

Monte Carlo Methods
for Econometric Inference I
Institute on Computational Economics

July 19, 2006

John Geweke, University of Iowa

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

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

Direct Sampling

Some familiar tools

Strong law of large numbers:

If $\omega^{(m)} \stackrel{iid}{\sim} p(\omega)$ and $\mathbf{E}[h(\omega)] = \int h(\omega) p(\omega) d\nu(\omega)$ exists, then

$$\bar{h}^{(M)} = M^{-1} \sum_{m=1}^M h(\omega^{(m)}) \xrightarrow{a.s.} \mathbf{E}[h(\omega)] = \bar{h};$$

that is, $P \left\{ \lim_{M \rightarrow \infty} \bar{h}^{(M)} = \bar{h} \right\} = 1$, a stronger condition than $\bar{h}^{(M)} \xrightarrow{p} \bar{h}$.

Lindeberg-Levy central limit theorem:

If in addition $\sigma_h^2 = \text{var} [h(\omega)]$ exists, then

$$M^{1/2} \left[\bar{h}^{(M)} - \bar{h} \right] \xrightarrow{d} N \left(\mathbf{0}, \sigma_h^2 \right).$$

Generic notation for posterior simulation

$$\begin{aligned}\boldsymbol{\theta} &\sim p(\boldsymbol{\theta} | I); \boldsymbol{\theta} \in \Theta \\ \boldsymbol{\omega} &\sim p(\boldsymbol{\omega} | \boldsymbol{\theta}, I); \boldsymbol{\omega} \in \Omega\end{aligned}$$

“ I ” denotes the distribution of interest.

Direct sampling is possible if there are practical algorithms for

$$\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | I) \quad \text{and} \quad \boldsymbol{\omega}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\omega} | \boldsymbol{\theta}^{(m)}, I).$$

Example:

(1) Mother of all random number generators (Geweke 1996):

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

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

$$\begin{aligned} \theta^{(m)} &= F^{-1}(u^{(m)}) \\ \theta^{(m)} \leq c &\iff u^{(m)} = F(\theta^{(m)}) \leq F(c) \\ \implies P(\theta^{(m)} \leq c) &= P[u^{(m)} \leq F(c)] = F(c) \end{aligned}$$

General case:

$$\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} \mid I) \quad \text{and} \quad \boldsymbol{\omega}^{(m)} \mid \boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\omega} \mid \boldsymbol{\theta}^{(m)}, I)$$

where $\boldsymbol{\theta}^{(m)} \in \Theta$, $\boldsymbol{\omega}^{(m)} \in \Omega$, and all draws are independent.

For given h and M , let $\bar{h}^{(M)} = M^{-1} \sum_{m=1}^M h(\boldsymbol{\omega}^{(m)})$.

Theorem 4.1(a): $\boldsymbol{\omega}^{(m)} \stackrel{i.i.d.}{\sim} p(\boldsymbol{\omega} \mid I)$

Theorem 4.1(b): If $\bar{h} = \mathbb{E}[h(\boldsymbol{\omega}) \mid I]$ exists then $\bar{h}^{(M)} \xrightarrow{a.s.} \bar{h}$ (strong law of large numbers).

Theorem 4.1(c): If $\bar{h} = E[h(\omega) | I]$ and $\text{var}[h(\omega) | I] = \sigma^2$ exist,

then $M^{1/2} (\bar{h}^{(M)} - \bar{h}) \xrightarrow{d} N(0, \sigma^2)$ and

$$\hat{\sigma}^{2(M)} = M^{-1} \sum_{m=1}^M [h(\omega^{(m)}) - \bar{h}^{(M)}]^2 \xrightarrow{a.s.} \sigma^2.$$

(Lindeberg-Levy central limit theorem and strong law of large numbers)

Theorem 4.1(d): If for given $p \in (0, 1)$, there is a unique h_p such that the statements

$$P [h(\omega) \leq h_p \mid I] \geq p \text{ and } P [h(\omega) \geq h_p \mid I] \geq 1 - p$$

are both true, then

$$\hat{h}_p^{(M)} \xrightarrow{a.s.} h_p$$

where $\hat{h}_p^{(M)}$ is any real number such that

$$M^{-1} \sum_{m=1}^M I_{(-\infty, \hat{h}_p^{(M)}} [h(\omega^{(m)})] \geq p; \quad M^{-1} \sum_{m=1}^M I_{[\hat{h}_p^{(M)}, \infty)} [h(\omega^{(m)})] \geq 1 - p$$

(Follows from in Rao (1965), result 6f.2(i).)

Theorem 4.1(e): If for given $p \in (0, 1)$, there is a unique h_p such that the statements

$$P[h(\omega) \leq h_p \mid I] \geq p \text{ and } P[h(\omega) \geq h_p \mid I] \geq 1 - p$$

are both true, and moreover $p, p[h(\omega) = h_p \mid I] > 0$, then

$$M^{1/2} \left[\hat{h}_p^{(M)} - h_p \right] \xrightarrow{d} N \left\{ 0, p(1-p) / p[h(\omega) = h_p \mid I]^2 \right\}.$$

(Follows from in Rao (1965), result 6f.2(i).)

Approximation of Bayes actions by direct sampling – General case

The setting:

- $\{\boldsymbol{\theta}^{(m)}, \boldsymbol{\omega}^{(m)}\}$ i.i.d. with $\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} | I)$, $\boldsymbol{\omega}^{(m)} | (\boldsymbol{\theta}^{(m)}, I) \sim p(\boldsymbol{\omega} | \boldsymbol{\theta}^{(m)}, I)$.
- Loss function $L(\mathbf{a}, \boldsymbol{\omega}) \geq 0$ defined on $\Omega \times A$, A an open subset of \mathbb{R}^q .
- The risk function

$$R(\mathbf{a}) = \int_{\Omega} \int_{\Theta} L(\mathbf{a}, \boldsymbol{\omega}) p(\boldsymbol{\theta} | I) p(\boldsymbol{\omega} | \boldsymbol{\theta}, I) d\boldsymbol{\theta} d\boldsymbol{\omega}$$

has a strict global minimum at $\hat{\mathbf{a}} \in A \subseteq \mathbb{R}^m$.

Let A_M be the set of all roots of $M^{-1} \sum_{m=1}^M \partial L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a} = \mathbf{0}$.

Theorem 4.2(a): If

(1) $M^{-1} \sum_{m=1}^M L(\mathbf{a}, \boldsymbol{\omega}^{(m)})$ converges uniformly to $R(\mathbf{a})$ for all $\mathbf{a} \in N(\hat{\mathbf{a}})$, almost surely;

(2) $\partial L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a}$ exists and is a continuous function of \mathbf{a} , for all $\boldsymbol{\omega} \in \Omega$ and all $\mathbf{a} \in N(\hat{\mathbf{a}})$;

then for any $\varepsilon > 0$

$$\lim_{M \rightarrow \infty} P \left[\inf_{\mathbf{a} \in A_M} (\mathbf{a} - \hat{\mathbf{a}})' (\mathbf{a} - \hat{\mathbf{a}}) > \varepsilon \mid I \right] = 0.$$

Theorem 4.2 (b): If, in addition,

(3) $\partial^2 L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a} \partial \mathbf{a}'$ exists and is a continuous function of \mathbf{a} , for all $\boldsymbol{\omega} \in \Omega$ and all $\mathbf{a} \in \mathbf{N}(\hat{\mathbf{a}})$;

(4) $\mathbf{B} = \text{var} [\partial L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a} |_{\mathbf{a}=\hat{\mathbf{a}}} | I]$ exists and is finite;

(5) $\mathbf{H} = E [\partial^2 L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a} \partial \mathbf{a}' |_{\mathbf{a}=\hat{\mathbf{a}}} | I]$ exists and is finite and nonsingular;

(6) For any $\varepsilon > 0$, there exists M_ε such that

$$P \left[\sup_{\mathbf{a} \in \mathbf{N}(\hat{\mathbf{a}})} \left| \partial^3 L(\mathbf{a}, \boldsymbol{\omega}) / \partial a_i \partial a_j \partial a_k \right| < M_\varepsilon \mid I \right] \geq 1 - \varepsilon$$

for all $i, j, k = 1, \dots, m$;

then:

$$1. M^{1/2} (\hat{\mathbf{a}}_M - \hat{\mathbf{a}}) \xrightarrow{d} N(\mathbf{0}, \mathbf{H}^{-1} \mathbf{B} \mathbf{H}^{-1}),$$

$$2. M^{-1} \sum_{m=1}^M \partial L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a} |_{\mathbf{a}=\hat{\mathbf{a}}_M} \cdot \partial L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a}' |_{\mathbf{a}=\hat{\mathbf{a}}_M} \xrightarrow{p} \mathbf{B},$$

$$3. M^{-1} \sum_{m=1}^M \partial^2 L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a} \partial \mathbf{a}' |_{\mathbf{a}=\hat{\mathbf{a}}_M} \xrightarrow{p} \mathbf{H}.$$

... This all follows from Amemiya (1985), Chapter 4, in straightforward fashion.

Acceptance and Importance Sampling

Motivation and approach

We have no practical method to simulate $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | I)$.

We can simulate $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | S)$, where $p(\boldsymbol{\theta} | S)$ is “similar” to $p(\boldsymbol{\theta} | I)$.

Can we use $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | S)$ to simulate from $p(\boldsymbol{\theta} | I)$?

Yes, **but** the sense in which $p(\boldsymbol{\theta} | S)$ is similar to $p(\boldsymbol{\theta} | I)$ is critical.

Acceptance Sampling

Density of interest: $p(\boldsymbol{\theta} | I)$

We cannot simulate $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | I)$,

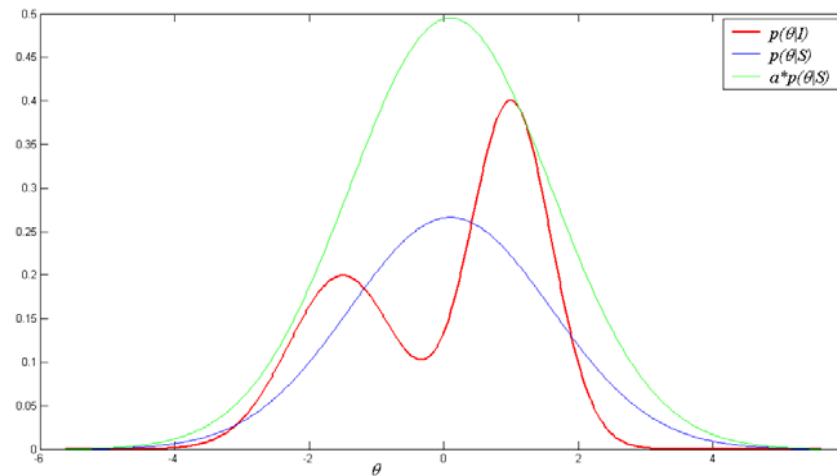
We can evaluate $k(\boldsymbol{\theta} | I) \propto p(\boldsymbol{\theta} | I)$

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

We can simulate $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | S)$

We can evaluate $k(\boldsymbol{\theta} | S) \propto p(\boldsymbol{\theta} | S)$

Acceptance sampling: Key idea



$$\sup_{\theta} p(\theta | I) / p(\theta | S) = p(1.16 | I) / p(1.16 | S) = a = 1.86$$

Draw θ^* from $p(\theta|S)$, and accept the draw with probability $p(\theta^*|I) / [a \cdot p(\theta^*|S)]$

Theorem 4.2.1 Acceptance sampling. Let

$$k(\boldsymbol{\theta} | I) = c_I \cdot p(\boldsymbol{\theta} | I)$$

be a kernel of the density of interest $p(\boldsymbol{\theta} | I)$ and let

$$k(\boldsymbol{\theta} | S) = c_S \cdot p(\boldsymbol{\theta} | S)$$

be a kernel of the source density $p(\boldsymbol{\theta} | S)$. Let $r = \sup_{\boldsymbol{\theta} \in \Theta} k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S) < \infty$. Suppose that $\boldsymbol{\theta}^{(m)}$ is drawn as follows.

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

If the draws in steps 1 and 2 are independent, then $\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} | I)$.

$$k(\boldsymbol{\theta} | I) = c_I \cdot p(\boldsymbol{\theta} | I), \quad k(\boldsymbol{\theta} | S) = c_S \cdot p(\boldsymbol{\theta} | S), \quad \Theta^* = \{\boldsymbol{\theta} : p(\boldsymbol{\theta} | S) > 0\}$$

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

$$P[(3) \rightarrow (4)] = \int_{\Theta^*} \{k(\boldsymbol{\theta} | I) / [rk(\boldsymbol{\theta} | S)]\} p(\boldsymbol{\theta} | S) d\nu(\boldsymbol{\theta}) = c_I / rc_S$$

$$\begin{aligned} P[(3) \rightarrow (4), \boldsymbol{\theta} \in A] &= \int_A \{k(\boldsymbol{\theta} | I) / [rk(\boldsymbol{\theta} | S)]\} p(\boldsymbol{\theta} | S) d\nu(\boldsymbol{\theta}) \\ &= \int_A k(\boldsymbol{\theta} | I) d\nu(\boldsymbol{\theta}) / rc_S \end{aligned}$$

$$\begin{aligned} P[\boldsymbol{\theta} \in A | (3) \rightarrow (4)] &= \int_A k(\boldsymbol{\theta} | I) d\nu(\boldsymbol{\theta}) / c_I = \int_A p(\boldsymbol{\theta} | I) d\nu(\boldsymbol{\theta}) \\ &= P(\boldsymbol{\theta} \in A | I) \end{aligned}$$

Efficiency of acceptance sampling

$$k(\boldsymbol{\theta} | I) = c_I p(\boldsymbol{\theta} | I), \quad k(\boldsymbol{\theta} | S) = c_S p(\boldsymbol{\theta} | S), \quad r = \sup_{\boldsymbol{\theta} \in \Theta} \frac{k(\boldsymbol{\theta} | I)}{k(\boldsymbol{\theta} | S)}, \quad a = \sup_{\boldsymbol{\theta} \in \Theta} \frac{p(\boldsymbol{\theta} | I)}{p(\boldsymbol{\theta} | S)}$$

$$\begin{aligned} P[(3) \rightarrow (4)] &= \int_{\Theta^*} \{k(\boldsymbol{\theta} | I) / [rk(\boldsymbol{\theta} | S)]\} p(\boldsymbol{\theta} | S) d\nu(\boldsymbol{\theta}) = \frac{c_I}{rc_S} \\ &= \frac{c_I}{c_S \cdot \sup_{\boldsymbol{\theta} \in \Theta} k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S)} = \inf_{\boldsymbol{\theta} \in \Theta} \frac{c_I}{c_S \cdot k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S)} \\ &= \inf_{\boldsymbol{\theta} \in \Theta} \frac{p(\boldsymbol{\theta} | S)}{p(\boldsymbol{\theta} | I)} = a^{-1} \end{aligned}$$

$$\frac{\text{Draws from } p(\boldsymbol{\theta} | S)}{\text{Draws from } p(\boldsymbol{\theta} | I)} \rightarrow a$$

Example:

Sampling from a standard normal distribution truncated to the interval (a, b)

Characteristics of the truncation			Source distribution
$a < 0 < b < \infty$	$a \geq -t_1$ and $b \leq t_1$		Uniform(a, b)
	$a < -t_1$ or $b > t_1$		$N(0, 1)$
$0 \leq a < b < \infty$	$\phi(a) / \phi(b) \leq t_2$		Uniform(a, b)
	$\phi(a) / \phi(b) > t_2$	$a \leq t_3$	$ N(0, 1) $
		$a > t_3$	$a + \exp(a^{-1})$
$b = \infty$	$a \leq t_4$		$N(0, 1)$
	$a > t_4$		$a + \exp(a^{-1})$
$t_1 = 0.375, t_2 = 2.18, t_3 = 0.725, \text{ and } t_4 = 0.45$			
2 to 6 times as fast as inverse c.d.f. method			

Importance Sampling

Density of interest: $p(\boldsymbol{\theta} | I)$

We cannot simulate $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | I)$,

We can evaluate $k(\boldsymbol{\theta} | I) \propto p(\boldsymbol{\theta} | I)$

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

We can simulate $\boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | S)$

We can evaluate $k(\boldsymbol{\theta} | S) \propto p(\boldsymbol{\theta} | S)$

Weighting function: $w(\boldsymbol{\theta}) = k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S)$

Theorem 4.2 (a,b) Approximation of moments by importance sampling.

Suppose $E[h(\omega) | I] = \bar{h}$ exists, the support of $p(\theta | S)$ includes Θ , and $\omega^{(m)} \sim p(\omega | \theta^{(m)}, I)$. Then

$$\bar{h}^{(M)} = \frac{\sum_{m=1}^M w(\theta^{(m)}) h(\omega^{(m)})}{\sum_{m=1}^M w(\theta^{(m)})} \xrightarrow{a.s.} \bar{h}.$$

Proof.

$$\begin{aligned} w(\theta^{(m)}) \text{ i.i.d., } E[w(\theta) | S] &= \int_{\Theta} \frac{k(\theta | I)}{k(\theta | S)} p(\theta | S) d\nu(\theta) = \frac{c_I}{c_S} = \bar{w} \\ \implies \bar{h}^{(M)} &= M^{-1} \sum_{m=1}^M w(\theta^{(m)}) \xrightarrow{a.s.} \bar{w}. \end{aligned}$$

Proof (concluded) So far we have

$$\bar{h}^{(M)} = \frac{\sum_{m=1}^M w(\boldsymbol{\theta}^{(m)}) h(\boldsymbol{\omega}^{(m)})}{\sum_{m=1}^M w(\boldsymbol{\theta}^{(m)})}, \quad \bar{w}^{(M)} = M^{-1} \sum_{m=1}^M w(\boldsymbol{\theta}^{(m)}) \xrightarrow{a.s.} \bar{w}$$

Then $w(\boldsymbol{\theta}) = k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S)$, $\{w(\boldsymbol{\theta}^{(m)}), h(\boldsymbol{\omega}^{(m)})\}$ i.i.d. \implies

$$\begin{aligned} \mathbb{E}[w(\boldsymbol{\theta}) h(\boldsymbol{\omega}) | S] &= \int_{\Theta} w(\boldsymbol{\theta}) \left[\int_{\Omega} h(\boldsymbol{\omega}) p(\boldsymbol{\omega} | \boldsymbol{\theta}, I) d\nu(\boldsymbol{\omega}) \right] p(\boldsymbol{\theta} | S) d\nu(\boldsymbol{\theta}) \\ &= (c_I / c_S) \int_{\Theta} \int_{\Omega} h(\boldsymbol{\omega}) p(\boldsymbol{\omega} | \boldsymbol{\theta}, I) p(\boldsymbol{\theta} | I) d\nu(\boldsymbol{\omega}) d\nu(\boldsymbol{\theta}) \\ &= (c_I / c_S) \mathbb{E}[h(\boldsymbol{\omega}) | I] = \bar{w} \cdot \bar{h} \\ \implies \bar{wh}^{(M)} &= M^{-1} \sum_{m=1}^M w(\boldsymbol{\theta}^{(m)}) h(\boldsymbol{\omega}^{(m)}) \xrightarrow{a.s.} \bar{w} \cdot \bar{h} \end{aligned}$$

Theorem 4.2 (c,d) Suppose

$$E[h(\omega) | I] = \bar{h} \quad \text{and} \quad \text{var}[h(\omega) | I] = \sigma^2$$

exist, and

$$w(\theta) = k(\theta | I) / k(\theta | S)$$

is bounded above on Θ . Then

$$M^{1/2} \left(\bar{h}^{(M)} - \bar{h} \right) \xrightarrow{d} N(\mathbf{0}, \tau^2)$$

and

$$\hat{\tau}^2(M) = \frac{M \sum_{m=1}^M \left[h(\omega^{(m)}) - \bar{h}^{(M)} \right]^2 w(\theta^{(m)})^2}{\left[\sum_{m=1}^M w(\theta^{(m)}) \right]^2} \xrightarrow{a.s.} \tau^2.$$

Proof.

$$\begin{aligned}
 \mathbf{E} \left[w(\boldsymbol{\theta})^2 h(\boldsymbol{\omega})^2 \mid S \right] &= \int_{\Theta} w(\boldsymbol{\theta})^2 \left[\int_{\Omega} h(\boldsymbol{\omega})^2 p(\boldsymbol{\omega} \mid \boldsymbol{\theta}, I) d\nu(\boldsymbol{\omega}) \right] p(\boldsymbol{\theta} \mid S) d\nu(\boldsymbol{\theta}) \\
 &= \int_{\Theta} w(\boldsymbol{\theta}) \frac{k(\boldsymbol{\theta} \mid I)}{k(\boldsymbol{\theta} \mid S)} \left[\int_{\Omega} h(\boldsymbol{\omega})^2 p(\boldsymbol{\omega} \mid \boldsymbol{\theta}, I) d\nu(\boldsymbol{\omega}) \right] p(\boldsymbol{\theta} \mid S) d\nu(\boldsymbol{\theta}) \\
 &= (c_I/c_S) \int_{\Theta} w(\boldsymbol{\theta}) \left[\int_{\Omega} h(\boldsymbol{\omega})^2 p(\boldsymbol{\omega} \mid \boldsymbol{\theta}, I) d\nu(\boldsymbol{\omega}) \right] p(\boldsymbol{\theta} \mid I) d\nu(\boldsymbol{\theta}) \\
 &\leq (c_I/c_S) \mathbf{E} \left[h(\boldsymbol{\omega})^2 \mid I \right] \sup_{\boldsymbol{\theta} \in \Theta} w(\boldsymbol{\theta}) < \infty
 \end{aligned}$$

Substitute $h(\boldsymbol{\omega}) = 1$ and conclude $\mathbf{E} \left[w(\boldsymbol{\theta})^2 \mid S \right] < \infty$ as well.

Proof (continued). Now we know

$$\mathbf{V} = \text{var} \begin{bmatrix} w(\boldsymbol{\theta}) \\ w(\boldsymbol{\theta}) h(\boldsymbol{\omega}) \mid S \end{bmatrix}$$

is finite. Apply Lindeberg-Levy Central Limit Theorem:

$$M^{1/2} \left[\begin{pmatrix} \bar{w}^{(M)} \\ \overline{wh}^{(M)} \end{pmatrix} - \begin{pmatrix} \bar{w} \\ \overline{wh} \end{pmatrix} \right] \xrightarrow{d} N(\mathbf{0}, \mathbf{V})$$

Proof (continued)

$$M^{1/2} \left[\begin{pmatrix} \bar{w}^{(M)} \\ \overline{wh}^{(M)} \end{pmatrix} - \begin{pmatrix} \bar{w} \\ \bar{w}\bar{h} \end{pmatrix} \right] \xrightarrow{d} N(\mathbf{0}, \mathbf{V}),$$

$$\frac{\overline{wh}^{(M)}}{\bar{w}^{(M)}} = \frac{\bar{w}\bar{h}}{\bar{w}} + \frac{\overline{wh}^{(M)} - \bar{w}\bar{h}}{\bar{w}} - \frac{\bar{w}\bar{h} \cdot (\bar{w}^{(M)} - \bar{w})}{\bar{w}^2} + o_p(M^{-1/2})$$

$$\implies M^{1/2} \left(\frac{\overline{wh}^{(M)}}{\bar{w}^{(M)}} - \bar{h} \right) \xrightarrow{d} N(0, \tau^2)$$

with

$$\tau^2 = \bar{w}^{-2} \left\{ \text{var} [w(\boldsymbol{\theta}) h(\boldsymbol{\omega}) \mid S] - 2\bar{h} \text{cov}_S [w(\boldsymbol{\theta}) h(\boldsymbol{\omega}), w(\boldsymbol{\theta}) \mid S] + \bar{h}^2 \text{var}_S [w(\boldsymbol{\theta}) \mid S] \right\} = \bar{w}^{-2} \text{var} \left\{ [w(\boldsymbol{\theta}) h(\boldsymbol{\omega}) - \bar{h}w(\boldsymbol{\theta})] \mid S \right\}.$$

Proof (concluded)

$$\begin{aligned} \tau^2 &= \bar{w}^{-2} \left\{ \text{var} [w(\boldsymbol{\theta}) h(\boldsymbol{\omega}) \mid S] - 2\bar{h} \text{cov}_S [w(\boldsymbol{\theta}) h(\boldsymbol{\omega}), w(\boldsymbol{\theta}) \mid S] \right. \\ &\quad \left. + \bar{h}^2 \text{var}_S [w(\boldsymbol{\theta}) \mid S] \right\} = \bar{w}^{-2} \text{var} \left\{ [w(\boldsymbol{\theta}) h(\boldsymbol{\omega}) - \bar{h}w(\boldsymbol{\theta})] \mid S \right\} \end{aligned}$$

$$\begin{aligned} \hat{\tau}^2(M) &= \frac{M \sum_{m=1}^M \left[h(\boldsymbol{\omega}^{(m)}) - \bar{h}^{(M)} \right]^2 w(\boldsymbol{\theta}^{(m)})^2}{\left[\sum_{m=1}^M w(\boldsymbol{\theta}^{(m)}) \right]^2} \\ &= \frac{M^{-1} \sum_{m=1}^M \left[w(\boldsymbol{\theta}^{(m)}) h(\boldsymbol{\omega}^{(m)}) - \bar{h}^{(M)} w(\boldsymbol{\theta}^{(m)}) \right]^2}{\left[M^{-1} \sum_{m=1}^M w(\boldsymbol{\theta}^{(m)}) \right]^2} \xrightarrow{a.s.} \tau^2 \end{aligned}$$

Example: Reweighting to a different prior distribution

Model	Prior	Likelihood
A_1	$p(\boldsymbol{\theta}_A A_1)$	$p(\mathbf{y} \boldsymbol{\theta}_A, A)$
A_2	$p(\boldsymbol{\theta}_A A_2)$	$p(\mathbf{y} \boldsymbol{\theta}_A, A)$

$$p(\boldsymbol{\theta}_A | \mathbf{y}^o, A_2) = p(\boldsymbol{\theta}_A | I)$$

$$\boldsymbol{\theta}_A^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta}_A | \mathbf{y}^o, A_1) = p(\boldsymbol{\theta}_A | S)$$

$$w(\boldsymbol{\theta}_A) = \frac{p(\boldsymbol{\theta}_A | \mathbf{y}^o, A_2)}{p(\boldsymbol{\theta}_A | \mathbf{y}^o, A_1)} = \frac{p(\boldsymbol{\theta}_A | A_1) p(\mathbf{y} | \boldsymbol{\theta}_A, A)}{p(\boldsymbol{\theta}_A | A_2) p(\mathbf{y} | \boldsymbol{\theta}_A, A)} = \frac{p(\boldsymbol{\theta}_A | A_1)}{p(\boldsymbol{\theta}_A | A_2)}$$

A useful summary measure of computational efficiency

$$\text{var} [h(\omega) | I] = \sigma^2 \quad \text{var} \left[\bar{h}^{(M)} | S \right] = \frac{\tau^2}{M}$$

$$w(\theta) \propto 1 \quad \forall \theta \iff p(\theta | S) = p(\theta | I) \implies \tau^2 = \sigma^2$$

Relative numerical efficiency (RNE) of the simulator is

$$RNE = \frac{\sigma^2}{\tau^2}.$$

$$\frac{\sigma^2}{M^*} = \frac{\tau^2}{M} \implies \frac{M^*}{M} = \frac{\sigma^2}{\tau^2} = RNE$$

Theorem 4.3 (a) Approximation of Bayes actions by importance sampling.

Conditions:

$$(1) \quad \boldsymbol{\theta}^{(m)} \stackrel{iid}{\sim} p(\boldsymbol{\theta} | S), \quad \boldsymbol{\omega}^{(m)} | \boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\omega} | \boldsymbol{\theta}^{(m)}, I)$$

$$(2) \quad w(\boldsymbol{\theta}) = k(\boldsymbol{\theta} | I) / k(\boldsymbol{\theta} | S) < c < \infty \quad \forall \boldsymbol{\theta}$$

$$(3) \quad L(\mathbf{a}, \boldsymbol{\omega}) \geq 0 \text{ defined on } \Omega \times A, \quad A \subseteq \mathbb{R}^q$$

$$(4) \quad R(\mathbf{a}) = \int_{\Omega} \int_{\Theta} L(\mathbf{a}, \boldsymbol{\omega}) p(\boldsymbol{\theta} | I) p(\boldsymbol{\omega} | \boldsymbol{\theta}, I) d\boldsymbol{\theta} d\boldsymbol{\omega}$$

$$(5) \quad R(\hat{\mathbf{a}}) < R(\mathbf{a}) \quad \forall \mathbf{a} \in A$$

Theorem 4.3 (a) (concluded)

More conditions:

(A1) $M^{-1} \sum_{m=1}^M L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) \rightarrow R(\mathbf{a})$ uniformly $\forall \mathbf{a} \in N(\hat{\mathbf{a}})$ (almost surely)

(A2) $\partial L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a}$ is a continuous function of $\mathbf{a} \forall \mathbf{a} \in N(\hat{\mathbf{a}})$

Denote

$$A_M = \left\{ \mathbf{a} : M^{-1} \sum_{m=1}^M \left[\partial L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a} \right] w(\boldsymbol{\theta}^{(m)}) = \mathbf{0} \right\}$$

Then for any $\varepsilon > 0$,

$$\lim_{M \rightarrow \infty} P \left[\inf_{\mathbf{a} \in A_M} (\mathbf{a} - \hat{\mathbf{a}})' (\mathbf{a} - \hat{\mathbf{a}}) > \varepsilon \mid S \right] = 0.$$

Theorem 4.3 (b) Approximation of Bayes actions by importance sampling.

Additional conditions:

(B1) $\partial^2 L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a} \partial \mathbf{a}'$ is a continuous function of $\mathbf{a} \forall \mathbf{a} \in N(\hat{\mathbf{a}}), \forall \boldsymbol{\omega} \in \Omega$

(B2) $\mathbf{B} = \text{var} \left[\partial L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a} \Big|_{\mathbf{a}=\hat{\mathbf{a}}} w(\boldsymbol{\theta})^{1/2} \mid S \right]$ exists and is finite

(B3) $\mathbf{H} = \mathbf{E} \left[\partial^2 L(\mathbf{a}, \boldsymbol{\omega}) / \partial \mathbf{a} \partial \mathbf{a}' \Big|_{\mathbf{a}=\hat{\mathbf{a}}} \mid S \right]$ exists and is finite and nonsingular

(B4) For any $\varepsilon > 0$, there exists M_ε such that

$$P \left[\sup_{\mathbf{a} \in N(\hat{\mathbf{a}})} \left| \partial^3 L(\mathbf{a}, \boldsymbol{\omega}) / \partial a_i \partial a_j \partial a_k \right| < M_\varepsilon \mid S \right] \geq 1 - \varepsilon$$

for all $i, j, k = 1, \dots, m$.

Theorem 4.3 (b) (concluded)

Then if $\hat{\mathbf{a}}_M$ is any element of A_M such that $\hat{\mathbf{a}}_M \xrightarrow{p} \hat{\mathbf{a}}$,

$$(1) \quad M^{1/2} (\hat{\mathbf{a}}_M - \hat{\mathbf{a}}) \xrightarrow{d} N(\mathbf{0}, \mathbf{H}^{-1} \mathbf{B} \mathbf{H}^{-1})$$

$$(2) \quad M^{-1} \sum_{m=1}^M w(\boldsymbol{\theta}^{(m)})^2 \partial L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a} |_{\mathbf{a}=\hat{\mathbf{a}}_M} \cdot \partial L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a}' |_{\mathbf{a}=\hat{\mathbf{a}}_M} \xrightarrow{p} \mathbf{B}$$

$$(3) \quad M^{-1} \sum_{m=1}^M w(\boldsymbol{\theta}^{(m)}) \partial^2 L(\mathbf{a}, \boldsymbol{\omega}^{(m)}) / \partial \mathbf{a} \partial \mathbf{a}' |_{\mathbf{a}=\hat{\mathbf{a}}_M} \xrightarrow{p} \mathbf{H}.$$

Markov Chain Monte Carlo Methods

The central idea

$$\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(m-1)}, C) \quad (m = 1, 2, 3, \dots)$$

If C is specified correctly, then

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

Better yet, if

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

then $J = I$. And even better,

$$p(\boldsymbol{\theta}^{(m)} \mid \boldsymbol{\theta}^{(0)}, C) \xrightarrow{d} p(\boldsymbol{\theta} \mid I) \quad \forall \boldsymbol{\theta}^{(0)} \in \Theta.$$

The central idea (continued)

If $p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(0)}, C) \xrightarrow{d} p(\boldsymbol{\theta} | I) \quad \forall \boldsymbol{\theta}^{(0)} \in \Theta$,

then we can approximate $E[h(\boldsymbol{\omega}) | I]$ by

- (1) iterating the chain B (“burn-in”) times;
- (2) drawing $\boldsymbol{\omega}^{(m)} \sim p(\boldsymbol{\omega} | \boldsymbol{\theta}^{(m)})$ ($m = 1, \dots, M$);
- (3) computing

$$\bar{h}^{(M)} = M^{-1} \sum_{m=1}^M h(\boldsymbol{\omega}^{(m)}).$$

The Gibbs sampler

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

Some notation: Corresponding to any subvector $\boldsymbol{\theta}_{(b)}$,

$$\boldsymbol{\theta}'_{<(b)} = (\boldsymbol{\theta}'_{(1)}, \dots, \boldsymbol{\theta}'_{(b-1)}) \quad (b = 2, \dots, B), \quad \boldsymbol{\theta}_{<(1)} = \emptyset$$

$$\boldsymbol{\theta}'_{>(b)} = (\boldsymbol{\theta}'_{(b+1)}, \dots, \boldsymbol{\theta}'_{(B)}) \quad (b = 1, \dots, B-1), \quad \boldsymbol{\theta}_{>(B)} = \emptyset$$

$$\boldsymbol{\theta}'_{-(b)} = (\boldsymbol{\theta}'_{<(b)}, \boldsymbol{\theta}'_{>(b)})$$

Very important: Choose the blocking so that

$$\boldsymbol{\theta}_{(b)} \sim p(\boldsymbol{\theta}_{(b)} \mid \boldsymbol{\theta}_{-(b)}, I)$$

is possible.

Intuitive argument for the Gibbs sampler

Imagine $\boldsymbol{\theta}^{(0)} \sim p(\boldsymbol{\theta} | I)$, and then in succession

$$\begin{aligned} \boldsymbol{\theta}_{(1)}^{(1)} &\sim p\left(\boldsymbol{\theta}_{(1)} \mid \boldsymbol{\theta}_{-(1)}^{(0)}, I\right), \\ \boldsymbol{\theta}_{(2)}^{(1)} &\sim p\left(\boldsymbol{\theta}_{(2)} \mid \boldsymbol{\theta}_{<(2)}^{(1)}, \boldsymbol{\theta}_{>(2)}^{(0)}, I\right), \\ \boldsymbol{\theta}_{(3)}^{(1)} &\sim p\left(\boldsymbol{\theta}_{(3)} \mid \boldsymbol{\theta}_{<(3)}^{(1)}, \boldsymbol{\theta}_{>(3)}^{(0)}, I\right), \\ &\vdots \\ \boldsymbol{\theta}_{(b)}^{(1)} &\sim p\left(\boldsymbol{\theta}_{(b)} \mid \boldsymbol{\theta}_{<(b)}^{(1)}, \boldsymbol{\theta}_{>(b)}^{(0)}, I\right) \\ &\vdots \\ \boldsymbol{\theta}_{(B)}^{(1)} &\sim p\left(\boldsymbol{\theta}_{(B)} \mid \boldsymbol{\theta}_{-(B)}^{(1)}, I\right) \end{aligned}$$

We have $\boldsymbol{\theta}^{(1)} \sim p(\boldsymbol{\theta} | I)$.

Now repeat

$$\begin{aligned}
 \boldsymbol{\theta}_{(1)}^{(2)} &\sim p\left(\boldsymbol{\theta}_{(1)} \mid \boldsymbol{\theta}_{-(1)}^{(1)}, I\right), \\
 \boldsymbol{\theta}_{(2)}^{(2)} &\sim p\left(\boldsymbol{\theta}_{(2)} \mid \boldsymbol{\theta}_{<(2)}^{(2)}, \boldsymbol{\theta}_{>(2)}^{(1)}, I\right), \\
 \boldsymbol{\theta}_{(3)}^{(2)} &\sim p\left(\boldsymbol{\theta}_{(3)} \mid \boldsymbol{\theta}_{<(3)}^{(2)}, \boldsymbol{\theta}_{>(3)}^{(1)}, I\right), \\
 &\vdots \\
 \boldsymbol{\theta}_{(b)}^{(2)} &\sim p\left(\boldsymbol{\theta}_{(b)} \mid \boldsymbol{\theta}_{<(b)}^{(2)}, \boldsymbol{\theta}_{>(b)}^{(1)}, I\right) \\
 &\vdots \\
 \boldsymbol{\theta}_{(B)}^{(2)} &\sim p\left(\boldsymbol{\theta}_{(B)} \mid \boldsymbol{\theta}_{-(B)}^{(2)}, I\right).
 \end{aligned}$$

We have $\boldsymbol{\theta}^{(2)} \sim p(\boldsymbol{\theta} \mid I)$.

The general step in the Gibbs sampler is

$$\boldsymbol{\theta}_{(b)}^{(m)} \sim p\left(\boldsymbol{\theta}_{(b)} \mid \boldsymbol{\theta}_{<(b)}^{(m)}, \boldsymbol{\theta}_{>(b)}^{(m-1)}, I\right)$$

for $b = 1, \dots, B$ and $m = 1, 2, \dots$

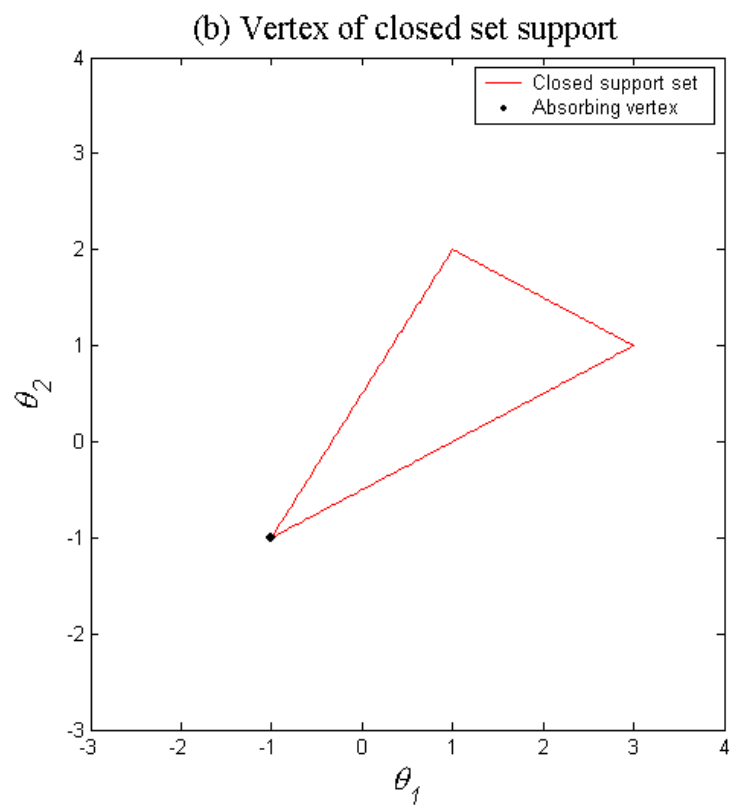
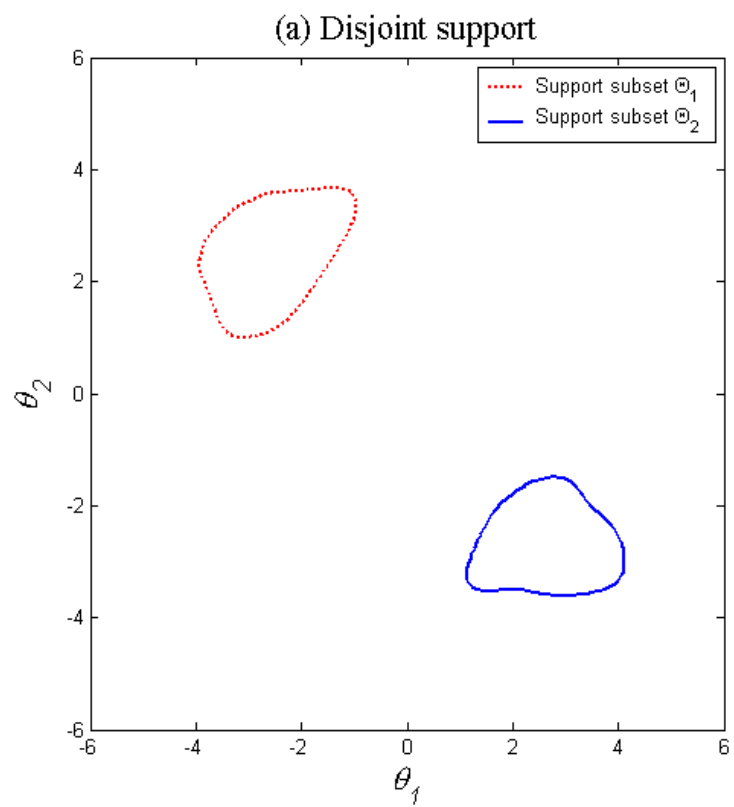
This defines the Markov chain

$$p\left(\boldsymbol{\theta}^{(m)} \mid \boldsymbol{\theta}^{(m-1)}, G\right) = \prod_{b=1}^B p\left[\boldsymbol{\theta}_{(b)}^{(m)} \mid \boldsymbol{\theta}_{<(b)}^{(m)}, \boldsymbol{\theta}_{>(b)}^{(m-1)}, I\right].$$

Key property:

$$\boldsymbol{\theta}^{(0)} \sim p(\boldsymbol{\theta} \mid I) \implies \boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta} \mid I)$$

Potential problems:



The Metropolis-Hastings Algorithm

What it does: $\boldsymbol{\theta}^* \sim q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H)$

Then

$$\begin{aligned} P(\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^*) &= \alpha(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H), \\ P(\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^{(m-1)}) &= 1 - \alpha(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H) \end{aligned}$$

where

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

Some aspects of the Metropolis-Hastings algorithm

If we define

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

then

$$P(\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^{(m-1)} | \boldsymbol{\theta}^{(m-1)} = \boldsymbol{\theta}, H) = r(\boldsymbol{\theta} | H) = 1 - \int_{\Theta} u(\boldsymbol{\theta}^* | \boldsymbol{\theta}, H) d\nu(\boldsymbol{\theta}^*).$$

Notice that

$$P(\boldsymbol{\theta}^{(m)} \in A | \boldsymbol{\theta}^{(m-1)} = \boldsymbol{\theta}, H) = \int_A u(\boldsymbol{\theta}^* | \boldsymbol{\theta}, H) d\nu(\boldsymbol{\theta}^*) + r(\boldsymbol{\theta} | H) I_A(\boldsymbol{\theta}).$$

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

We can write the transition density in one line making use of the Dirac delta function, an operator with the property

$$\int_A \delta_{\boldsymbol{\theta}}(\boldsymbol{\theta}^*) f(\boldsymbol{\theta}^*) d\nu(\boldsymbol{\theta}^*) = f(\boldsymbol{\theta}) I_A(\boldsymbol{\theta}).$$

Then

$$p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) = u(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H) + r(\boldsymbol{\theta}^{(m-1)} | H) \delta_{\boldsymbol{\theta}^{(m-1)}}(\boldsymbol{\theta}^{(m)}).$$

Special cases of the Metropolis-Hastings algorithm

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

Special case 1, original Metropolis (1953):

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

$$\implies \alpha(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H) = \min \left[p(\boldsymbol{\theta}^* | I) / p(\boldsymbol{\theta}^{(m-1)} | I), 1 \right]$$

Important example: *random walk Metropolis chain*

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

where $q(\cdot | H)$ is symmetric about zero.

Special cases of the Metropolis-Hastings algorithm

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

Special case 2, *Metropolis independence chain*:

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

where $w(\boldsymbol{\theta}) = p(\boldsymbol{\theta} | I) / q(\boldsymbol{\theta} | H)$.

Why does the Metropolis-Hastings algorithm work?

A two part argument – Part 1:

Suppose any transition probability density function $p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, T)$ satisfies the *reversibility condition*

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

with respect to $p(\boldsymbol{\theta} | I)$. Then

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

and so $p(\boldsymbol{\theta} | I)$ is an *invariant density* of the Markov chain.

Part 2 of the argument (How Hastings did it):

Suppose we don't know the probability $\alpha(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(m-1)}, H)$, but we want

$p(\boldsymbol{\theta}^{(m)} | \boldsymbol{\theta}^{(m-1)}, H)$ to be reversible with respect to $p(\boldsymbol{\theta} | I)$:

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

Trivial if $\boldsymbol{\theta}^{(m-1)} = \boldsymbol{\theta}^{(m)}$. For $\boldsymbol{\theta}^{(m-1)} \neq \boldsymbol{\theta}^{(m)}$ we need

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

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

Suppose without loss of generality that

$$p\left(\boldsymbol{\theta}^{(m-1)} \mid I\right) q\left(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(m-1)}, H\right) > p\left(\boldsymbol{\theta}^* \mid I\right) q\left(\boldsymbol{\theta}^{(m-1)} \mid \boldsymbol{\theta}^*, H\right).$$

Set $\alpha\left(\boldsymbol{\theta}^{(m-1)} \mid \boldsymbol{\theta}^*, H\right) = 1$ and

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