \quad \lambda_i > 0 \quad \text{where} \quad (6)$$

$$\lambda_i = \exp(\mathbf{x}_i' \boldsymbol{\beta})$$

The interpretation of the elements of $\boldsymbol{\beta}$ is thus a fractional change in expected outcome due to a unit change in the corresponding explanatory variable. Multiplying by 100 gives the percentage change.

The likelihood function is given as

$$p(\mathbf{y} | \boldsymbol{\beta}) = \prod_{i=1}^N \frac{\exp(-\lambda_i) \lambda_i^{y_i}}{y_i!} = \prod_{i=1}^N \frac{\exp(-\exp(\mathbf{x}_i' \boldsymbol{\beta})) (\exp(\mathbf{x}_i' \boldsymbol{\beta}))^{y_i}}{y_i!} \quad (7)$$

In absence of any theoretical restrictions on marginal effects, we can choose the usual multivariate normal priors for $\boldsymbol{\beta}$:

$$p(\boldsymbol{\beta}) = (2\pi)^{-k/2} |\mathbf{V}_0|^{-1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) \quad (8)$$

where k is the dimension of $\boldsymbol{\beta}$.

The posterior kernel is thus given as

$$p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) \propto \tilde{p}(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) * \prod_{i=1}^N \frac{\exp(-\exp(\mathbf{x}_i' \boldsymbol{\beta})) (\exp(\mathbf{x}_i' \boldsymbol{\beta}))^{y_i}}{y_i!} \quad (9)$$

This is not the kernel of a well-understood density, so we cannot obtain direct draws from it or derive analytical results.

Instead, we can apply a M-H procedure where we draw a candidate vector $\boldsymbol{\beta}_c$ from some known candidate-generating function $q(\boldsymbol{\beta}_c)$, and retain it with acceptance probability $\alpha(\boldsymbol{\beta}_o \rightarrow \boldsymbol{\beta}_c)$, where $\boldsymbol{\beta}_o$ is the "old" or "current" draw. In the most general sense, this CDG can be a function of $\boldsymbol{\beta}_o$, and / or the data, i.e. $q(\boldsymbol{\beta}_c) = q(\boldsymbol{\beta}_c | \boldsymbol{\beta}_o, \mathbf{y}, \mathbf{X})$. Different "flavors" of the M-H step evolve from different choices of $q(\cdot)$.

Version 1: Random Walk Chain MH

For the RWMH we choose a multivariate normal CGD with mean equal to the current draw $\boldsymbol{\beta}_o$ (which, initially will be the starting draw), and variance matrix equal to cV_c , where c is some arbitrary scalar (initially set to something small, like 0.01) and V_c is initially set to an identity matrix.

Therefore, we have $q(\boldsymbol{\beta}_c | \boldsymbol{\beta}_o) = n(\boldsymbol{\beta}_o, cV_c) = q(\boldsymbol{\beta}_o | \boldsymbol{\beta}_c)$ due to the symmetry property of the normal density.

The formula for the log (acceptance probability) thus reduces to

$$\begin{aligned} \log(\alpha(\boldsymbol{\beta}_o \rightarrow \boldsymbol{\beta}_c)) &= \min\left(\log\left(\frac{\tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X})q(\boldsymbol{\beta}_o | \boldsymbol{\beta}_c)}{\tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X})q(\boldsymbol{\beta}_c | \boldsymbol{\beta}_o)}\right), 0\right) = \min\left(\log\left(\frac{\tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X})}{\tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X})}\right), 0\right) = \\ &= \min(\log \tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) - \log \tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}), 0), \quad \text{where} \\ \log \tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) &= \sum_{i=1}^N (-\exp(\mathbf{x}'_i \boldsymbol{\beta}_c) + y_i (\mathbf{x}'_i \boldsymbol{\beta}_c) - \log y_i!) + \log\left((\boldsymbol{\beta}_c - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta}_c - \boldsymbol{\mu}_0)\right) \\ \log \tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}) &= \sum_{i=1}^N (-\exp(\mathbf{x}'_i \boldsymbol{\beta}_o) + y_i (\mathbf{x}'_i \boldsymbol{\beta}_o) - \log y_i!) + \log\left((\boldsymbol{\beta}_o - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta}_o - \boldsymbol{\mu}_0)\right) \end{aligned} \tag{10}$$

We first run a "prep" script " with a generous number of burn-ins that simply aims at getting non-ridiculous acceptance rates (say 10% or more). You may run it a few times, experimenting with the setting for c . As a general rule, acceptance rates will increase with smaller settings for c .

Matlab script `mod7_Poisson_rwc_prep` and function `gs_Poisson_rwc` illustrate this process.

Then run the main script, this time using the estimated posterior variance matrix for $\boldsymbol{\beta}$ from the prep – results for V_c . Start with $c = 1$ and experiment with different settings for c until you have a reasonable acceptance rate, ideally in the 20-25% range.

Matlab script `mod7_Poisson_rwc` illustrates this process.

You may want to repeat this step a few times –working with an ever more refined estimate of V_c . Also, don't be shy with your burn-ins – often several 100,000 's may be required until CD diagnostics check out.

Version 2: Independence Chain MH

In the *IMH* the CGD is no longer a function of the current draw $\boldsymbol{\beta}_o$, but rather – at least indirectly – of the data \mathbf{y} , \mathbf{X} . A popular version is a multivariate t distribution, with mean equal to the mode of the posterior kernel $\tilde{p}(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X})$, call it $\hat{\boldsymbol{\beta}}$, and variance ("scale matrix") equal to inverted negative Hessian flowing from the MLE routine, evaluated at $\hat{\boldsymbol{\beta}}$ and possibly scaled by some arbitrary scalar "c" (similar to the scalar in the RWMH version). Also, the degree-of-freedom parameter for the *mvt* (call it ν) is specified at the onset along with the model priors. A smaller value of ν (say in the 4-8 range) implies heavier tails, which is generally desirable. However, if ν is chosen too small, too many candidate draws will be rejected.

Thus, in each round of the sampler, we first have to run a sub-routine MLE procedure to find $\hat{\boldsymbol{\beta}}$ as

$$\hat{\boldsymbol{\beta}} = \max_{\boldsymbol{\beta}} \left[\sum_{i=1}^N (-\exp(\mathbf{x}'_i \boldsymbol{\beta}) + y_i (\mathbf{x}'_i \boldsymbol{\beta}) - \log y_i!) + \log \left((\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0) \right) \right] \quad (11)$$

It usually helps in terms of speed and precision to supply the analytical gradient and Hessian for this step. However, in most cases software-supplied numerical approximations for these constructs will be sufficient to obtain $\hat{\boldsymbol{\beta}}$ in a few iterations.

The MLE routine thus produces $\hat{\boldsymbol{\beta}}$ and $V_c = \left(-H(\hat{\boldsymbol{\beta}}) \right)^{-1}$ where $H(\cdot)$ is the Hessian matrix evaluated at the MLE solution. We can now draw a candidate vector $\boldsymbol{\beta}_c$ from the *mvt* density, i.e.

$$\boldsymbol{\beta}_c \sim mvt(\hat{\boldsymbol{\beta}}, c\mathbf{V}_c, \nu) = q(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) \quad (12)$$

Since this CDG is no longer symmetric, we need to explicitly consider it in the acceptance probability, i.e.

$$\begin{aligned} \log(\alpha(\boldsymbol{\beta}_o \rightarrow \boldsymbol{\beta}_c)) &= \min \left(\log \left(\frac{\tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) q(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X})}{\tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}) q(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X})} \right), 0 \right) = \\ &= \min \left((\log \tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) - \log q(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X})) - (\log \tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}) - \log q(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}))), 0 \right), \quad \text{where} \\ \log \tilde{p}(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) &= \sum_{i=1}^N (-\exp(\mathbf{x}'_i \boldsymbol{\beta}_c) + y_i (\mathbf{x}'_i \boldsymbol{\beta}_c) - \log y_i!) + \log \left((\boldsymbol{\beta}_c - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta}_c - \boldsymbol{\mu}_0) \right) \\ \log \tilde{p}(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}) &= \sum_{i=1}^N (-\exp(\mathbf{x}'_i \boldsymbol{\beta}_o) + y_i (\mathbf{x}'_i \boldsymbol{\beta}_o) - \log y_i!) + \log \left((\boldsymbol{\beta}_o - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta}_o - \boldsymbol{\mu}_0) \right) \\ \log q(\boldsymbol{\beta}_c | \mathbf{y}, \mathbf{X}) &= \log \left(mvt(\boldsymbol{\beta}_c | \hat{\boldsymbol{\beta}}, c\mathbf{V}_c, \nu) \right) \\ \log q(\boldsymbol{\beta}_o | \mathbf{y}, \mathbf{X}) &= \log \left(mvt(\boldsymbol{\beta}_o | \hat{\boldsymbol{\beta}}, c\mathbf{V}_c, \nu) \right) \end{aligned} \quad (13)$$

Matlab script `mod7_Poisson_ic` with functions `gs_Poisson_ic` and `Poisson_beta_mle` show the implementation of this basic IMH.

You can see that compared to the RWC version, the IC version takes much longer per iteration, but converges much faster, such that overall runtime is comparable between the two approaches, at least for this application.

There are clear efficiency gains from using the IC approach, as is evident by inspection of the IEF statistics and the autocorrelation plots. Thus, breaking the direct dependence of the old and new draw via the RWC proposal density has the expected effect of reducing correlation in the Markov Chain.

Example 2: Hierarchical Poisson Model

Knowledge of the MH technique also enables us to tackle a broader variety of hierarchical models. Specifically, we can combine elements of the preceding Poisson model with those of the HNRM of Module 6 to specify and estimate a Hierarchical Poisson . For a recent application of this model see Davis and Moeltner (2010).

Here is the outline of Davis and Moeltner's hierarchical Poisson model with normally distributed random coefficients in the mean function – also called the Poisson-lognormal model. A detailed discussion on the Bayesian estimation of this model is given in Chib et al. (1998).

Consider an individual i that takes trips to $j=1 \dots J$ sites during a specific time period. Let the density function for trips be Poisson with the usual parameterized mean function. The hierarchical structure is added by letting some parameters in the mean function be "random", i.e. follow a (normal) hierarchical distribution. Thus, the structural model and likelihood are given as:

$$f(y_{ij} | \lambda_{ij}) = \frac{\exp(-\lambda_{ij}) \lambda_{ij}^{y_{ij}}}{y_{ij}!} \quad \text{where}$$

$$\lambda_{ij} = \exp(\mathbf{x}'_{ij} \boldsymbol{\beta} + \mathbf{h}'_{ij} \boldsymbol{\gamma}_i) \quad \text{and} \quad \boldsymbol{\gamma}_i \sim mvn(\boldsymbol{\gamma}, \boldsymbol{\Sigma}). \quad (14)$$

$$p(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\Sigma}) = \prod_{i=1}^N \left(\int_{\boldsymbol{\gamma}_i} \left(\prod_{j=1}^J \frac{\exp(-\lambda_{ij}) \lambda_{ij}^{y_{ij}}}{y_{ij}!} \right) f(\boldsymbol{\gamma}_i | \boldsymbol{\gamma}, \boldsymbol{\Sigma}) d\boldsymbol{\gamma}_i \right)$$

We choose the usual multivariate normal priors for the fixed effects ($\boldsymbol{\beta}$) and the hierarchical expectation of the random effects ($\boldsymbol{\gamma}$), and an IW prior for the variance-covariance matrix of the random effects ($\boldsymbol{\Sigma}$).

$$p(\boldsymbol{\beta}) = (2\pi)^{-k_f/2} |\mathbf{V}_\beta|^{-1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)' \mathbf{V}_\beta^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)\right)$$

$$p(\boldsymbol{\gamma}) = (2\pi)^{-k_r/2} |\mathbf{V}_\gamma|^{-1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)' \mathbf{V}_\gamma^{-1} (\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)\right) \quad (15)$$

$$p(\boldsymbol{\Sigma}) = \left(2^{v_0 k_r / 2} \pi^{k_r(k_r-1)/4} \prod_{i=1}^{k_r} \Gamma\left(\frac{v_0 + 1 - i}{2}\right) \right)^{-1} * |\mathbf{S}_0|^{v_0/2} |\boldsymbol{\Sigma}|^{-(v_0+k_r+1)/2} \exp\left(-\frac{1}{2} tr(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right)$$

It would be difficult to break the resulting posterior into easy-to-draw-from conditional components, due to the multi-dimensional integral over the random effects in the likelihood function. We can again facilitate posterior simulation by augmenting the parameter space with the individual-level random effects, γ_i , $i=1 \dots N$. The augmented posterior takes the following generic form:

$$\begin{aligned}
& p(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\gamma}, \{\gamma_i\}_{i=1}^N | \mathbf{y}, \mathbf{X}) \propto \\
& p(\boldsymbol{\beta}) p(\{\gamma_i\}_{i=1}^N, \boldsymbol{\gamma}, \boldsymbol{\Sigma}) p(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\gamma}, \gamma_i) = \\
& p(\boldsymbol{\beta}) p(\boldsymbol{\gamma}) p(\boldsymbol{\Sigma}) \prod_{i=1}^N p(\gamma_i | \boldsymbol{\gamma}, \boldsymbol{\Sigma}) p(\mathbf{y} | \boldsymbol{\beta}, \gamma_i)
\end{aligned} \tag{16}$$

The last component of (32) (I call it the conditional likelihood, as it conditions on the random effects) has a well-known closed form (the basic Poisson density). This circumvents the approximation of the multi-dimensional integral in the original likelihood function (as would be necessary in a classical estimation routine, e.g. via MLE).

Here is the explicit form of the augmented posterior:

$$\begin{aligned}
& p(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\gamma}, \gamma_i | \mathbf{y}, \mathbf{X}) \propto \\
& \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)' \mathbf{V}_\beta^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)\right) * \exp\left(-\frac{1}{2}(\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)' \mathbf{V}_\gamma^{-1} (\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)\right) * \\
& |\boldsymbol{\Sigma}|^{-(v_0 + k_r + 1)/2} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1})\right) * \\
& \prod_{i=1}^N |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\gamma_i - \boldsymbol{\gamma})' \boldsymbol{\Sigma}^{-1} (\gamma_i - \boldsymbol{\gamma})\right) * \prod_{i=1}^N \prod_{j=1}^J \frac{\exp(-\lambda_{ij}) \lambda_{ij}^{y_{ij}}}{y_{ij}!} \quad \text{with} \\
& \lambda_{ij} = \exp(\mathbf{x}'_{ij} \boldsymbol{\beta} + \mathbf{h}'_{ij} \boldsymbol{\gamma}_i)
\end{aligned} \tag{17}$$

The Gibbs Sampler proceeds in four blocks, two of which require a built-in Metropolis-Hastings (MH) routine. Details and options for the MH steps can also be found in Chib et al. (1998).

Step 1: Draw $\boldsymbol{\beta}$:

$$\begin{aligned}
& p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}, \boldsymbol{\Sigma}, \gamma_i) \propto \\
& \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)' \mathbf{V}_\beta^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)\right) * \prod_{i=1}^N \prod_{j=1}^J \frac{\exp(-\lambda_{ij}) \lambda_{ij}^{y_{ij}}}{y_{ij}!} \quad \text{with} \\
& \lambda_{ij} = \exp(\mathbf{x}'_{ij} \boldsymbol{\beta} + \mathbf{h}'_{ij} \boldsymbol{\gamma}_i)
\end{aligned} \tag{18}$$

➔ MH (Davis & Moeltner use a Independence-Chain MH with a *mvt* proposal function)

Step 2: Draw the individual vectors of random effects (repeat N times)

$$p(\boldsymbol{\gamma}_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\gamma}, \mathbf{y}_i, \mathbf{X}_i) \propto \exp\left(-\frac{1}{2}(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})\right) \prod_{i=1}^N \prod_{j=1}^J \frac{\exp(-\lambda_{ij}) \lambda_{ij}^{y_{ij}}}{y_{ij}!} \quad (19)$$

➔ MH (Davis & Moeltner use a Independence-Chain MH with a *mvt* proposal function)

Step 3: Draw the vector of random effect expectations

$$p(\boldsymbol{\gamma} | \boldsymbol{\gamma}_i, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2}(\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)' \mathbf{V}_\gamma^{-1}(\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)\right) * \prod_{i=1}^N \exp\left(-\frac{1}{2}(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})\right) = \exp\left(-\frac{1}{2}\left((\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma)' \mathbf{V}_\gamma^{-1}(\boldsymbol{\gamma} - \boldsymbol{\mu}_\gamma) + \sum_{i=1}^N (\boldsymbol{\gamma}_i - \boldsymbol{\gamma})' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})\right)\right) \quad (20)$$

Thus:

$$\boldsymbol{\gamma} | \boldsymbol{\gamma}_i, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_1, \mathbf{V}_1)$$

$$\mathbf{V}_1 = (\mathbf{V}_\gamma^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left(\mathbf{V}_\gamma^{-1} \boldsymbol{\mu}_\gamma + \boldsymbol{\Sigma}^{-1} \sum_{i=1}^N \boldsymbol{\gamma}_i \right)$$

Step 4: Draw the random effect variance matrix

$$p(\boldsymbol{\Sigma} | \boldsymbol{\gamma}, \boldsymbol{\gamma}_i) \propto |\boldsymbol{\Sigma}|^{-(v_0 + k_r + 1 + N)/2} \exp\left(-\frac{1}{2}\left(\text{tr}(\mathbf{S}_0 \cdot \boldsymbol{\Sigma}^{-1}) + \sum_{i=1}^N (\boldsymbol{\gamma}_i - \boldsymbol{\gamma})' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})\right)\right) \quad (21)$$

Thus:

$$\boldsymbol{\Sigma} | \boldsymbol{\gamma}, \boldsymbol{\gamma}_i \sim IW(v_1, \mathbf{S}_1)$$

$$v_1 = v_0 + N \quad \mathbf{S}_1 = \mathbf{S}_0 + \sum_{i=1}^N (\boldsymbol{\gamma}_i - \boldsymbol{\gamma})(\boldsymbol{\gamma}_i - \boldsymbol{\gamma})'$$

As can be seen from the output file, the IEF scores are uncomfortable high for some parameters, and longer run-time plus thinning would have likely been a good idea in this case. However, posterior results for the key parameters are plausible, as are those for the posterior predictive densities (see paper).

Example 3: AR-1 Model

(see separate pdf file)

Some caveats when working with MH:

- Any *MH* step in your GS will generally hurt the efficiency of your posterior sampler, so expect higher IEF scores and lower m^* 's. This effect is especially pronounced if you use an *RWMH* algorithm.
- IEF scores in the 10-30 range may often be as good as it gets, especially for the *RWMH*. As a general rule you can expect higher autocorrelation for the *RWMH* compared to the *IMH*. This trade-off comes with a price tag: The *IMH* usually takes much, much, longer to run. However, you may get away with a substantially fewer burn-ins! (see Matlab application).
- Try running the *IMH* without analytical gradient and Hessian in the MLE sub-routine. If the Sampler has trouble converging or is simply moving too slowly, you may need to spend some extra “pencil & paper” time to derive the analytical gradient and Hessian.
- Save time and avoid frustration by first running the *IMH* with just a few hundred iterations, then check your acceptance rates. There is no “fast & steady” rule on *acceptance rates for the IMH*. In theory, since the chain is independent, you want all draws to be accepted. However, this may still lead to high autocorrelation in the sequence even in the *IMH*. Let the IEF scores guide you in this respect. I have found that for most applications acceptance rates between 0.5 and 0.8 lead to best efficiency.
- Other “remedies” include experimenting with different parameter blockings and re-parameterization of the structural model.

Appendix A: Why / How MH works and the Equality of Joint Probabilities

(KPT p. 153)

Suppose the discrete RV θ has pmf $p(\theta)$. Let θ^{t-1} be the current value of the MH chain and assume that we can be certain that $\theta^{t-1} \sim p$. Now assume the next draw produced by the MH algorithm is θ^t . Show that $\theta^t \sim p$.¹

Let θ^a and θ^b be two distinct points in the support of θ . Then the MH acceptance probability of *moving from θ^b to θ^a* is given as

$$\alpha(\theta^a | \theta^b) = \min \left\{ 1, \frac{p(\theta^a)q(\theta^b | \theta^a)}{p(\theta^b)q(\theta^a | \theta^b)} \right\} \quad (22)$$

¹ Naturally, this would be trivial if we had $\theta^t = \theta^{t-1}$.

Without loss in generality assume that $\theta^t \neq \theta^{t-1}$ and that $\alpha(\theta^a | \theta^b) > 1$ such that its reciprocal $\alpha(\theta^b | \theta^a) < 1$. Now consider the joint probability of observing the sequence θ^b, θ^a in the MH algorithm, i.e.

$$\Pr(\theta^{t-1} = \theta^b, \theta^t = \theta^a) = \Pr(\theta^t = \theta^a | \theta^{t-1} = \theta^b) \Pr(\theta^{t-1} = \theta^b) \quad (23)$$

The first term on the right hand side is the conditional probability of arriving at θ^a when the current draw is θ^b . This conditional can be expressed as the probability of *drawing* θ^a when the current draw is θ^b times the probability of accepting θ^a once it's drawn, i.e.

$$\Pr(\theta^t = \theta^a | \theta^{t-1} = \theta^b) = q(\theta^a | \theta^b) \alpha(\theta^a | \theta^b) = q(\theta^a | \theta^b) \quad (24)$$

where the last equality follows from our assumption that $\alpha(\theta^a | \theta^b) > 1$. Thus we get

$$\Pr(\theta^{t-1} = \theta^b, \theta^t = \theta^a) = q(\theta^a | \theta^b) \Pr(\theta^{t-1} = \theta^b) = q(\theta^a | \theta^b) p(\theta^b) \quad (25)$$

Now consider a different joint probability

$$\begin{aligned} \Pr(\theta^{t-1} = \theta^a, \theta^t = \theta^b) &= \Pr(\theta^t = \theta^b | \theta^{t-1} = \theta^a) \Pr(\theta^{t-1} = \theta^a) = \\ &= q(\theta^b | \theta^a) \alpha(\theta^b | \theta^a) p(\theta^a) = \\ &= q(\theta^b | \theta^a) \frac{p(\theta^b) q(\theta^a | \theta^b)}{p(\theta^a) q(\theta^b | \theta^a)} p(\theta^a) = p(\theta^b) q(\theta^a | \theta^b) \end{aligned} \quad (26)$$

Thus the two joint probabilities are indeed equal. This *equality of joint probabilities* implies that

$$\Pr(\theta^t = \theta^a | \theta^{t-1} = \theta^b) \Pr(\theta^{t-1} = \theta^b) = \Pr(\theta^t = \theta^b | \theta^{t-1} = \theta^a) \Pr(\theta^{t-1} = \theta^a) \quad (27)$$

Now summing over θ^a (i.e. marginalizing the joint pmf over θ^a) yields

$$\Pr(\theta^{t-1} = \theta^b) \sum_{\theta^a = \theta_{\min}}^{\theta_{\max}} \Pr(\theta^t = \theta^a | \theta^{t-1} = \theta^b) = \sum_{\theta^a = \theta_{\min}}^{\theta_{\max}} \Pr(\theta^t = \theta^b | \theta^{t-1} = \theta^a) \Pr(\theta^{t-1} = \theta^a) \quad (28)$$

The second term on the left hand side is just the full *cdf* of a conditional density for a discrete variable (summing to 1), whereas the right hand side is the definition for a marginal density for discrete variables. Thus we get:

$$\Pr(\theta^{t-1} = \theta^b) = \Pr(\theta^t = \theta^b) \quad (29)$$

Since $\theta^{t-1} \sim p(\theta)$ it follows that $\theta^t \sim p(\theta)$ as well.

Appendix B: Marginal Likelihood Estimation for Models with MH Components

Chib and Jeliazkov (2001)

To compare models that use MH in their posterior simulator we need to derive their marginal likelihood. This requires a modification of the Chib (1995) method, as illustrated in Chib and Jeliazkov (2001), henceforth *C/J*.

As before, we can write the log mL , evaluated at point $\bar{\theta}$ (e.g. the posterior mean) as

$$\log p(\mathbf{y}) = \log p(\mathbf{y} | \bar{\theta}) + \log p(\bar{\theta}) - \log p(\bar{\theta} | \mathbf{y}) \quad (1)$$

The trick is to find an estimate of the posterior ordinate $p(\bar{\theta} | \mathbf{y})$. Assume that one or more of the normalizing constants of full conditional densities are not known (that's why we're using an *MH* routine in the posterior sampler). This is where the *C/J* method comes in.

One Block Sampling

Suppose the posterior density $p(\theta | \mathbf{y}) \propto p(\theta)p(\mathbf{y} | \theta)$ is sampled in one block via *MH*². Our goal is again to estimate the posterior ordinate at the posterior mean, i.e. $p(\bar{\theta} | \mathbf{y})$.

Let $q(\theta^a, \theta^b | \mathbf{y})$ be the proposal density for the transition from θ^a to θ^b (the dependence of q on \mathbf{y} denotes the most general case and may not always apply). Then, as discussed before, the probability of move (of accepting θ^b) is given by

$$\alpha(\theta^a, \theta^b | \mathbf{y}) = \min \left\{ 1, \frac{p(\mathbf{y} | \theta^b)p(\theta^b)q(\theta^b, \theta^a | \mathbf{y})}{p(\mathbf{y} | \theta^a)p(\theta^a)q(\theta^a, \theta^b | \mathbf{y})} \right\} \quad (2)$$

For any two points θ^0 and $\bar{\theta}$ denote the joint probability of observing the sequence $\theta^0, \bar{\theta}$ as

$$\Pr(\theta^0, \bar{\theta} | \mathbf{y}) = \Pr(\bar{\theta} | \mathbf{y}, \theta^0) \Pr(\theta^0 | \mathbf{y}) = q(\theta^0, \bar{\theta} | \mathbf{y}) \alpha(\theta^0, \bar{\theta} | \mathbf{y}) p(\theta^0 | \mathbf{y}) = \frac{p(\mathbf{y} | \bar{\theta}) p(\bar{\theta}) q(\bar{\theta}, \theta^0 | \mathbf{y})}{p(\mathbf{y} | \theta^0) p(\theta^0)} p(\theta^0 | \mathbf{y}) \quad (3)$$

where it is assumed that θ^0 is a draw from the posterior density $p(\theta^0 | \mathbf{y})$. From KPT, p. 153ff (see appendix Lecture notes 8) we know that this joint density is reversible, i.e. $\Pr(\theta^0, \bar{\theta} | \mathbf{y}) = \Pr(\bar{\theta}, \theta^0 | \mathbf{y})$, or

$$q(\theta^0, \bar{\theta} | \mathbf{y}) \alpha(\theta^0, \bar{\theta} | \mathbf{y}) p(\theta^0 | \mathbf{y}) = q(\bar{\theta}, \theta^0 | \mathbf{y}) \alpha(\bar{\theta}, \theta^0 | \mathbf{y}) p(\bar{\theta} | \mathbf{y}) \quad (4)$$

Thus, for a *specific* draw of θ^0 we have

² We encountered this case in the generic Poisson model.

$$p(\bar{\boldsymbol{\theta}} | \mathbf{y}, \boldsymbol{\theta}^0) = \frac{q(\boldsymbol{\theta}^0, \bar{\boldsymbol{\theta}} | \mathbf{y}) \alpha(\boldsymbol{\theta}^0, \bar{\boldsymbol{\theta}} | \mathbf{y}) p(\boldsymbol{\theta}^0 | \mathbf{y})}{q(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^0 | \mathbf{y}) \alpha(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^0 | \mathbf{y})} \quad (5)$$

Marginalizing over all possible values of $\boldsymbol{\theta}^0$, i.e. over $\boldsymbol{\theta}$, this yields

$$p(\bar{\boldsymbol{\theta}} | \mathbf{y}) = \frac{\int q(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}} | \mathbf{y}) \alpha(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}} | \mathbf{y}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}}{\int \alpha(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} | \mathbf{y}) q(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}} \quad (6)$$

To highlight the estimation strategy, we can write this in terms of expectations as

$$p(\bar{\boldsymbol{\theta}} | \mathbf{y}) = \frac{E_{p(\boldsymbol{\theta} | \mathbf{y})} \{q(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}} | \mathbf{y}) \alpha(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}} | \mathbf{y})\}}{E_{q(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} | \mathbf{y})} \{\alpha(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} | \mathbf{y})\}} \quad (7)$$

This yields a simulation-consistent estimate of

$$\hat{p}(\bar{\boldsymbol{\theta}} | \mathbf{y}) = \frac{M^{-1} \sum_{g=1}^M q(\boldsymbol{\theta}^g, \bar{\boldsymbol{\theta}} | \mathbf{y}) \alpha(\boldsymbol{\theta}^g, \bar{\boldsymbol{\theta}} | \mathbf{y})}{J^{-1} \sum_{j=1}^J \alpha(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta}^j | \mathbf{y})} \quad (8)$$

where $\boldsymbol{\theta}^g$ are the sampled draws from the posterior distribution (from the original sampler run) and $\boldsymbol{\theta}^j$ are draws from $q(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} | \mathbf{y})$ given the fixed value $\bar{\boldsymbol{\theta}}$, using a follow-up routine to the main sampler (since we won't know $\bar{\boldsymbol{\theta}}$ until we're done with the main run). Note that in practice we usually set $M = J$.

Note: If S (the area over which the posterior of $\boldsymbol{\theta}$ is defined) is a proper *subset* of \mathfrak{R}^d (the dimension of $\boldsymbol{\theta}$) then all $\boldsymbol{\theta}^j$'s from $q(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} | \mathbf{y})$ that fall outside of S are included in the average of the denominator of (8), with $\alpha(\cdot) = 0$. (Example: some of the elements of $\boldsymbol{\theta}$ might be truncated).

For an example of this single-block application see `mod7_Poisson_rwc_chib`.

Two Parameter Blocks and Multiple Latent Variable Blocks

For simplicity assume the normalizing constant of only the *first* full conditional is not known, and that this density is sampled by *MH*.

The log-ML is given by

$$\log m(\mathbf{y}) = \log p(\mathbf{y} | \bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2) + \log p(\bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2) - \log p(\bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2 | \mathbf{y}) \quad (9)$$

and the objective is to estimate $p(\bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2 | \mathbf{y})$. We assume that the likelihood ordinate is readily available either by direct computation or by a Monte Carlo Integration method. As in Chib (1995) we decompose the posterior ordinate as

$$p(\bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2 | \mathbf{y}) = p(\bar{\boldsymbol{\theta}}_1 | \mathbf{y}) p(\bar{\boldsymbol{\theta}}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1) \quad (10)$$

where the marginal density ordinate $p(\bar{\theta}_1 | \mathbf{y})$ cannot be estimated as usual via Monte Carlo integration since, by assumption, the normalizing constant of $p(\theta_1 | \mathbf{y}, \theta_2, \mathbf{z})$ is not known.

Let $q(\theta_1^a, \theta_1^b | \mathbf{y}, \theta_2, \mathbf{z})$ be the proposal density for transition from θ_1^a to θ_1^b (where for generality q is allowed to depend on data and the two remaining blocks). The probability of move is

$$\alpha(\theta_1^a, \theta_1^b | \mathbf{y}, \theta_2, \mathbf{z}) = \min \left\{ 1, \frac{p(\mathbf{y} | \theta_1^b, \theta_2, \mathbf{z}) p(\theta_1^b, \theta_2) q(\theta_1^b, \theta_1^a | \mathbf{y}, \theta_2, \mathbf{z})}{p(\mathbf{y} | \theta_1^a, \theta_2, \mathbf{z}) p(\theta_1^a, \theta_2) q(\theta_1^a, \theta_1^b | \mathbf{y}, \theta_2, \mathbf{z})} \right\} \quad (11)$$

Consider again the joint probability

$$\Pr(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) = q(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) \quad (12)$$

Using again the rule of reversibility (or “local reversibility” in this case) we get

$$\begin{aligned} q(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) &= \\ q(\bar{\theta}_1, \theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\bar{\theta}_1, \theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) p(\bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) & \end{aligned} \quad (13)$$

Now multiply both sides by $p(\theta_2, \mathbf{z} | \mathbf{y})$ to get the full augmented posterior $p(\theta_1, \theta_2, \mathbf{z} | \mathbf{y})$

$$\begin{aligned} q(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_2, \mathbf{z} | \mathbf{y}) &= \\ q(\bar{\theta}_1, \theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\bar{\theta}_1, \theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) p(\bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_2, \mathbf{z} | \mathbf{y}) & \quad \text{or} \end{aligned} \quad (14)$$

$$\begin{aligned} q(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\theta_1^0, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_1^0, \theta_2, \mathbf{z} | \mathbf{y}) &= \\ q(\bar{\theta}_1, \theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\bar{\theta}_1, \theta_1^0 | \mathbf{y}, \theta_2, \mathbf{z}) p(\bar{\theta}_1, \theta_2, \mathbf{z} | \mathbf{y}) & \end{aligned}$$

Then integrate over $\boldsymbol{\psi} = [\theta_1^0 \quad \theta_2^0 \quad \mathbf{z}^0]'$, i.e. all parameters and the augmented data (since we're now integrating over all possible values of θ_1^0 we're switching notation from θ_1^0 to θ_1):

$$\begin{aligned} \int q(\theta_1, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\theta_1, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_1, \theta_2, \mathbf{z} | \mathbf{y}) d\boldsymbol{\psi} &= \\ \int q(\bar{\theta}_1, \theta_1 | \mathbf{y}, \theta_2, \mathbf{z}) \alpha(\bar{\theta}_1, \theta_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\bar{\theta}_1 | \mathbf{y}) p(\theta_2, \mathbf{z} | \mathbf{y}, \bar{\theta}_1) d\boldsymbol{\psi} & \end{aligned} \quad (15)$$

since

$$p(\bar{\theta}_1, \theta_2, \mathbf{z} | \mathbf{y}) = p(\bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_2, \mathbf{z} | \mathbf{y})$$

It follows immediately that

$$\begin{aligned} p(\bar{\theta}_1 | \mathbf{y}) &= \frac{\int \alpha(\theta_1, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) q(\theta_1, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_1, \theta_2, \mathbf{z} | \mathbf{y}) d\boldsymbol{\psi}}{\int \alpha(\bar{\theta}_1, \theta_1 | \mathbf{y}, \theta_2, \mathbf{z}) q(\bar{\theta}_1, \theta_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_2, \mathbf{z} | \mathbf{y}, \bar{\theta}_1) d\boldsymbol{\psi}} = \\ &= \frac{E_{p(\theta_1, \theta_2, \mathbf{z} | \mathbf{y})} \left\{ \alpha(\theta_1, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) q(\theta_1, \bar{\theta}_1 | \mathbf{y}, \theta_2, \mathbf{z}) \right\}}{E_{q(\bar{\theta}_1, \theta_1 | \mathbf{y}, \theta_2, \mathbf{z}) p(\theta_2, \mathbf{z} | \mathbf{y}, \bar{\theta}_1)} \left\{ \alpha(\bar{\theta}_1, \theta_1 | \mathbf{y}, \theta_2, \mathbf{z}) \right\}} \end{aligned} \quad (16)$$

To estimate the numerator by MC take the draws $\{\boldsymbol{\theta}_1^g, \boldsymbol{\theta}_2^g, \mathbf{z}^g\}_{g=1}^M$ from the full (original) run and average the quantity $\alpha(\boldsymbol{\theta}_1^g, \bar{\boldsymbol{\theta}}_1 | \mathbf{y}, \boldsymbol{\theta}_2^g, \mathbf{z}^g) q(\bar{\boldsymbol{\theta}}_1 | \boldsymbol{\theta}_1^g, \mathbf{y}, \boldsymbol{\theta}_2^g, \mathbf{z}^g)$. The expectation in the denominator is conditional on $\bar{\boldsymbol{\theta}}_1$ and the original GS does (obviously) not generate draws conditional on $\bar{\boldsymbol{\theta}}_1$. So *continue* the MCMC simulation for an additional J iterations with the two full conditional densities $p(\boldsymbol{\theta}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1, \mathbf{z})$ and $p(\mathbf{z} | \mathbf{y}, \bar{\boldsymbol{\theta}}_1, \boldsymbol{\theta}_2)$. At each iteration, given $\boldsymbol{\theta}_2^j, \mathbf{z}^j$ we generate a variate

$$\boldsymbol{\theta}_1^j \sim q(\bar{\boldsymbol{\theta}}_1, \boldsymbol{\theta}_1 | \mathbf{y}, \boldsymbol{\theta}_2^j, \mathbf{z}^j) \quad (17)$$

This leads to the triple $\boldsymbol{\theta}_1^j, \boldsymbol{\theta}_2^j, \mathbf{z}^j$ that is a draw from $q(\bar{\boldsymbol{\theta}}_1, \boldsymbol{\theta}_1 | \mathbf{y}, \boldsymbol{\theta}_2, \mathbf{z}) p(\boldsymbol{\theta}_2, \mathbf{z} | \mathbf{y}, \bar{\boldsymbol{\theta}}_1)$. The marginal ordinate can now be estimated as

$$\hat{p}(\bar{\boldsymbol{\theta}}_1 | \mathbf{y}) = \frac{M^{-1} \sum_{g=1}^M \alpha(\boldsymbol{\theta}_1^g, \bar{\boldsymbol{\theta}}_1 | \mathbf{y}, \boldsymbol{\theta}_2^g, \mathbf{z}^g) q(\boldsymbol{\theta}_1^g, \bar{\boldsymbol{\theta}}_1 | \mathbf{y}, \boldsymbol{\theta}_2^g, \mathbf{z}^g)}{J^{-1} \sum_{j=1}^J \alpha(\bar{\boldsymbol{\theta}}_1, \boldsymbol{\theta}_1^j | \mathbf{y}, \boldsymbol{\theta}_2^j, \mathbf{z}^j)} \quad (18)$$

Next, the values \mathbf{z}^j from the preceding reduced run, which are marginally from $p(\mathbf{z} | \mathbf{y}, \bar{\boldsymbol{\theta}}_1)$ are used to form the average

$$\hat{p}(\bar{\boldsymbol{\theta}}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1) = J^{-1} \sum_{j=1}^J p(\bar{\boldsymbol{\theta}}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1, \mathbf{z}^j) \quad (19)$$

which is a simulation-consistent estimate of $p(\bar{\boldsymbol{\theta}}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1)$. Thus, at the conclusion of the reduced run both ordinates are available, and we get

$$\log \hat{p}(\mathbf{y}) = \log p(\mathbf{y} | \bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2) + \log p(\bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2) - (\log \hat{p}(\bar{\boldsymbol{\theta}}_1 | \mathbf{y}) + \log \hat{p}(\bar{\boldsymbol{\theta}}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1)) \quad (20)$$

Three remarks: (i) A single reduced run, augmented to sample $\boldsymbol{\theta}_1$ from the proposal density, delivers the variates that are used in the calculation of the ordinates $p(\bar{\boldsymbol{\theta}}_1 | \mathbf{y})$ and $p(\bar{\boldsymbol{\theta}}_2 | \mathbf{y}, \bar{\boldsymbol{\theta}}_1)$. (ii) If one places the recalcitrant ordinate first in the decomposition of the posterior ordinate, the reduced run does not involve any *MH* steps. (iii) This same approach can also be applied when the full conditional density of \mathbf{z} is sampled by a sequence of *MH* steps.

Multiple Parameter blocks

We will abstract from latent data in this section, as it can be handled as outlined in the previous section. Suppose we have B blocks, with $\boldsymbol{\theta}_k \in S_k \subseteq \mathfrak{R}^{d_k}$. Write the posterior ordinate at $\bar{\boldsymbol{\theta}}$ as

$$p(\bar{\boldsymbol{\theta}}_1, \bar{\boldsymbol{\theta}}_2, \dots, \bar{\boldsymbol{\theta}}_B | \mathbf{y}) = \prod_{i=1}^B p(\bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\theta}}_1, \dots, \bar{\boldsymbol{\theta}}_{i-1})$$

and consider the estimation of the reduced ordinate $p(\bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\theta}}_1, \dots, \bar{\boldsymbol{\theta}}_{i-1})$. Let $\boldsymbol{\psi}_{i-1}$ denote the blocks up to i and $\boldsymbol{\psi}_{i+1}$ the blocks beyond i . (So we have $\boldsymbol{\theta}_{-i} = (\boldsymbol{\psi}'_{-i} \quad \boldsymbol{\psi}'_{+i})'$).

Suppose that the full conditional density $p(\boldsymbol{\theta}_i | \mathbf{y}, \boldsymbol{\theta}_{-i}) \propto p(\boldsymbol{\theta}) p(\mathbf{y} | \boldsymbol{\theta})$, $i=1 \dots B$ is sampled by the *MH* algorithm with proposal density $q(\boldsymbol{\theta}_i^a, \boldsymbol{\theta}_i^b | \mathbf{y}, \boldsymbol{\Psi}_{i-1}, \boldsymbol{\Psi}_{i+1})$ and probability of move

$$\alpha(\boldsymbol{\theta}_i^a, \boldsymbol{\theta}_i^b | \mathbf{y}, \boldsymbol{\Psi}_{i-1}, \boldsymbol{\Psi}_{i+1}) = \min \left\{ 1, \frac{p(\mathbf{y} | \boldsymbol{\theta}_i^b, \boldsymbol{\Psi}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\theta}_i^b, \boldsymbol{\theta}_{-i}) q(\boldsymbol{\theta}_i^b, \boldsymbol{\theta}_i^a | \mathbf{y}, \boldsymbol{\Psi}_{i-1}, \boldsymbol{\Psi}_{i+1})}{p(\mathbf{y} | \boldsymbol{\theta}_i^a, \boldsymbol{\Psi}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\theta}_i^a, \boldsymbol{\theta}_{-i}) q(\boldsymbol{\theta}_i^a, \boldsymbol{\theta}_i^b | \mathbf{y}, \boldsymbol{\Psi}_{i-1}, \boldsymbol{\Psi}_{i+1})} \right\} \quad (21)$$

By completely analogous arguments from above we get

$$\begin{aligned} & q(\boldsymbol{\theta}_i^0, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \alpha(\boldsymbol{\theta}_i^0, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\theta}_i^0 | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) = \\ & q(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i^0 | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \alpha(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i^0 | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \end{aligned} \quad (22)$$

Where the use of $\bar{\boldsymbol{\Psi}}_{i-1}$ indicates that the elements in $\boldsymbol{\Psi}_{i-1}$ have been set to their posterior mean. Now multiply both sides by $p(\boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1})$ and integrate over $(\boldsymbol{\theta}_i, \boldsymbol{\Psi}_{i+1})$ to obtain

$$\begin{aligned} p(\bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\theta}}_1, \dots, \bar{\boldsymbol{\theta}}_{i-1}) &= \frac{\int \alpha(\boldsymbol{\theta}_i, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) q(\boldsymbol{\theta}_i, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\theta}_i, \boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}) d(\boldsymbol{\theta}_i, \boldsymbol{\Psi}_{i+1})}{\int \alpha(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) q(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\theta}}_i, \bar{\boldsymbol{\Psi}}_{i-1}) d(\boldsymbol{\theta}_i, \boldsymbol{\Psi}_{i+1})} = \\ & \frac{E_{p(\boldsymbol{\theta}_i, \boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1})} \left\{ \alpha(\boldsymbol{\theta}_i, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) q(\boldsymbol{\theta}_i, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \right\}}{E_{q(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\theta}}_i, \bar{\boldsymbol{\Psi}}_{i-1})} \left\{ \alpha(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \right\}} \end{aligned} \quad (23)$$

These two integrals can be estimated as before from the output of reduced MCMC runs as follows:

Step 1 (Numerator):

Set $\boldsymbol{\Psi}_{i-1} = \bar{\boldsymbol{\Psi}}_{i-1}$ and sample the reduced set of full conditional distributions $p(\boldsymbol{\theta}_k | \mathbf{y}, \boldsymbol{\theta}_{-k}, \bar{\boldsymbol{\Psi}}_{i-1})$, $k=i, \dots, B$.

Let the generated draws be $\{\boldsymbol{\theta}_i^g, \dots, \boldsymbol{\theta}_B^g\}$, $g=1 \dots M$. We can now approximate the numerator in (23) by

$$\begin{aligned} & E_{p(\boldsymbol{\theta}_i, \boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1})} \left(\alpha(\boldsymbol{\theta}_i, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) q(\boldsymbol{\theta}_i, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \right) = \\ & M^{-1} \sum_{i=1}^M \alpha(\boldsymbol{\theta}_i^g, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}^g) q(\boldsymbol{\theta}_i^g, \bar{\boldsymbol{\theta}}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}^g) \end{aligned} \quad (24)$$

Step 2 (Denominator):

In contrast to the numerator, the expectation in the denominator is conditioned on $\bar{\boldsymbol{\theta}}_i$. Thus, we need to go into another reduced run, drawing $p(\boldsymbol{\theta}_k | \mathbf{y}, \boldsymbol{\theta}_{-k}, \bar{\boldsymbol{\Psi}}_{i-1}, \bar{\boldsymbol{\theta}}_i)$, $k=i+1, \dots, B$. This yields

$\{\boldsymbol{\theta}_{i+1}^j, \dots, \boldsymbol{\theta}_B^j\}$, $j=1 \dots J$. At each step also draw $\boldsymbol{\theta}_i^j \sim q(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}^j)$. Together, the draws of $\{\boldsymbol{\theta}_i^j, \boldsymbol{\theta}_{i+1}^j, \dots, \boldsymbol{\theta}_B^j\}$ are draws from $q(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\theta}}_i, \bar{\boldsymbol{\Psi}}_{i-1})$. We can now approximate the denominator as

$$\hat{E}_{q(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) p(\boldsymbol{\Psi}_{i+1} | \mathbf{y}, \bar{\boldsymbol{\theta}}_i, \bar{\boldsymbol{\Psi}}_{i-1})} \left(\alpha(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}) \right) = J^{-1} \sum_{j=1}^J \alpha(\bar{\boldsymbol{\theta}}_i, \boldsymbol{\theta}_i^j | \mathbf{y}, \bar{\boldsymbol{\Psi}}_{i-1}, \boldsymbol{\Psi}_{i+1}^j) \quad (25)$$

where the average in the denominator may include zeros if there are values of θ_i^j that lie outside the support of the posterior S . Thus, the reduced ordinate in question can be estimated as

$$\hat{p}(\bar{\theta}_i | \mathbf{y}, \bar{\theta}_1, \dots, \bar{\theta}_{i-1}) = \frac{M^{-1} \sum_{i=1}^M \alpha(\theta_i^s, \bar{\theta}_i | \mathbf{y}, \bar{\psi}_{i-1}, \psi_{i+1}^s) q(\theta_i^s, \bar{\theta}_i | \mathbf{y}, \bar{\psi}_{i-1}, \psi_{i+1}^s)}{J^{-1} \sum_{j=1}^J \alpha(\bar{\theta}_i, \theta_i^j | \mathbf{y}, \bar{\psi}_{i-1}, \psi_{i+1}^j)} \quad (26)$$

Note: The draws of $\{\theta_{i+1}^j, \dots, \theta_B^j\}$, $j = 1 \dots J$ from Step 2 are also used to estimate the *numerator* of the next reduced posterior ordinate $p(\bar{\theta}_{i+1} | \mathbf{y}, \bar{\theta}_1, \dots, \bar{\theta}_i)$.

References:

- Chib, S. 2001. Markov Chain Monte Carlo methods: Computation and inference. In *Handbook of Econometrics*, J. J. Heckman and E. Leamer (eds.). Elsevier Science, p. 3569 – 3649
- Chib, S. and E. Greenberg. 1995. Understanding the Metropolis-Hastings algorithm. *The American Statistician* **49**: 327-335.
- Chib, S., E. Greenberg and R. Winkelmann. 1998. "Posterior simulation and Bayes factors in panel count data models." *Journal of Econometrics* 86, 33-54.
- Chib, Siddhartha and Bradley P. Carlin. 1999. "On MCMC sampling in hierarchical longitudinal models." *Statistics and Computing* 9, 17-26.
- Chib, S. and I. Jeliazkov. 2001. Marginal likelihood from the Metropolis-Hastings output. *Journal of the American Statistical Association* **95**: 270-281.
- Davis, A. and K. Moeltner. 2010. "Valuing the Prevention of an Infestation: The Threat of the New Zealand Mud Snail in Northern Nevada. ." *Agricultural and Resource Economics Review*.
- Kass, R. E. and A. E. Raftery. 1995. "Bayes Factors." *Journal of the American Statistical Association* 90, 773-795.