

# Normal linear regression model via Gibbs Sampling

## Gibbs Sampler Diagnostics

(K Ch. 4, KPT Ch. 11)

R scripts mod6s3a - mod6s3g

As mentioned previously conjugate priors may be overly restrictive in many Bayesian applications. Here we use a popular combination of independent priors for the regression model – normal for  $\boldsymbol{\beta}$  and inverse-gamma for  $\sigma^2$ . This also implies that we parameterize the error variance directly in terms of  $\sigma^2$  and not in term of its inverse, i.e. the precision. This will “feel” more natural to a classically trained analyst.

The enhanced flexibility in prior modeling comes at the price of abandoning analytical results for the posterior distribution. Instead, we will use posterior simulation via Gibbs Sampling to obtain draws from the joint and marginal posteriors.

The structural model is the same as for the conjugate model. Here, nothing is gained by introducing  $\mathbf{b}$  and SSE into the LHF, so we start simply with:

$$p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) = (2\pi)^{-n/2} (\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \left( (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right)\right) \quad (1)$$

Note the change in notation from  $h$  to  $\sigma^2$ . The priors are given as follows:

$$\begin{aligned} p(\boldsymbol{\beta}, \sigma^2) &= p(\boldsymbol{\beta}) p(\sigma^2) \text{ where} \\ \boldsymbol{\beta} &\sim n(\boldsymbol{\mu}_0, \mathbf{V}_0), \quad \sigma^2 \sim ig(v_0, \tau_0) \\ p(\boldsymbol{\beta}) &= (2\pi)^{-k/2} |\mathbf{V}_0|^{-1/2} \exp\left(-\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) \\ p(\sigma^2) &= \frac{\tau_0^{v_0}}{\Gamma(v_0)} (\sigma^2)^{-(v_0+1)} \exp\left(-\frac{\tau_0}{\sigma^2}\right), \quad \text{with } E(\sigma^2) = \frac{\tau_0}{v_0 - 1}, \quad V(\sigma^2) = \frac{\tau_0^2}{(v_0 - 1)^2 (v_0 - 2)} \end{aligned} \quad (2)$$

Note that  $\sigma^2$  does not enter the prior density of  $\boldsymbol{\beta}$ . Also, amongst the many possible parameterizations of the inverse-gamma (*ig*) density, we choose the form given in Gelman et al. (2004), where  $v_0$  is the shape parameter and  $\tau_0$  is the scale parameter. The *ig* density becomes more diffuse (flatter) with smaller values for shape and scale. However, for the density to have a defined mean, we need  $v_0 > 1$ , and for a well-defined variance we need  $v_0 > 2$ .

Combining the priors with the likelihood, and dropping all terms that are multiplicatively unrelated to our parameters of interest yields the posterior kernel

$$\begin{aligned} p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &\propto \\ (\sigma^2)^{-\frac{n-2v_0-2}{2}} \exp\left(-\frac{1}{2\sigma^2} (2\tau_0)\right) &\exp\left(-\frac{1}{2} \left( \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0) \right)\right). \end{aligned} \quad (3)$$

We first aim to find the posterior density for  $\boldsymbol{\beta}$ , conditional on  $\sigma^2$  (i.e. treating  $\sigma^2$  as a constant). Thus, we will first focus on the components of the posterior kernel that cannot be multiplicatively separated from  $\boldsymbol{\beta}$ . This leaves

$$p(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right)\right). \quad (4)$$

Note the conditionality on  $\sigma^2$  on both sides of (4). We'll now use a similar transformation trick as for the previous model. Let

$$\mathbf{V}_1 = \left(\mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{X}\right)^{-1} \quad \text{and} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left(\mathbf{V}_0^{-1}\boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{y}\right) = \mathbf{V}_1 \left(\mathbf{V}_0^{-1}\boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{y}\right) \quad (5)$$

Then:

$$\begin{aligned} & \frac{1}{\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\mu}_0) = \\ & \frac{1}{\sigma^2} \mathbf{y}'\mathbf{y} + \boldsymbol{\beta}' \frac{1}{\sigma^2} \mathbf{X}'\mathbf{X}\boldsymbol{\beta} - \frac{1}{\sigma^2} \mathbf{y}'\mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}' \frac{1}{\sigma^2} \mathbf{X}'\mathbf{y} + \boldsymbol{\beta}' \mathbf{V}_0^{-1} \boldsymbol{\beta} - \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 = \\ & \boldsymbol{\beta}' \left(\mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{X}\right) \boldsymbol{\beta} - \boldsymbol{\beta}' \left(\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{y}\right) - \left(\boldsymbol{\mu}_0' \mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{y}'\mathbf{X}\right) \boldsymbol{\beta} + \frac{1}{\sigma^2} \mathbf{y}'\mathbf{y} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 = \\ & \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\beta} + \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 + \frac{1}{\sigma^2} \mathbf{y}'\mathbf{y} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 = \\ & (\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1) - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 + \frac{1}{\sigma^2} \mathbf{y}'\mathbf{y} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 \end{aligned} \quad (6)$$

Thus, we can re-write the conditional posterior kernel in (4) as

$$p(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1)\right) \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2} \mathbf{y}'\mathbf{y} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1\right)\right). \quad (7)$$

None of the terms in the second exponent include  $\boldsymbol{\beta}$ , so this conditional further simplifies to

$$p(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1)\right). \quad (8)$$

Thus, we have again the kernel of a multivariate normal density, and can state that

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_1, \mathbf{V}_1) \quad \text{with} \quad \mathbf{V}_1 = \left(\mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{X}\right)^{-1} \quad \text{and} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left(\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{y}\right) \quad (9)$$

To derive the conditional posterior density for  $\sigma^2$ , we return to our original form for the joint posterior given in (3). Ignoring terms that are not related to  $\sigma^2$ , we have

$$p(\sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X}) \propto (\sigma^2)^{-\frac{n-2}{2}} \exp\left(-\frac{1}{2\sigma^2} \left(2\tau_0 + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)\right). \quad (10)$$

Comparing this expression to the kernel of the *ig* prior in (2), we recognize this as the kernel of another *ig* density. Specifically:

$$\sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X} \sim ig(v_1, \tau_1) \quad \text{with} \quad (11)$$

$$v_1 = \frac{2v_0 + n}{2} \quad \text{and} \quad \tau_1 = \frac{2\tau_0 + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2}$$

Thus, it will be straightforward to draw  $\boldsymbol{\beta}$  conditional on  $\sigma^2$  and vice versa.

## Gibbs Sampler

The Gibbs Sampler (GS) has become the “workhorse” of Bayesian posterior simulation in recent years. The general idea is simple: break the joint posterior into conditional posteriors for which the analytical form of its density is known. Then sample sequentially and repeatedly from these conditionals. After a number of draws the joint sequence of conditional draws will converge to the desired joint posterior densities for all parameter. In addition, each individual sequence can be interpreted as the marginal posterior for a given parameter. The GS is an example of a “Markov Chain Monte Carlo”, or  $MC^2$  procedure.

Formally, assume our parameter vector of interest is  $\boldsymbol{\theta}$ , with posterior kernel  $p(\boldsymbol{\theta} | \mathbf{y}) \propto p(\boldsymbol{\theta})p(\mathbf{y} | \boldsymbol{\theta})$ . Assume further that this kernel has no known analytical form. However, we can split  $\boldsymbol{\theta}$  into two components, say  $\boldsymbol{\theta}_1$  and  $\boldsymbol{\theta}_2$ . For example, in our current application,  $\boldsymbol{\theta}_1 = \boldsymbol{\beta}$  and  $\boldsymbol{\theta}_2 = \sigma^2$ . In other applications, we may want to split  $\boldsymbol{\theta}$  into 3, 4, or even more components. The key notion is that we know the analytical form of the resulting *full conditional posterior distributions*, i.e.

$$p(\boldsymbol{\theta}_1 | \mathbf{y}, \boldsymbol{\theta}_2) \quad \text{and} \quad p(\boldsymbol{\theta}_2 | \mathbf{y}, \boldsymbol{\theta}_1). \quad (3.12)$$

All we need to get the GS started is an initial value for  $\boldsymbol{\theta}_2$ , call it  $\boldsymbol{\theta}_2^0$ . This can be chosen arbitrarily, or one can use OLS results or results from previous analyses. We assume this starting value comes directly from the marginal posterior  $p(\boldsymbol{\theta}_2 | \mathbf{y})$ . Next, we draw  $\boldsymbol{\theta}_1$  conditional on  $\boldsymbol{\theta}_2^0$  from  $p(\boldsymbol{\theta}_1 | \mathbf{y}, \boldsymbol{\theta}_2^0)$ . Call this draw  $\boldsymbol{\theta}_1^1$ . Next, we draw another value of  $\boldsymbol{\theta}_2$  conditional on  $\boldsymbol{\theta}_1^1$  from  $p(\boldsymbol{\theta}_2 | \mathbf{y}, \boldsymbol{\theta}_1^1)$ . Call this draw  $\boldsymbol{\theta}_2^1$ . We repeat this process  $R$  times. In essence, we use the basic rule of conditional probabilities, i.e.

$$p(\boldsymbol{\theta} | \mathbf{y}) = p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 | \mathbf{y}) = p(\boldsymbol{\theta}_1 | \mathbf{y}, \boldsymbol{\theta}_2) p(\boldsymbol{\theta}_2 | \mathbf{y}) = p(\boldsymbol{\theta}_2 | \mathbf{y}, \boldsymbol{\theta}_1) p(\boldsymbol{\theta}_1 | \mathbf{y}) \quad (3.13)$$

over and over again. As a caveat we should note that, naturally, there is no guarantee that our starting value  $\boldsymbol{\theta}_2^0$  really came from the marginal posterior  $p(\boldsymbol{\theta}_2 | \mathbf{y})$ . However, under relatively weak conditions (see Koop Ch. 4 for details) the starting value(s) will not matter and the GS will indeed converge to draws from  $p(\boldsymbol{\theta} | \mathbf{y})$ . To assure that the effect of the starting value has truly “faded away”, we usually discard the first  $r_1$  draws of the sequence, and keep only the remaining  $r_2 = R - r_1$  draws. The discarded draws are often referred to as “burn-ins”.

## Monte Carlo Integration

The derivation of moments from this simulated posterior is accomplished via Monte Carlo Integration. For example, the analytical expressions for the mean (or expectation) and variance of a given element of  $\boldsymbol{\beta}$ , say  $\beta_j$ , are given by

$$E(\beta_j) = \int \beta_j p(\beta_j | \mathbf{y}, \mathbf{X}) d\beta_j \quad V(\beta_j) = \int (\beta_j - E(\beta_j))^2 p(\beta_j | \mathbf{y}, \mathbf{X}) d\beta_j$$

A convenient alternative expression for the variance is

$$V(\beta_j) = E(\beta_j^2) - (E(\beta_j))^2 \quad \text{where} \quad E(\beta_j^2) = \int \beta_j^2 p(\beta_j | \mathbf{y}, \mathbf{X}) d\beta_j$$

Also, the expectation and variance of any other function of  $\beta_j$ , say  $g(\beta_j)$ , take the form of

$$E(g(\beta_j)) = \int g(\beta_j) p(\beta_j | \mathbf{y}, \mathbf{X}) d\beta_j \quad V(g(\beta_j)) = \int (g(\beta_j) - E(g(\beta_j)))^2 p(\beta_j | \mathbf{y}, \mathbf{X}) d\beta_j$$

For any of these expressions, MCI approximates the integral with averaging over draws. Thus

$$E(\beta_j) \approx \frac{1}{R} \sum \beta_{j,r}$$

$$E(\beta_j^2) \approx \frac{1}{R} \sum \beta_{j,r}^2$$

$$E(g(\beta_j)) \approx \frac{1}{R} \sum g(\beta_{j,r})$$

## Implementation with Simulated Data

**R** script `mod6s3a` provides an opportunity to examine *ig* distributions for different settings for the shape and scale parameters. In general, lower settings for  $\nu_0$  and higher settings for  $\tau_0$  produce a flatter (more diffuse) priors. However, the most diffuse and most commonly used prior settings are  $\frac{1}{2}$  for both parameters. While this does not produce identified moments, the resulting pdf is still proper, in the sense of integrating to one.

**R** script `mod6s3b` is the main script for this model. It opens the log file, loads data, specifies priors and other settings for the posterior simulator, and runs the Gibbs Sampler (GS). The script also applies diagnostic tools (discussed later) to the posterior draws and generates plots comparing the prior and posterior distributions for our parameters. As before, the posterior densities are much tighter than the priors in all cases.

## Convergence Plots

Convergence plots are a simple visual tool to examine if the simulator has converged to the posterior distribution for a given parameter. The plot simply shows all draws of a given parameter in chronological order, as generated by the sampler. "Convergence" usually implies that the draws are tightly clustered around a flat line (the posterior mean), and do not wander widely. This is similar to assessing stationarity for time series data.

Convergence plots can be used to assess the sensitivity of the algorithm to starting draws (the chain should converge to the same distribution under different starting draws), the speed of convergence (and thus the efficiency of the sampler), and the sufficiency of the chosen number of discarded draws (burn-ins).

Script `mod6s3c` provides a few examples using simulated data. There are three parts: Full data, data truncated to 100 observations, and data truncated to 10 observations. For each case, we first use two sets of starting draws for  $\beta$  and  $\sigma^2$ . The first set comes from the OLS output and is thus "right on target", i.e. close to the posterior mean (and the true parameter values underlying our simulated data). The second set is deliberately located quite far from the true values.

For all cases we use a slightly modified GS. As opposed to the previous version, this algorithm preserves the burn-in draws and sends them back to the main script along with the "keepers" (i.e. retained draws from before).

You can see that with the full data set of 10,000 observations, the choice of starting draws virtually doesn't matter – the chain essentially converges after one or two iterations. Furthermore, parameter draws fluctuate very tightly around the posterior mean. A reduction in sample size implicitly assigns more weight to our diffuse priors. As a result, the chain exhibits more "noise", i.e. wider fluctuations around the mean. This will translate into a larger posterior standard deviation. However, even with just a handful of observations, convergence is virtually immediate even with off-target starting draws. This is an indication that the sampler itself is efficient, i.e. "mixes rapidly".

Posterior noise (i.e. the posterior standard deviation) is also driven by the information content of the data, regardless of sample size. Highly collinear data implies poor information content and will generate wider fluctuations around the posterior mean. This is illustrated in script `mod6s3d`. If you compare the plots from this script to the ones from before for any parameter and sample size, you will notice the increase in variability around the posterior mean for the collinear data.

### Application: Female Earnings

*R* script `mod6s3e` implements the normal linear regression model with independent priors using Mroz's (1987) labor data.

## Gibbs Sampler Diagnostics

This section describes in detail the output component generated by the "diagnostics" section of `mod6s3b`, and all subsequent scripts using a Gibbs Sampler.

### Numerical standard error (nse)

The *nse* captures simulation noise for a value of interest generated via posterior simulation. Usually this value of interest is a measure of central tendency, such as the posterior mean. Consider the mean

$\bar{\theta} = \frac{1}{M} \sum_{m=1}^M \theta_m$  of a sequence of  $m$  draws of (some generic) parameter  $\theta$ . Assume these draws were

generated by a Gibbs Sampler (GS) or some other Markov-Chain Monte Carlo (MC<sup>2</sup>) algorithm. If these draws were perfectly independently and identically distributed (i.i.d.) with sample variance  $s^2$ , we could quickly derive the *nse* for  $\bar{\theta}$  using the basic formula for the standard error of a sample mean, i.e.

$$nse(\bar{\theta}) = \sqrt{V(\bar{\theta})} = \sqrt{\frac{1}{M^2} Ms^2} = \frac{s}{\sqrt{M}} \quad (1)$$

However, as with all MC<sup>2</sup> procedures, we have to ex ante assume that these sequential draws will have a considerable degree of *correlation*. This means we have to consider all covariance terms between all draws of  $\theta$ . After some straightforward analytical simplifications (shown in detail in KPT, p. 145), we end up with a more general expression for the *nse* that allows for correlation across all draws:

$$nse(\bar{\theta}) = \sqrt{V(\bar{\theta})} = \sqrt{\frac{s^2}{M} \left( 1 + 2 \sum_{j=1}^{M-1} \left( 1 - \frac{j}{M} \right) \rho_j \right)}, \quad (2)$$

where  $\rho_j = \frac{s_j}{s^2}$  is the lag-correlation between some draw  $\theta_i$  and a draw that was obtained  $j$  iterations prior to  $\theta_i$ , i.e.  $\theta_{i-j}$ , with  $s_j$  denoting the associated sample covariance. It can be easily seen from (2) that under perfect independence ( $\rho_j = 0, \forall j$ ) we arrive again at the basic formula given in (1). In most MC<sup>2</sup> applications, lag-correlations will be positive but declining in magnitude. Thus, the second term under the square root in (2) will exceed 1, and the *nse* under correlation will be larger than the *nse* under independence. Note that in either case the *nse* can be made arbitrarily small by increasing  $M$ , the number of draws. However, this can be very costly in terms of computation time. Thus, we are always interested in devising a posterior sampler that is “efficient” in the sense of generating draws with low lag-correlation. There are many “tricks” for accomplishing this – we’ll touch upon a few in this course.

In summary, the first purpose of the *nse* is to provide a measure of “simulation error” or “simulation noise” surrounding a posterior construct of interest, usually the posterior mean of a given parameter. Thus, the *nse* has a similar function to the standard error (s.e.) in Classical Analysis. However, its intuition is very different – it simply captures simulation noise, i.e. the penalty for having to approximate the joint posterior via simulations (since its analytical form is unknown). The Classical s.e. conveys the notion of *sampling error* – i.e. the variability of the statistical construct of interest under (hypothetical) re-sampling. You can also think of this as the penalty from working with a small sample relative to a large population.

### Inefficiency factors (IEFs)

The second purpose of the *nse* is to provide a summary measure of the “efficiency” of the posterior simulator. An efficient simulator will have low correlation across draws, and thus will be able to “tell the same story with fewer draws” – saving valuable computing time. As discussed in Chib (2001,

section 3.2) the ratio of the squared  $nse$  under correlation over the squared  $nse$  under independence can be interpreted as “inefficiency factor” (IEF), also known as “autocorrelation time” for a given parameter, i.e.

$$IEF(\theta) = \frac{nse^2(\bar{\theta})}{nse^2(\bar{\theta}; \rho_j = 0, \forall j)} = 1 + 2 \sum_{j=1}^{M-1} \left(1 - \frac{j}{M}\right) \rho_j \quad (3)$$

A well-designed posterior simulator will generate sequences of parameter draws with low IEFs, the ideal being an IEF close to 1. Sometimes the inverse of the IEF is used to measure posterior efficiency. This quantity is called “numerical efficiency” (See Geweke, 1992). A third quantity to assess efficiency is the “i.i.d - equivalent number of iterations”, labeled  $M^*$  in the following, i.e. the number of i.i.d. draws that contain the same amount of information about  $\theta$  as the observed number of draws under correlation. It is easily derived via

$$\frac{s^2}{M^*} = \frac{s^2}{M} \left(1 + 2 \sum_{j=1}^{M-1} \left(1 - \frac{j}{M}\right) \rho_j\right) \rightarrow M^* = \frac{M}{\left(1 + 2 \sum_{j=1}^{M-1} \left(1 - \frac{j}{M}\right) \rho_j\right)} = \frac{M}{IEF} \quad (4)$$

It follows that under perfect efficiency, we have  $M^* = M$ , but usually we observe  $M^* < M$ .

It should be noted that a very high IEF (very low  $M^*$ ) can be indicative of identification problems in your model. How high is “very high”? From personal experience, I would say that IEFs in the 1-5 range indicate good efficiency, in the 6-20 range they’re still “tolerable”, but anything above 20 deserves closer inspection. Certainly, IEFs of 100 and higher are almost a sure-bet indication of an identification or specification problem in the underlying structural model.

### Geweke’s (1992) Convergence Diagnostics(CD)

Another important sampler diagnostic is Geweke’s (1992)  $CD$  score. It is based on the simple intuition that if the entire sequence of retained  $\theta$ 's (focusing on a single parameter for simplicity) can truly be interpreted as random draws from the same posterior density  $p(\theta | \mathbf{y})$ , and we divide the sequence of  $R$  draws into three segments, the mean of the first segment of  $r = 1 \cdots R_1$  draws should be “not too different” from the mean of the last segment of  $r = R_2 + 1 \cdots R$  draws. As stated in Koop Ch. 4, setting  $R_1 = 0.1R$  and  $R_2 = 0.6R$  produces adequate results for most applications. Define the two means as  $\bar{\theta}_1$  and  $\bar{\theta}_2$ , respectively. Then, asymptotically, the difference between these means, weighted by their respective numerical standard errors, converges to a standard normal (“z”) variate, i.e.

$$CD = \frac{\bar{\theta}_1 - \bar{\theta}_2}{\sqrt{nse_1^2 + nse_2^2}} \overset{a}{\sim} n(0,1). \quad (4.5)$$

Thus, a CD value that clearly exceeds 1.96 for a specific parameter  $\theta$  would raise a flag – it indicates that the sequence of posterior draws may not have converged to  $p(\theta | \mathbf{y})$ . In practice, if a few CD values in the 2-2.5 range for a model with many parameters would hardly raise concerns. However, if your posterior simulator generates CD values of 3 or higher, an increase in the number of burn-in draws may be warranted. Similarly to IEFs, grossly inflated CD values may also indicate identification and / or mis-specification problems in the underlying model.

## Autocorrelation (AC) plots

There are two additional noteworthy diagnostics tools: autocorrelation plots (“AC plots”) and re-running the posterior sampler with different starting values. An AC plot provides a simple visual inspection of the lag-correlation terms ( $\rho_j$  in (2)) for a given parameter. An example for AC plots is provided in the `R` script `mod6s3f`. A "well-behaved" AC plot will exhibit small correlation effects that randomly fluctuate around zero. A plot with high (usually positive) correlations that taper off only very slowly with increasing lag would be indicative of inefficiencies in the posterior simulator.

## Blocking

As discussed in previous Sessions of this course, a Gibbs Sampler operates by splitting the full set of parameters into different groups or "*blocks*", which are then drawn sequentially and repeatedly, conditional on *all other* blocks.

The main consideration in designing these blocks for a standard Gibbs Sampler is that the conditional posterior density for each block is known, else we wouldn't be able to take any draws from it. However, this requirement can be relaxed by employing other posterior simulation techniques, such as the Metropolis Hastings (MH) algorithm, which does not require full knowledge of the conditional posterior density for a given parameter or block of parameters. Thus, the "optimal blocking" of the full set of parameters becomes a more general question.

As discussed *inter alia* in Chib (2001, section 7.1) it is generally recommended that parameters be drawn in as few blocks as possible, and that parameters that tend to be highly correlated be collected in the same block.

Identifying what constitutes a full-fledged block in a given posterior algorithm can be tricky. This is because blocks can be combined by the method of composition, i.e. by exploiting the fact that *partially* conditional posterior distributions may be known (or can be approximated at low computational cost) for some parameters. For example, consider an initial blocking of the full parameter vector into three groups,  $\theta_1, \theta_2$ , and  $\theta_3$ . A "naïve" posterior sampler will then operate as follows:

1. Draw  $\theta_1$  from  $p(\theta_1 | \theta_2, \theta_3, \mathbf{y})$ , where  $\mathbf{y}$  represents the available data.
2. Draw  $\theta_2$  from  $p(\theta_2 | \theta_1, \theta_3, \mathbf{y})$
3. Draw  $\theta_3$  from  $p(\theta_3 | \theta_1, \theta_2, \mathbf{y})$
4. Repeat.

Now suppose the partially conditional density  $p(\theta_1 | \theta_3, \mathbf{y})$  is known or can be approximated at low computational cost. A (likely) more efficient posterior sampler would then collect  $\theta_1$  and  $\theta_2$  in a single block and operate as follows:

1. Draw  $\theta_1, \theta_2$  from  $p(\theta_1, \theta_2 | \theta_3, \mathbf{y}) = p(\theta_2 | \theta_1, \theta_3, \mathbf{y}) * p(\theta_1 | \theta_3, \mathbf{y})$  as follows:
  - a. Draw  $\theta_1$  from  $p(\theta_1 | \theta_3, \mathbf{y})$
  - b. Draw  $\theta_2$  from  $p(\theta_2 | \theta_1, \theta_3, \mathbf{y})$
2. Draw  $\theta_3$  from  $p(\theta_3 | \theta_1, \theta_2, \mathbf{y})$

Thus, even though step 1 involves 2 sub-steps, it is considered a single block. The key notion in identifying the number of blocks is that for a set of parameters to constitute a self-standing block, it needs to be drawn conditional on *all other* blocks, and *vice versa*, i.e. all other blocks also need to be conditioned on the first block.

### Blocking in the normal linear regression model

For the normal regression model we have considered so far the parameter blocking was quite obvious and, as it turns out, efficient: We grouped the constant term and all slope parameters into a single block ( $\boldsymbol{\beta}$ ), which left the regression variance  $\sigma^2$  as the only other (single-parameter) block. Our GS proceeded as follows:

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

Matlab script `mod6s3g` proposes a different approach based on 4 blocks. Specifically, we split  $\boldsymbol{\beta}$  into three parts of equal length, labeled  $\boldsymbol{\beta}_1$ ,  $\boldsymbol{\beta}_2$ , and  $\boldsymbol{\beta}_3$ . We then proceed as follows:

1. Draw  $\boldsymbol{\beta}_1$  from  $p(\boldsymbol{\beta}_1 | \boldsymbol{\beta}_2, \boldsymbol{\beta}_3, \sigma^2, \mathbf{y}, \mathbf{X})$
2. Draw  $\boldsymbol{\beta}_2$  from  $p(\boldsymbol{\beta}_2 | \boldsymbol{\beta}_1, \boldsymbol{\beta}_3, \sigma^2, \mathbf{y}, \mathbf{X})$
3. Draw  $\boldsymbol{\beta}_3$  from  $p(\boldsymbol{\beta}_3 | \boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \sigma^2, \mathbf{y}, \mathbf{X})$
4. Draw  $\sigma^2$  from  $p(\sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X})$

In practice, draws of  $\boldsymbol{\beta}_j, j = 1 \dots 3$  can be obtained as follows:

We know that for the basic linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{6}$$

we obtain conditional draws of  $\boldsymbol{\beta}$  via

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_1, \mathbf{V}_1) \quad \text{with} \quad \mathbf{V}_1 = \left( \mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{X} \right)^{-1} \quad \text{and} \quad \boldsymbol{\mu}_1 = \mathbf{V}_1 \left( \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}'\mathbf{y} \right) \tag{7}$$

Now partition  $\mathbf{X}$  and  $\boldsymbol{\beta}$  into three parts corresponding to our new blocking, i.e.

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{X}_3\boldsymbol{\beta}_3 + \boldsymbol{\varepsilon} \tag{8}$$

Consider the "transformed" dependent variable  $\tilde{\mathbf{y}} = \mathbf{y} - \mathbf{X}_2\boldsymbol{\beta}_2 - \mathbf{X}_3\boldsymbol{\beta}_3$  and the resulting modified regression model

$$\tilde{\mathbf{y}} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\varepsilon} \tag{9}$$

We can thus derive the conditional posterior for  $\boldsymbol{\beta}_1$  as

$$\beta_1 | \beta_2, \beta_3, \sigma^2, y, \mathbf{X} \sim n(\boldsymbol{\mu}_1, \mathbf{V}_1) \quad \text{with} \quad \mathbf{V}_1 = \left( \mathbf{V}_{01}^{-1} + \frac{1}{\sigma^2} \mathbf{X}'_1 \mathbf{X}_1 \right)^{-1} \quad \text{and} \quad (10)$$

$$\boldsymbol{\mu}_1 = \mathbf{V}_1 \left( \mathbf{V}_{01}^{-1} \boldsymbol{\mu}_{01} + \frac{1}{\sigma^2} \mathbf{X}'_1 \tilde{y} \right) = \mathbf{V}_1 \left( \mathbf{V}_{01}^{-1} \boldsymbol{\mu}_{01} + \frac{1}{\sigma^2} \mathbf{X}' (y - \mathbf{X}_2 \beta_2 - \mathbf{X}_3 \beta_3) \right)$$

where  $\boldsymbol{\mu}_{01}$  and  $\mathbf{V}_{01}$  are the prior mean and variance for  $\beta_1$ . Draws of  $\beta_2$  and  $\beta_3$  can be obtained in analogous fashion.

The posterior output clearly shows efficiency losses compared to the original version, as judged by IEF and M\* scores, which is evident from comparing autocorrelation plots for the original and the "excessively-blocked" version for selected parameters.

Also note the loss in speed – the inefficient sampler takes about twice as long to take the same number of draws.

We can also compare the relative performance of the two samplers based on AC and convergence plots. The plotted chain of draws for "we" is especially illustrative: the chain wanders widely, with a clear autocorrelation pattern. This is confirmed by the AC plot for "we", which shows high correlations for neighboring draws that are reluctant to taper off.

The main drawback of such a highly correlated chain is that it takes much longer to "visit" the entire posterior distribution with appropriate frequencies. This may result in misleading posterior inference, based on overly tight or otherwise "incomplete" distributions. (In fact, the posterior standard deviations flowing from the inefficient sampler are actually slightly smaller than those generated by the efficient sampler).

## References

- Chib, Siddhartha. 2001. "Markov Chain Monte Carlo Methods: Computation and Inference," in J. J. Heckman and E. Leamer (eds), *Handbook of Econometrics*: Elsevier.
- Geweke, J. 1992. "Evaluating the Accuracy of Sampling-based Approaches to the Calculation of Posterior Moments," in J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith (eds), *Bayesian Statistics 4* Oxford, UK: Oxford University Press.