

Markov Chain Monte Carlo Methods

John Geweke
University of Iowa, USA

2005 Institute on Computational Economics
University of Chicago - Argonne National Laboratories

July 22, 2005

The problem

$p(\boldsymbol{\theta}, \boldsymbol{\omega} | I)$ Distribution of interest

$\boldsymbol{\theta} \in \Theta, \boldsymbol{\omega} \in \Omega; I = \text{“Information”}$

$$p(\boldsymbol{\theta}, \boldsymbol{\omega} | I) = p(\boldsymbol{\theta} | I) \cdot p(\boldsymbol{\omega} | \boldsymbol{\theta}, I)$$

$p(\boldsymbol{\omega} | \boldsymbol{\theta}, I)$ Tractable for simulation

$p(\boldsymbol{\theta} | I)$ Not so easy

Leading example:

$$I = \{\text{Model specification}\} \cup \{\text{Data}\}$$

Origins of the problem in econometric inference

Complete model

$$\left. \begin{array}{l} p(\mathbf{y} | \boldsymbol{\theta}, A) \\ p(\boldsymbol{\theta} | A) \end{array} \right\} \implies p(\boldsymbol{\theta} | \mathbf{y}^o, A)$$

Vector of interest $\boldsymbol{\omega}$

$$p(\boldsymbol{\omega} | \mathbf{y}^o, \boldsymbol{\theta}, A)$$

Simulation problem

$$\begin{array}{l} \boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} | \mathbf{y}^o, A) \\ \boldsymbol{\omega}^{(m)} \sim p(\boldsymbol{\omega} | \mathbf{y}^o, \boldsymbol{\theta}, A) \end{array}$$

Some background: Random number generation

Mother of all random number generators (Geweke 1996):

$$u^{(m)} \stackrel{iid}{\sim} \text{uniform}(0, 1)$$

Inverse c.d.f. transformation to simulate $\theta : P(\theta \leq c) = F(c)$:

$$\theta^{(m)} = F^{-1}(u^{(m)})$$

$$\theta^{(m)} \leq c \iff u^{(m)} = F(\theta^{(m)}) \leq F(c)$$

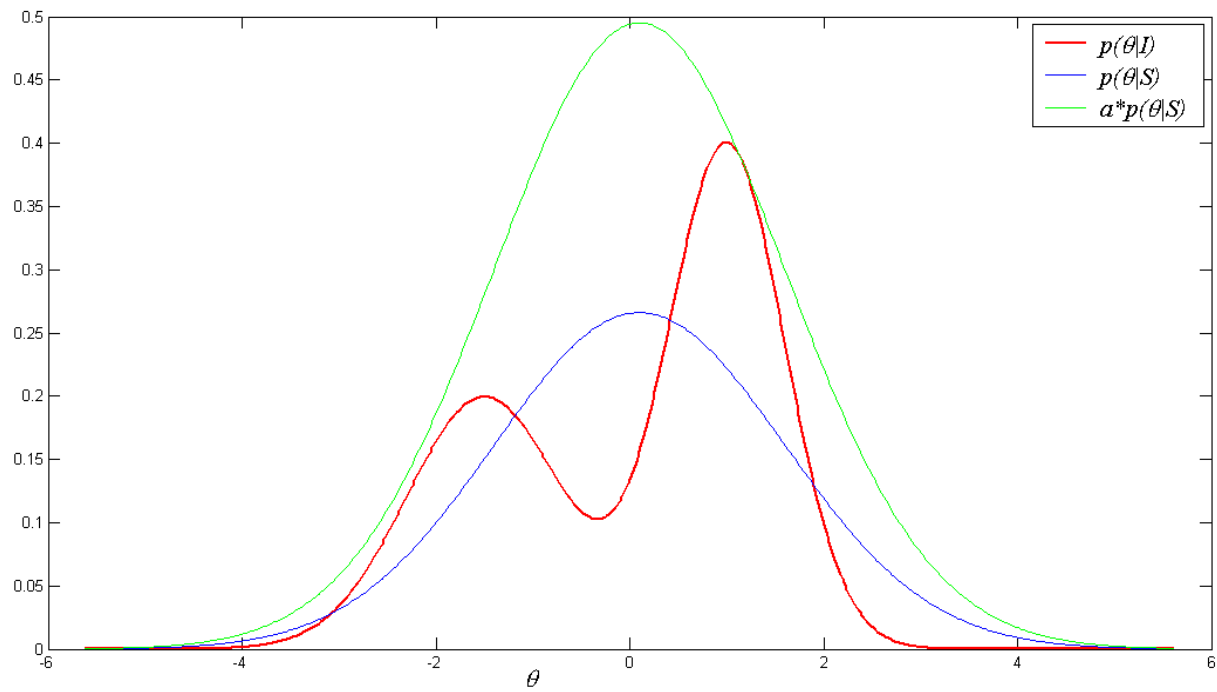
$$\implies P(\theta^{(m)} \leq c) = P[u^{(m)} \leq F(c)] = F(c)$$

Examples:

$$\text{Sensible: } \theta^{(m)} = -\log(u^{(m)}) / \alpha \implies F(\theta) = 1 - \exp(-\alpha\theta)$$

$$\text{Questionable: } \theta^{(m)} = \Phi^{-1}(u^{(m)}) \implies F(\theta) = \Phi(\theta)$$

Some background: Acceptance sampling (Motivation)



Some background: Acceptance sampling (Algorithm)

Target density: $p(\boldsymbol{\theta} | I)$

Source density: $p(\boldsymbol{\theta} | S)$

Corresponding kernels: $k(\boldsymbol{\theta} | I) \propto p(\boldsymbol{\theta} | I)$, $k(\boldsymbol{\theta} | S) \propto p(\boldsymbol{\theta} | S)$

$$r = \sup [k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S)] < \infty$$

1. $u \stackrel{iid}{\sim} \text{uniform}(0, 1)$
2. $\boldsymbol{\theta}^* \sim p(\boldsymbol{\theta} | S)$
3. If $u > k(\boldsymbol{\theta}^* | I) / [r \cdot k(\boldsymbol{\theta}^* | S)]$ then return to step 1
4. Set $\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^*$

Essence of Markov Chain Monte Carlo (MCMC)

$$\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(m-1)}, C)$$

Convergence in distribution:

$$\boldsymbol{\theta}^{(m)} \xrightarrow{d} p(\boldsymbol{\theta} \mid I)$$

Ergodicity:

$$M^{-1} \sum_{m=1}^M g(\boldsymbol{\theta}^{(m)}) \xrightarrow{a.s.} E[g(\boldsymbol{\theta}) \mid I]$$

If $\boldsymbol{\omega}^{(m)} \sim p(\boldsymbol{\omega} \mid \boldsymbol{\theta}^{(m)})$, then

$$\{\boldsymbol{\theta}^{(m)}, \boldsymbol{\omega}^{(m)}\} \xrightarrow{d} p(\boldsymbol{\theta}, \boldsymbol{\omega} \mid I) \text{ and } M^{-1} \sum_{m=1}^M h(\boldsymbol{\omega}^{(m)}) \xrightarrow{a.s.} E[h(\boldsymbol{\omega}) \mid I]$$

Gibbs sampler: Bivariate distribution

$$\boldsymbol{\theta}' = (\theta_{(1)}, \theta_{(2)})$$

$$\theta_{(1)}^{(0)} \sim ??$$

$$\theta_{(2)}^{(0)} \sim p\left(\theta_{(2)} \mid \theta_{(1)}^{(0)}, I\right)$$

For $m = 1, 2, \dots$

$$\theta_{(1)}^{(m)} \sim p\left(\theta_{(1)} \mid \theta_{(2)}^{(m-1)}, I\right)$$

$$\theta_{(2)}^{(m)} \sim p\left(\theta_{(2)} \mid \theta_{(1)}^{(m)}, I\right)$$

Gibbs sampler: Bivariate distribution example

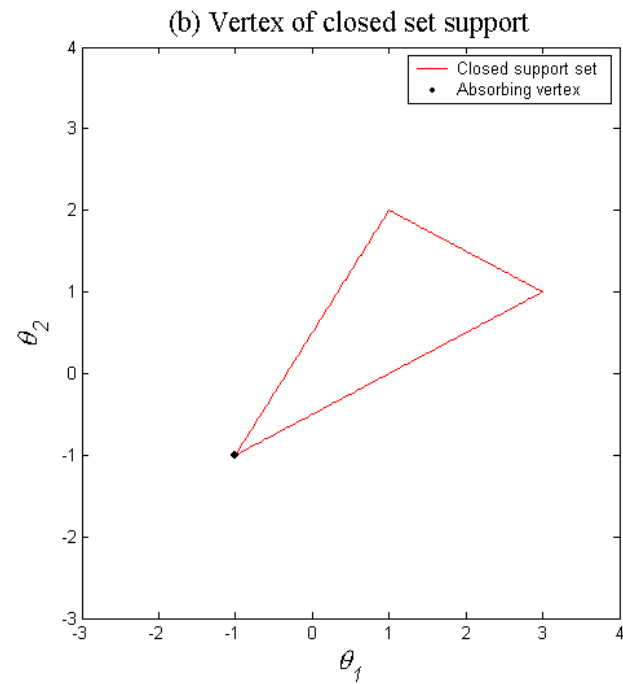
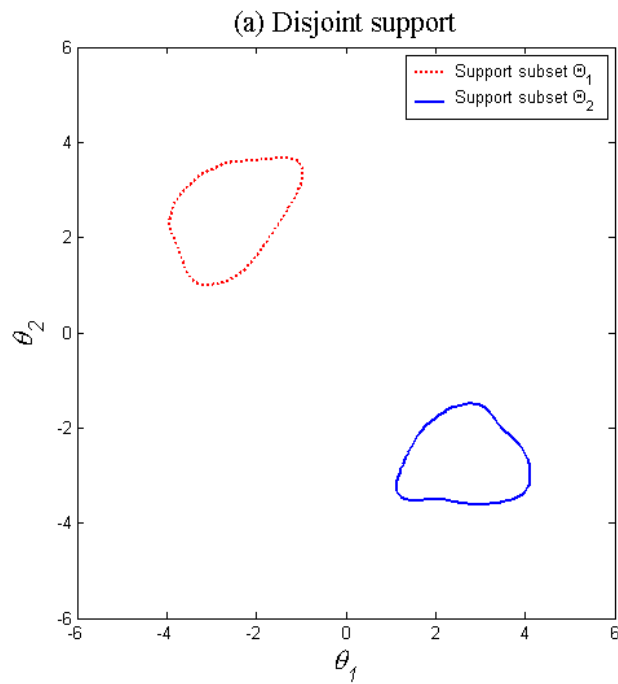
$$\boldsymbol{\theta}_{2 \times 1} \mid I \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

subject to

$$\theta_1 \in (\underline{\theta}_1, \bar{\theta}_1) \text{ and } \theta_2 \in (\underline{\theta}_2, \bar{\theta}_2)$$

$$\theta_1 \mid (\theta_2, I) \sim N \left[\mu_1 + \frac{\sigma_{12}(\theta_2 - \mu_2)}{\sigma_{22}}, \sigma_{11} - \frac{\sigma_{12}^2}{\sigma_{22}} \right] \text{ subject to } \theta_1 \in (\underline{\theta}_1, \bar{\theta}_1)$$

Gibbs sampler: Potential problems with a bivariate distribution



Gibbs sampler: General case

Blocking of θ :

$$\theta' = (\theta'_{(1)}, \dots, \theta'_{(B)}); \theta_{(b)} \mid [\theta_{(a)} \ (a \neq b), I] \text{ is tractable}$$

Given $\theta^{(0)}$,

$$\theta_{(b)}^{(1)} \sim p \left(\theta_{(b)} \mid \left[\theta_{(a)}^{(1)} \ (a < b), \theta_{(a)}^{(0)} \ (a > b), I \right] \right) \quad (b = 1, \dots, B)$$

defines

$$p \left(\theta^{(1)} \mid \theta^{(0)}, G \right) = \prod_{b=1}^B p \left(\theta_{(b)} \mid \left[\theta_{(a)}^{(1)} \ (a < b), \theta_{(a)}^{(0)} \ (a > b), I \right] \right)$$

Gibbs sampler: Convergence result

Suppose that for every point $\boldsymbol{\theta} \in \Theta$ and every measurable $A \subseteq \Theta$,

$$\int_A p(\boldsymbol{\theta} | I) d\nu(\boldsymbol{\theta}) > 0 \Rightarrow \int_A p(\boldsymbol{\theta}^* | \boldsymbol{\theta}, G) d\nu(\boldsymbol{\theta}^*) > 0.$$

Then the transition kernel G of the Gibbs sampler is ergodic: if $E[g(\boldsymbol{\theta}), I]$ exists, then for all $\boldsymbol{\theta}^{(0)} \in \Theta$,

$$M^{-1} \sum_{m=1}^M g(\boldsymbol{\theta}^{(m)}) \xrightarrow{a.s.} \int_{\Theta} g(\boldsymbol{\theta}) p(\boldsymbol{\theta} | I) d\boldsymbol{\theta}.$$

For alternative conditions see Geweke (2005) and for weaker conditions still see Tierney (1994).

Metropolis-Hastings algorithm: Structure

Begin with an arbitrary transition density function

$$q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(m-1)}, H)$$

Set $\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^*$ with probability

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(m-1)}, H) \in (0, 1).$$

Otherwise set $\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^{(m-1)}$.

Density kernel of accepted candidates $\boldsymbol{\theta}^*$:

$$q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(m-1)}, H) \cdot \alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(m-1)}, H)$$

Reversibility

Suppose a transition T is reversible with respect to $p(\boldsymbol{\theta} | I)$:

$$p(\boldsymbol{\theta}^{(m-1)} | I) p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, T) = p(\boldsymbol{\theta}^{(m)} | I) p(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, T).$$

Then $p(\boldsymbol{\theta} | I)$ is an invariant density for T :

$$\begin{aligned} & \int_{\Theta} p(\boldsymbol{\theta}^{(m-1)} | I) p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, T) d\nu(\boldsymbol{\theta}^{(m-1)}) \\ &= \int_{\Theta} p(\boldsymbol{\theta}^{(m)} | I) p(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, T) d\nu(\boldsymbol{\theta}^{(m-1)}) \\ &= p(\boldsymbol{\theta}^{(m)} | I) \int_{\Theta} p(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, T) d\nu(\boldsymbol{\theta}^{(m-1)}) = p(\boldsymbol{\theta}^{(m)} | I) \end{aligned}$$

Note that reversibility holds trivially when $\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^{(m-1)}$.

Metropolis-Hastings algorithm: Making it work

For $\boldsymbol{\theta}^{(m)} \neq \boldsymbol{\theta}^{(m-1)}$ we require

$$p(\boldsymbol{\theta}^{(m-1)} | I) p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) = p(\boldsymbol{\theta}^{(m)} | I) p(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, H).$$

where $p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) \propto q(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) \cdot \alpha(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H)$.

Hence

$$\begin{aligned} & p(\boldsymbol{\theta}^{(m-1)} | I) q(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) \alpha(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) \\ = & p(\boldsymbol{\theta}^{(m)} | I) q(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, H) \alpha(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, H) \\ \implies & \frac{\alpha(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H)}{\alpha(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, H)} = \frac{p(\boldsymbol{\theta}^{(m)} | I) q(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m)}, H)}{p(\boldsymbol{\theta}^{(m-1)} | I) q(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H)} \end{aligned}$$

... and the condition

$$\frac{\alpha(\boldsymbol{\theta}^{(m)} \mid \boldsymbol{\theta}^{(m-1)}, H)}{\alpha(\boldsymbol{\theta}^{(m-1)} \mid \boldsymbol{\theta}^{(m)}, H)} = \frac{p(\boldsymbol{\theta}^{(m)} \mid I) q(\boldsymbol{\theta}^{(m-1)} \mid \boldsymbol{\theta}^{(m)}, H)}{p(\boldsymbol{\theta}^{(m-1)} \mid I) q(\boldsymbol{\theta}^{(m)} \mid \boldsymbol{\theta}^{(m-1)}, H)}$$

is satisfied if

$$\alpha(\boldsymbol{\theta}^{(m)} \mid \boldsymbol{\theta}^{(m-1)}, H) = \min \left\{ \frac{p(\boldsymbol{\theta}^{(m)} \mid I) q(\boldsymbol{\theta}^{(m-1)} \mid \boldsymbol{\theta}^{(m)}, H)}{p(\boldsymbol{\theta}^{(m-1)} \mid I) q(\boldsymbol{\theta}^{(m)} \mid \boldsymbol{\theta}^{(m-1)}, H)}, 1 \right\}.$$

Metropolis-Hastings algorithm: Convergence result

Suppose that for every point $\boldsymbol{\theta} \in \Theta$ and every measurable $A \subseteq \Theta$,

$$\int_A p(\boldsymbol{\theta} | I) d\nu(\boldsymbol{\theta}) > 0 \Rightarrow \int_A q(\boldsymbol{\theta}^* | \boldsymbol{\theta}, H) d\nu(\boldsymbol{\theta}^*) > 0.$$

Then the transition kernel H of the Metropolis-Hastings algorithm is ergodic: if $E[g(\boldsymbol{\theta}), I]$ exists, then for all $\boldsymbol{\theta}^{(0)} \in \Theta$,

$$M^{-1} \sum_{m=1}^M g(\boldsymbol{\theta}^{(m)}) \xrightarrow{a.s.} \int_{\Theta} g(\boldsymbol{\theta}) p(\boldsymbol{\theta} | I) d\boldsymbol{\theta}.$$

For alternative conditions see Geweke (2005) and for weaker conditions see Tierney (1994).

Metropolis independence chain

$$q(\boldsymbol{\theta}^* | \boldsymbol{\theta}, H) = q(\boldsymbol{\theta}^* | H)$$

If $q(\boldsymbol{\theta}^* | H) > 0$ for all $\boldsymbol{\theta}^* \in \Theta$, then this chain must be ergodic. Note

$$\begin{aligned} \alpha(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H) &= \alpha(\boldsymbol{\theta}^* | H) = \min \left\{ \frac{p(\boldsymbol{\theta}^* | I) q(\boldsymbol{\theta}^{(m-1)} | H)}{p(\boldsymbol{\theta}^{(m-1)} | I) q(\boldsymbol{\theta}^* | H)}, 1 \right\} \\ &= \min \left\{ \frac{p(\boldsymbol{\theta}^* | I) / q(\boldsymbol{\theta}^* | H)}{p(\boldsymbol{\theta}^{(m-1)} | I) / q(\boldsymbol{\theta}^{(m-1)} | H)}, 1 \right\} \end{aligned}$$

Metropolis random walk chain

$$q(\boldsymbol{\theta}^* | \boldsymbol{\theta}, H) = q(\boldsymbol{\theta}^* - \boldsymbol{\theta}, H) = q(\boldsymbol{\theta} - \boldsymbol{\theta}^*, H)$$

Leading example: $\boldsymbol{\theta}^* | H \sim N(\boldsymbol{\theta}, \Sigma)$

This is an instance of a Metropolis chain because $q(\boldsymbol{\theta}^* | \boldsymbol{\theta}, H) = q(\boldsymbol{\theta} | \boldsymbol{\theta}^*, H)$ implying

$$\begin{aligned} \alpha(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H) &= \min \left\{ \frac{p(\boldsymbol{\theta}^* | I) q(\boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^*, H)}{p(\boldsymbol{\theta}^{(m-1)} | I) q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H)}, 1 \right\} \\ &= \min \left\{ p(\boldsymbol{\theta}^* | I) / p(\boldsymbol{\theta}^{(m-1)} | I), 1 \right\} \end{aligned}$$

Metropolis within Gibbs

Suppose that we block a problem for the Gibbs sampler

$$\boldsymbol{\theta}' = (\boldsymbol{\theta}'_{(1)}, \dots, \boldsymbol{\theta}'_{(B)})$$

but $\boldsymbol{\theta}_{(b)} \mid [\boldsymbol{\theta}_{(a)} \ (a \neq b), I]$ is intractable for some b .

Let $q(\boldsymbol{\theta}_{(b)}^* \mid [\boldsymbol{\theta}_{(a)} \ (a \neq b)], H_b)$ be an arbitrary transition density.

At step m of the MCMC algorithm, draw

$$\boldsymbol{\theta}_{(b)}^* \sim q\left(\boldsymbol{\theta}_{(b)}^* \mid \boldsymbol{\theta}_{(1)}^m, \dots, \boldsymbol{\theta}_{(b-1)}^{(m)}, \boldsymbol{\theta}_{(b)}^{(m-1)}, \boldsymbol{\theta}_{(b+1)}^{(m-1)}, \dots, \boldsymbol{\theta}_{(B)}^{(m-1)}, H_b\right)$$

... and accept with probability

$$\begin{aligned}
 & \alpha \left(\boldsymbol{\theta}_{(b)}^* \mid \boldsymbol{\theta}_{(1)}^m, \dots, \boldsymbol{\theta}_{(b-1)}^{(m)}, \boldsymbol{\theta}_{(b+1)}^{(m-1)}, \dots, \boldsymbol{\theta}_{(B)}^{(-1)}, H_b \right) \\
 = & \min \left\{ \frac{p \left(\boldsymbol{\theta}_{(1)}^m, \dots, \boldsymbol{\theta}_{(b-1)}^{(m)}, \boldsymbol{\theta}_{(b)}^*, \boldsymbol{\theta}_{(b+1)}^{(m-1)}, \dots, \boldsymbol{\theta}_{(B)}^{(-1)}, I \right)}{p \left(\boldsymbol{\theta}_{(1)}^m, \dots, \boldsymbol{\theta}_{(b-1)}^{(m)}, \boldsymbol{\theta}_{(b)}^{(m-1)}, \boldsymbol{\theta}_{(b+1)}^{(m-1)}, \dots, \boldsymbol{\theta}_{(B)}^{(-1)}, I \right)} \right. \\
 & \left. \cdot \frac{q \left(\boldsymbol{\theta}_{(b)}^{(m-1)} \mid \boldsymbol{\theta}_{(1)}^m, \dots, \boldsymbol{\theta}_{(b-1)}^{(m)}, \boldsymbol{\theta}_{(b+1)}^{(m-1)}, \dots, \boldsymbol{\theta}_{(B)}^{(m-1)}, H_b \right)}{q \left(\boldsymbol{\theta}_{(b)}^* \mid \boldsymbol{\theta}_{(1)}^m, \dots, \boldsymbol{\theta}_{(b-1)}^{(m)}, \boldsymbol{\theta}_{(b)}^{(m-1)}, \boldsymbol{\theta}_{(b+1)}^{(m-1)}, \dots, \boldsymbol{\theta}_{(B)}^{(m-1)}, H_b \right)}, 1 \right\}
 \end{aligned}$$

Chib and Greenberg (1995); Geweke (2005, Section 4.6.2)

Application to inference

Model

$$p(\boldsymbol{\theta} | A), \quad p(\mathbf{y} | \boldsymbol{\theta}, A)$$

$\boldsymbol{\theta}$ Vector of unobservables

\mathbf{y} Vector of observables

For observed values $\mathbf{y} = \mathbf{y}^o$,

$$p(\mathbf{y}^o | \boldsymbol{\theta}, A) \quad \text{Likelihood function}$$

$$p(\boldsymbol{\theta} | A) p(\mathbf{y}^o | \boldsymbol{\theta}, A) \quad \text{Kernel of posterior density function}$$

$$p(\boldsymbol{\theta} | \mathbf{y}^o, A) \propto p(\boldsymbol{\theta} | A) p(\mathbf{y}^o | \boldsymbol{\theta}, A)$$

MCMC and the posterior distribution

Some perspective:

$$\begin{array}{ll} \boldsymbol{\theta} \sim p(\boldsymbol{\theta} | A) & \text{Easy} \\ \mathbf{y} \sim p(\mathbf{y} | \boldsymbol{\theta}, A) & \text{Easy} \end{array}$$

We seek $p(\boldsymbol{\theta} | \boldsymbol{\theta}^{(m-1)}, \mathbf{y}^o, C)$ with the property

$$M^{-1} \sum_{m=1}^M g(\boldsymbol{\theta}^{(m)}) \xrightarrow{a.s.} E[g(\boldsymbol{\theta}) | \mathbf{y}^o, A]$$

which implies

$$\boldsymbol{\theta}^{(m)} \xrightarrow{d} p(\boldsymbol{\theta} | \mathbf{y}^o, A)$$

Getting it right: the problem

Errors in derivations and/or coding still produce a Markov chain

$$\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(m-1)}, C)$$

that will produce reasonable answers.

Most students and many researchers “check” their code by

$$\begin{aligned} \boldsymbol{\theta} &= \boldsymbol{\theta}^0 \\ \mathbf{y} &\sim p(\mathbf{y} \mid \boldsymbol{\theta}^0, A) \end{aligned}$$

and observing whether or not $\{\boldsymbol{\theta}^{(m)} \ (m = 1, \dots, M)\}$ is concentrated near $\boldsymbol{\theta}^0$.

This does not work.

Getting it right: the solution

Note there exists a joint distribution for θ and \mathbf{y} ,

$$p(\theta, \mathbf{y} | A) = p(\theta | A) p(\mathbf{y} | \theta, A).$$

There is an easy way to simulate from this distribution:

$$\begin{aligned}\theta^{(m)} &\sim p(\theta | A) \\ \mathbf{y}^{(m)} &\sim p(\mathbf{y} | \theta^{(m)}, A)\end{aligned}$$

There is a hard way to simulate from this distribution:

$$\begin{aligned}\tilde{\theta}^{(m)} &\sim p(\theta | \tilde{\theta}^{(m-1)}, \tilde{\mathbf{y}}^{(m-1)}, C) \\ \tilde{\mathbf{y}}^{(m)} &\sim p(\mathbf{y} | \tilde{\theta}^{(m)}, A)\end{aligned}$$

... if all derivations and coding are correct, then a standard HAC (Newey-West) test will show that

$$M^{-1} \sum_{m=1}^M g(\boldsymbol{\theta}^{(m)}) - M^{-1} \sum_{m=1}^M g(\tilde{\boldsymbol{\theta}}^{(m)})$$

has mean zero.

... if there are errors in derivations (even conceptual) or coding, then the test will fail.

Geweke (2004; 2005, Section 8.1)

References

Chib, S., and E. Greenberg (1995), “Understanding the Metropolis-Hastings algorithm”, *The American Statistician* 49: 327-335.

Geweke, J. (1996), “Monte Carlo simulation and numerical integration”, in: H. Amman, D. Kendrick and J. Rust, eds., *Handbook of Computational Economics* (North-Holland, Amsterdam) 731-800.

Geweke, J. (2004), “Getting It Right: Joint Distribution Tests of Posterior Simulators,” *Journal of the American Statistical Association* 99: 799-804.

Geweke, J. (2005), *Contemporary Bayesian Econometrics and Statistics* (Wiley, New York; forthcoming in October).

Tierney, L. (1994), “Markov chains for exploring posterior distributions”, *Annals of Statistics* 22: 1701-1762.