

Efficient Implementation of MCMC When Using An Unbiased Likelihood Estimator

Arnaud Doucet

University of Oxford

Joint work with M. Pitt, G. Deligiannidis & R. Kohn

Tokyo, 24/07/14

- Likelihood function $p_{\theta}(y)$ where $\theta \in \Theta \subseteq \mathbb{R}^d$.

Bayesian Inference

- Likelihood function $p_{\theta}(y)$ where $\theta \in \Theta \subseteq \mathbb{R}^d$.
- Prior distribution of density $p(\theta)$.

- Likelihood function $p_{\theta}(y)$ where $\theta \in \Theta \subseteq \mathbb{R}^d$.
- Prior distribution of density $p(\theta)$.
- Bayesian inference relies on the posterior

$$\pi(\theta) = p(\theta|y) = \frac{p_{\theta}(y) p(\theta)}{\int_{\Theta} p_{\theta'}(y) p(\theta') d\theta'}.$$

- Likelihood function $p_{\theta}(y)$ where $\theta \in \Theta \subseteq \mathbb{R}^d$.
- Prior distribution of density $p(\theta)$.
- Bayesian inference relies on the posterior

$$\pi(\theta) = p(\theta|y) = \frac{p_{\theta}(y) p(\theta)}{\int_{\Theta} p_{\theta'}(y) p(\theta') d\theta'}.$$

- For non-trivial models, inference relies typically on MCMC.

Metropolis-Hastings algorithm

Metropolis-Hastings algorithm

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.

Metropolis-Hastings algorithm

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- With probability

$$1 \wedge \frac{\pi(\vartheta) q(\vartheta_{i-1} | \vartheta)}{\pi(\vartheta_{i-1}) q(\vartheta | \vartheta_{i-1})} = 1 \wedge \frac{p_{\vartheta}(y) p(\vartheta) q(\vartheta_{i-1} | \vartheta)}{p_{\vartheta_{i-1}}(y) p(\vartheta_{i-1}) q(\vartheta | \vartheta_{i-1})},$$

set $\vartheta_i = \vartheta$, otherwise set $\vartheta_i = \vartheta_{i-1}$.

Metropolis-Hastings algorithm

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- With probability

$$1 \wedge \frac{\pi(\vartheta)}{\pi(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} = 1 \wedge \frac{p_{\vartheta}(y) p(\vartheta)}{p_{\vartheta_{i-1}}(y) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})},$$

set $\vartheta_i = \vartheta$, otherwise set $\vartheta_i = \vartheta_{i-1}$.

- **Problem:** Metropolis-Hastings (MH) cannot be implemented if $p_{\vartheta}(y)$ cannot be evaluated.

Intractable Likelihood Function

- For latent variable models, one has

$$p_{\theta}(y) = \int p_{\theta}(x, y) dx$$

where the integral cannot often be evaluated.

Intractable Likelihood Function

- For latent variable models, one has

$$p_{\theta}(y) = \int p_{\theta}(x, y) dx$$

where the integral cannot often be evaluated.

- A standard “solution” consists of using MCMC to sample from

$$p(\theta, x|y) \propto p_{\theta}(x, y) p(\theta)$$

by updating iterately x and θ .

Intractable Likelihood Function

- For latent variable models, one has

$$p_{\theta}(y) = \int p_{\theta}(x, y) dx$$

where the integral cannot often be evaluated.

- A standard “solution” consists of using MCMC to sample from

$$p(\theta, x|y) \propto p_{\theta}(x, y) p(\theta)$$

by updating iterately x and θ .

- Gibbs sampling strategies can be slow mixing and difficult to put in practice.

Intractable Likelihood Function

- For latent variable models, one has

$$p_{\theta}(y) = \int p_{\theta}(x, y) dx$$

where the integral cannot often be evaluated.

- A standard “solution” consists of using MCMC to sample from

$$p(\theta, x | y) \propto p_{\theta}(x, y) p(\theta)$$

by updating iterately x and θ .

- Gibbs sampling strategies can be slow mixing and difficult to put in practice.
- Could we use approximations of $p_{\theta}(y)$ within MH instead?

Pseudo-Marginal MH algorithm

- **Key Idea:** Replace $p_{\theta}(y)$ by an estimate $\hat{p}_{\theta}(y)$ in MH.

Pseudo-Marginal MH algorithm

- **Key Idea:** Replace $p_{\theta}(y)$ by an estimate $\hat{p}_{\theta}(y)$ in MH.

Pseudo-Marginal MH algorithm

- **Key Idea:** Replace $p_{\vartheta}(y)$ by an estimate $\hat{p}_{\vartheta}(y)$ in MH.

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.

Pseudo-Marginal MH algorithm

- **Key Idea:** Replace $p_{\vartheta}(y)$ by an estimate $\hat{p}_{\vartheta}(y)$ in MH.

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- Compute an estimate $\hat{p}_{\vartheta}(y)$ of $p_{\vartheta}(y)$.

Pseudo-Marginal MH algorithm

- **Key Idea:** Replace $p_{\vartheta}(y)$ by an estimate $\hat{p}_{\vartheta}(y)$ in MH.

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- Compute an estimate $\hat{p}_{\vartheta}(y)$ of $p_{\vartheta}(y)$.
- With probability

$$1 \wedge \frac{\hat{p}_{\vartheta}(y) p(\vartheta)}{\hat{p}_{\vartheta_{i-1}}(y) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})},$$

set $\vartheta_i = \vartheta$, $\hat{p}_{\vartheta_i}(y) = \hat{p}_{\vartheta}(y)$ otherwise set $\vartheta_i = \vartheta_{i-1}$,
 $\hat{p}_{\vartheta_i}(y) = \hat{p}_{\vartheta_{i-1}}(y)$.

Pseudo-Marginal MH algorithm

- **Key Idea:** Replace $p_{\theta}(y)$ by an estimate $\hat{p}_{\theta}(y)$ in MH.

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- Compute an estimate $\hat{p}_{\vartheta}(y)$ of $p_{\vartheta}(y)$.
- With probability

$$1 \wedge \underbrace{\frac{p(y; \vartheta)}{p(y; \vartheta_{i-1})} \frac{p(\vartheta)}{p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})}}_{\text{exact MH ratio}} \times \underbrace{\frac{\hat{p}(y; \vartheta) / p(y; \vartheta)}{\hat{p}(y; \vartheta_{i-1}) / p(y; \vartheta_{i-1})}}_{\text{noise}}$$

set $\vartheta_i = \vartheta$, $\hat{p}_{\vartheta_i}(y) = \hat{p}_{\vartheta}(y)$ otherwise set $\vartheta_i = \vartheta_{i-1}$,
 $\hat{p}_{\vartheta_i}(y) = \hat{p}_{\vartheta_{i-1}}(y)$.

- For latent variable models, one has

$$p_{\theta}(y) = \int p_{\theta}(x, y) dx = \int \frac{p_{\theta}(x, y)}{q_{\theta}(x)} q_{\theta}(x) dx$$

where $q_{\theta}(x)$ is an importance sampling density.

Importance Sampling Estimator

- For latent variable models, one has

$$p_{\theta}(y) = \int p_{\theta}(x, y) dx = \int \frac{p_{\theta}(x, y)}{q_{\theta}(x)} q_{\theta}(x) dx$$

where $q_{\theta}(x)$ is an importance sampling density.

- An unbiased estimator is given by

$$\hat{p}_{\theta}(y) = \frac{1}{N} \sum_{k=1}^N \frac{p_{\theta}(X^k, y)}{q_{\theta}(X^k)}, \quad X^k \stackrel{\text{i.i.d.}}{\sim} q_{\theta}(\cdot)$$

Sequential Monte Carlo Estimator

- $\{X_t\}_{t \geq 1}$ is a \mathbb{X} -valued latent Markov process with $X_1 \sim \mu(\cdot; \theta)$ and $X_{t+1} | X_t \sim f(\cdot | X_t; \theta)$.

Sequential Monte Carlo Estimator

- $\{X_t\}_{t \geq 1}$ is a \mathbb{X} -valued latent Markov process with $X_1 \sim \mu(\cdot; \theta)$ and $X_{t+1} | X_t \sim f(\cdot | X_t; \theta)$.
- Observations $\{Y_t\}_{t \geq 1}$ are conditionally independent given $\{X_t\}_{t \geq 0}$ with $Y_t | \{X_k\}_{k \geq 0} \sim g(\cdot | X_t, \theta)$.

Sequential Monte Carlo Estimator

- $\{X_t\}_{t \geq 1}$ is a \mathbb{X} -valued latent Markov process with $X_1 \sim \mu(\cdot; \theta)$ and $X_{t+1} | X_t \sim f(\cdot | X_t; \theta)$.
- Observations $\{Y_t\}_{t \geq 1}$ are conditionally independent given $\{X_t\}_{t \geq 0}$ with $Y_t | \{X_k\}_{k \geq 0} \sim g(\cdot | X_t, \theta)$.
- Likelihood of $y_{1:T} = (y_1, \dots, y_T)$ is

$$p(y_{1:T}; \theta) = \int_{\mathbb{X}^T} p(x_{1:T}, y_{1:T}; \theta) dx_{1:T}.$$

Sequential Monte Carlo Estimator

- $\{X_t\}_{t \geq 1}$ is a \mathbb{X} -valued latent Markov process with $X_1 \sim \mu(\cdot; \theta)$ and $X_{t+1} | X_t \sim f(\cdot | X_t; \theta)$.
- Observations $\{Y_t\}_{t \geq 1}$ are conditionally independent given $\{X_t\}_{t \geq 0}$ with $Y_t | \{X_k\}_{k \geq 0} \sim g(\cdot | X_t, \theta)$.
- Likelihood of $y_{1:T} = (y_1, \dots, y_T)$ is

$$p(y_{1:T}; \theta) = \int_{\mathbb{X}^T} p(x_{1:T}, y_{1:T}; \theta) dx_{1:T}.$$

- SMC provides an unbiased estimator of relative variance $\mathcal{O}(T/N)$ where N is the number of particles.

Main Result

- **Proposition:** Let $\hat{p}_\theta(y)$ be a non-negative unbiased estimator then the pseudo-marginal MH kernel admits an invariant distribution admitting $\pi(\theta)$ as a marginal.

Main Result

- **Proposition:** Let $\hat{p}_\theta(y)$ be a non-negative unbiased estimator then the pseudo-marginal MH kernel admits an invariant distribution admitting $\pi(\theta)$ as a marginal.
- “Proof”. Define $Z = \log \hat{p}(y; \theta) / p(y; \theta)$ and an auxiliary target density on $\Theta \times \mathbb{R}$

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z)g_\theta(z)}_{\text{unbiasedness} \Leftrightarrow \int(\cdot)dz=1}$$

where $Z \sim g_\theta$; e.g. importance sampling or particle filter.

Main Result

- **Proposition:** Let $\hat{p}_\vartheta(y)$ be a non-negative unbiased estimator then the pseudo-marginal MH kernel admits an invariant distribution admitting $\pi(\theta)$ as a marginal.
- “Proof”. Define $Z = \log \hat{p}(y; \theta) / p(y; \theta)$ and an auxiliary target density on $\Theta \times \mathbb{R}$

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g_\theta(z)}_{\text{unbiasedness} \Leftrightarrow \int(\cdot) dz = 1}$$

where $Z \sim g_\theta$; e.g. importance sampling or particle filter.

- Pseudo marginal MH is MH of target $\bar{\pi}(\theta, z)$ and proposal $q(\theta, \vartheta) g_\vartheta(z)$ as

$$\frac{\bar{\pi}(\vartheta, Z)}{\bar{\pi}(\vartheta_{i-1}, Z_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta) g_{\vartheta_{i-1}}(Z_{i-1})}{q(\vartheta | \vartheta_{i-1}) g_\vartheta(Z)} = \frac{\hat{p}(y; \vartheta)}{\hat{p}(y; \vartheta_{i-1})} \frac{p(\vartheta)}{p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})}$$

A Nonlinear State-Space Model

- Standard non-linear model

$$X_t = \frac{1}{2}X_{t-1} + 25\frac{X_{t-1}}{1+X_{t-1}^2} + 8\cos(1.2t) + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_V^2),$$

$$Y_t = \frac{1}{20}X_t^2 + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_W^2).$$

A Nonlinear State-Space Model

- Standard non-linear model

$$X_t = \frac{1}{2}X_{t-1} + 25\frac{X_{t-1}}{1+X_{t-1}^2} + 8\cos(1.2t) + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_V^2),$$

$$Y_t = \frac{1}{20}X_t^2 + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_W^2).$$

- $T = 200$ data points with $\theta = (\sigma_V^2, \sigma_W^2) = (10, 10)$.

A Nonlinear State-Space Model

- Standard non-linear model

$$X_t = \frac{1}{2}X_{t-1} + 25\frac{X_{t-1}}{1+X_{t-1}^2} + 8\cos(1.2t) + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_V^2),$$

$$Y_t = \frac{1}{20}X_t^2 + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_W^2).$$

- $T = 200$ data points with $\theta = (\sigma_V^2, \sigma_W^2) = (10, 10)$.
- Difficult to perform standard MCMC as $p(x_{1:T} | y_{1:T}, \theta)$ is highly multimodal.

A Nonlinear State-Space Model

- Standard non-linear model

$$X_t = \frac{1}{2}X_{t-1} + 25\frac{X_{t-1}}{1+X_{t-1}^2} + 8\cos(1.2t) + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_V^2),$$

$$Y_t = \frac{1}{20}X_t^2 + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_W^2).$$

- $T = 200$ data points with $\theta = (\sigma_V^2, \sigma_W^2) = (10, 10)$.
- Difficult to perform standard MCMC as $p(x_{1:T} | y_{1:T}, \theta)$ is highly multimodal.
- We sample from $p(\theta | y_{1:T})$ using a random walk pseudo-marginal MH where $p_\theta(y_{1:T})$ is estimated using SMC with N particles.

A Nonlinear State-Space Model

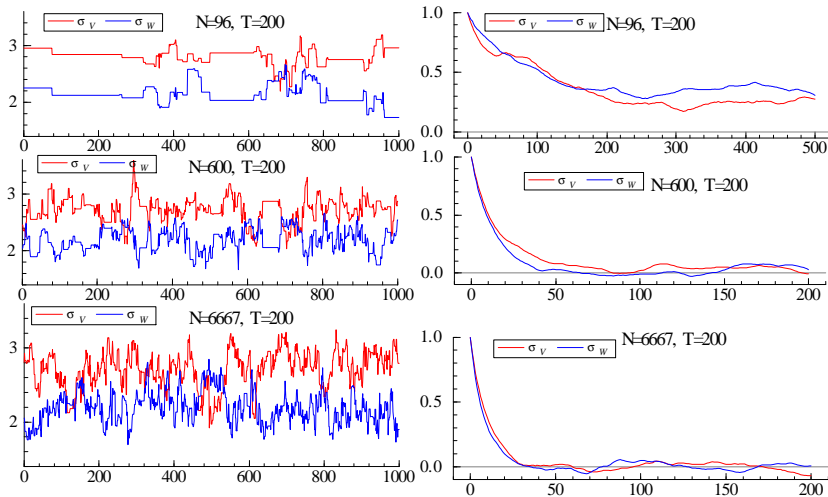


Figure: Autocorrelation of $\{\sigma_v^{(i)}\}$ and $\{\sigma_w^{(i)}\}$ of the MH sampler for various N .

Optimal Tuning of the Pseudo-marginal MH

- A key issue from a practical point of view is the selection of N .

Optimal Tuning of the Pseudo-marginal MH

- A key issue from a practical point of view is the selection of N .
- If N is too small, then the algorithm mixes poorly and will require many MCMC iterations.

Optimal Tuning of the Pseudo-marginal MH

- A key issue from a practical point of view is the selection of N .
- If N is too small, then the algorithm mixes poorly and will require many MCMC iterations.
- If N is too large, then each pseudo-marginal MH iteration is expensive.

Inefficiency of the Pseudo-marginal MH

- Consider the particle MH chain $\{\theta_i, Z_i\}$ of $\bar{\pi}$ -invariant transition kernel Q

$$Q\{(\theta, z), (d\vartheta, dw)\} = q(\vartheta|\theta)g_{\vartheta}(w) \min\{1, r(\theta, \vartheta) \exp(w - z)\} d\vartheta dw \\ + \{1 - \varrho_Q(\theta, z)\} \delta_{(\theta, z)}(d\vartheta, dw)$$

where $r(\theta, \vartheta) = \pi(\vartheta)q(\theta|\vartheta) / \{\pi(\theta)q(\vartheta|\theta)\}$.

Inefficiency of the Pseudo-marginal MH

- Consider the particle MH chain $\{\theta_i, Z_i\}$ of $\bar{\pi}$ -invariant transition kernel Q

$$Q\{(\theta, z), (d\vartheta, dw)\} = q(\vartheta|\theta)g_{\vartheta}(w) \min\{1, r(\theta, \vartheta) \exp(w - z)\} d\vartheta dw \\ + \{1 - \varrho_Q(\theta, z)\} \delta_{(\theta, z)}(d\vartheta, dw)$$

where $r(\theta, \vartheta) = \pi(\vartheta)q(\theta|\vartheta) / \{\pi(\theta)q(\vartheta|\theta)\}$.

- Proposition** (KV 1986). Let $h : \Theta \rightarrow \mathbb{R}$, $\pi(h) = \mathbb{E}_{\pi}[h(\theta)]$ and $\hat{\pi}_n(h) = n^{-1} \sum_{i=1}^n h(\theta_i)$. If $\{\theta_i, Z_i\}$ is stationary and ergodic, $\mathbb{V}_{\pi}[h(\theta)] < \infty$ and $IF_h^Q(\sigma) = 1 + 2 \sum_{\tau=1}^{\infty} \text{corr}_{\bar{\pi}, Q}\{h(\theta_0), h(\theta_{\tau})\} < \infty$ then

$$\sqrt{n} \{\hat{\pi}_n(h) - \pi(h)\} \rightarrow \mathcal{N}\left(0, \mathbb{V}_{\pi}[h(\theta)] IF_h^Q(\sigma)\right).$$

Inefficiency of the Pseudo-marginal MH

- Consider the particle MH chain $\{\theta_i, Z_i\}$ of $\bar{\pi}$ -invariant transition kernel Q

$$Q\{(\theta, z), (d\vartheta, dw)\} = q(\vartheta|\theta)g_\vartheta(w) \min\{1, r(\theta, \vartheta) \exp(w - z)\} d\vartheta dw \\ + \{1 - \varrho_Q(\theta, z)\} \delta_{(\theta, z)}(d\vartheta, dw)$$

where $r(\theta, \vartheta) = \pi(\vartheta)q(\theta|\vartheta) / \{\pi(\theta)q(\vartheta|\theta)\}$.

- Proposition** (KV 1986). Let $h : \Theta \rightarrow \mathbb{R}$, $\pi(h) = \mathbb{E}_\pi[h(\theta)]$ and $\hat{\pi}_n(h) = n^{-1} \sum_{i=1}^n h(\theta_i)$. If $\{\theta_i, Z_i\}$ is stationary and ergodic, $\mathbb{V}_\pi[h(\theta)] < \infty$ and $IF_h^Q(\sigma) = 1 + 2 \sum_{\tau=1}^{\infty} \text{corr}_{\bar{\pi}, Q}\{h(\theta_0), h(\theta_\tau)\} < \infty$ then

$$\sqrt{n} \{\hat{\pi}_n(h) - \pi(h)\} \rightarrow \mathcal{N}\left(0, \mathbb{V}_\pi[h(\theta)] IF_h^Q(\sigma)\right).$$

- The **Integrated Autocorrelation Time** IF_h^Q is a measure of **inefficiency** of Q which **we want to minimize for a fixed computational budget**.

Aim of the Analysis

- **Simplifying Assumption:** The noise Z is independent of θ and Gaussian; i.e. $Z \sim \mathcal{N}(-\sigma^2/2; \sigma^2)$:

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g_{\sigma}(z)}_{\pi_{Z,\sigma}(z)} = \pi(\theta) \mathcal{N}(z; \sigma^2/2; \sigma^2).$$

Aim of the Analysis

- **Simplifying Assumption:** The noise Z is independent of θ and Gaussian; i.e. $Z \sim \mathcal{N}(-\sigma^2/2; \sigma^2)$:

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g_{\sigma}(z)}_{\pi_{Z, \sigma}(z)} = \pi(\theta) \mathcal{N}(z; \sigma^2/2; \sigma^2).$$

- **Aim:** Minimize the computational cost

$$CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$$

as $\sigma^2 \propto 1/N$ and computational efforts proportional to N .

Aim of the Analysis

- **Simplifying Assumption:** The noise Z is independent of θ and Gaussian; i.e. $Z \sim \mathcal{N}(-\sigma^2/2; \sigma^2)$:

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g_{\sigma}(z)}_{\pi_{Z, \sigma}(z)} = \pi(\theta) \mathcal{N}(z; \sigma^2/2; \sigma^2).$$

- **Aim:** Minimize the computational cost

$$CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$$

as $\sigma^2 \propto 1/N$ and computational efforts proportional to N .

- **Special cases:**

Aim of the Analysis

- **Simplifying Assumption:** The noise Z is independent of θ and Gaussian; i.e. $Z \sim \mathcal{N}(-\sigma^2/2; \sigma^2)$:

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g_{\sigma}(z)}_{\pi_{Z, \sigma}(z)} = \pi(\theta) \mathcal{N}(z; \sigma^2/2; \sigma^2).$$

- **Aim:** Minimize the computational cost

$$CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$$

as $\sigma^2 \propto 1/N$ and computational efforts proportional to N .

- **Special cases:**
- ① When $q(\vartheta|\theta) = p(\vartheta|y)$, $\sigma_{\text{opt}} = 0.92$ (Pitt et al., 2012).

Aim of the Analysis

- **Simplifying Assumption:** The noise Z is independent of θ and Gaussian; i.e. $Z \sim \mathcal{N}(-\sigma^2/2; \sigma^2)$:

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g_{\sigma}(z)}_{\pi_{Z, \sigma}(z)} = \pi(\theta) \mathcal{N}(z; \sigma^2/2; \sigma^2).$$

- **Aim:** Minimize the computational cost

$$CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$$

as $\sigma^2 \propto 1/N$ and computational efforts proportional to N .

- **Special cases:**

- 1 When $q(\vartheta|\theta) = p(\vartheta|y)$, $\sigma_{\text{opt}} = 0.92$ (Pitt et al., 2012).
- 2 When $\pi(\theta) = \prod_{i=1}^d f(\theta_i)$ and $q(\vartheta|\theta)$ is an isotropic Gaussian random walk then, as $d \rightarrow \infty$, $\sigma_{\text{opt}} = 1.81$ (Sherlock, Thiery, Roberts & Rosenthal, 2014).

Sketch of the Analysis

- For general proposals and targets, direct minimization of $CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$ impossible so minimize an upper bound over it.

Sketch of the Analysis

- For general proposals and targets, direct minimization of $CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$ impossible so minimize an upper bound over it.
- We introduce an auxiliary $\bar{\pi}$ -invariant kernel

$$Q^* \{(\theta, z), (d\vartheta, dw)\} = q(\vartheta|\theta) g_\sigma(w) \alpha_{EX}(\theta, \vartheta) \alpha_Z(z, w) d\vartheta dw \\ + \{1 - \varrho_{EX}(\theta) \varrho_{Z,\sigma}(z)\} \delta_{(\theta,z)}(d\vartheta, dw)$$

where

$$\alpha_{EX}(\theta, \vartheta) = \min\{1, r(\theta, \vartheta)\}, \quad \alpha_Z(z, w) = \min\{1, \exp(w - z)\}$$

Sketch of the Analysis

- For general proposals and targets, direct minimization of $CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$ impossible so minimize an upper bound over it.
- We introduce an auxiliary $\bar{\pi}$ -invariant kernel

$$Q^* \{(\theta, z), (d\vartheta, dw)\} = q(\vartheta|\theta) g_\sigma(w) \alpha_{EX}(\theta, \vartheta) \alpha_Z(z, w) d\vartheta dw \\ + \{1 - \rho_{EX}(\theta) \rho_{Z,\sigma}(z)\} \delta_{(\theta,z)}(d\vartheta, dw)$$

where

$$\alpha_{EX}(\theta, \vartheta) = \min\{1, r(\theta, \vartheta)\}, \quad \alpha_Z(z, w) = \min\{1, \exp(w - z)\}$$

- Peskun's theorem (1973) guarantees that $IF_h^Q(\sigma) \leq IF_h^{Q^*}(\sigma)$ so that $CT_h^Q(\sigma) \leq CT_h^{Q^*}(\sigma)$.

Sketch of the Analysis

- Let $(\theta_i, Z_i)_{i \geq 1}$ be generated by Q^* .

Sketch of the Analysis

- Let $(\theta_i, Z_i)_{i \geq 1}$ be generated by Q^* .
- Denote $(\tilde{\theta}_i, \tilde{Z}_i)_{i \geq 1}$ the accepted proposals and $(\tau_i)_{i \geq 1}$ the associated sojourn times; i.e. $(\tilde{\theta}_1, \tilde{Z}_1) = (\theta_1, Z_1) = \dots = (\theta_{\tau_1}, Z_{\tau_1})$,
 $(\tilde{\theta}_2, \tilde{Z}_2) = (\theta_{\tau_1+1}, Z_{\tau_1+1}) = \dots = (\theta_{\tau_2}, Z_{\tau_2})$ and so on where
 $(\tilde{\theta}_{i+1}, \tilde{Z}_{i+1}) \neq (\tilde{\theta}_i, \tilde{Z}_i)$ a.s.

Sketch of the Analysis

- Let $(\theta_i, Z_i)_{i \geq 1}$ be generated by Q^* .
- Denote $(\tilde{\theta}_i, \tilde{Z}_i)_{i \geq 1}$ the accepted proposals and $(\tau_i)_{i \geq 1}$ the associated sojourn times; i.e. $(\tilde{\theta}_1, \tilde{Z}_1) = (\theta_1, Z_1) = \dots = (\theta_{\tau_1}, Z_{\tau_1})$, $(\tilde{\theta}_2, \tilde{Z}_2) = (\theta_{\tau_1+1}, Z_{\tau_1+1}) = \dots = (\theta_{\tau_2}, Z_{\tau_2})$ and so on where $(\tilde{\theta}_{i+1}, \tilde{Z}_{i+1}) \neq (\tilde{\theta}_i, \tilde{Z}_i)$ a.s.
- $IF_h^{Q^*}(\sigma)$ can be re-expressed in terms of $IF_{h/(q_{EX}q_Z)}^{\tilde{Q}^*}(\sigma)$ where

$$\begin{aligned} \tilde{Q}^* \{(\theta, z), (d\vartheta, dw)\} &= \tilde{Q}^{EX}(\theta, d\vartheta) \tilde{Q}_\sigma^Z(z, dw) \\ &= \frac{q(d\vartheta|\theta) \alpha_{EX}(\theta, \vartheta)}{q_{EX}(\theta)} \frac{g_\sigma(dw) \alpha_Z(z, w)}{q_{Z,\sigma}(z)} \end{aligned}$$

- **Proposition:** Under weak assumptions, we have $IF_h^Q(\sigma) \leq IF_h^{Q^*}(\sigma)$ where

$$\begin{aligned} IF_h^{Q^*}(\sigma) &= 2 \frac{\{1 + IF_h^{\text{EX}}\}}{1 + IF_{h/\varrho_{\text{EX}}}^{\tilde{Q}^{\text{EX}}}} \left\{ \pi_{Z,\sigma}(1/\varrho_{Z,\sigma}) - 1/\pi_{Z,\sigma}(\varrho_{Z,\sigma}) \right\} \\ &\quad \times \sum_{n=0}^{\infty} \phi_n\left(h/\varrho_{\text{EX}}, \tilde{Q}^{\text{EX}}\right) \phi_n\left(1/\varrho_Z, \tilde{Q}_\sigma^Z\right) \\ &\quad + \frac{1 + IF_h^{\text{EX}}}{\pi_{Z,\sigma}(\varrho_{Z,\sigma})} - 1, \end{aligned}$$

where $\phi_n(\varphi, P)$ denotes the autocorrelation at lag n under a Markov kernel P .

- **Proposition:** Under weak assumptions, we have $IF_h^Q(\sigma) \leq IF_h^{Q^*}(\sigma)$ where

$$\begin{aligned} IF_h^{Q^*}(\sigma) &= 2 \frac{\{1 + IF_h^{\text{EX}}\}}{1 + IF_{h/\varrho_{\text{EX}}}^{\tilde{Q}^{\text{EX}}}} \left\{ \pi_{Z,\sigma}(1/\varrho_{Z,\sigma}) - 1/\pi_{Z,\sigma}(\varrho_{Z,\sigma}) \right\} \\ &\quad \times \sum_{n=0}^{\infty} \phi_n\left(h/\varrho_{\text{EX}}, \tilde{Q}^{\text{EX}}\right) \phi_n\left(1/\varrho_Z, \tilde{Q}_\sigma^Z\right) \\ &\quad + \frac{1 + IF_h^{\text{EX}}}{\pi_{Z,\sigma}(\varrho_{Z,\sigma})} - 1, \end{aligned}$$

where $\phi_n(\varphi, P)$ denotes the autocorrelation at lag n under a Markov kernel P .

- This identity allows us to “decouple” the influence of the parameter and of the noise on $IF_h^{Q^*}(\sigma)$.

Simpler Bounds on the Relative Inefficiency

- If $IF_{h/\varrho_{EX}}^{\tilde{Q}^{EX}} \geq 1$, e.g. \tilde{Q}^{EX} is a positive kernel, then

$$\frac{IF_h^Q(\sigma)}{IF_h^{EX}} \leq \frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \leq \frac{1}{2} \left(1 + \frac{1}{IF_h^{EX}} \right) \pi_{Z,\sigma}(1/\varrho_{Z,\sigma}) - \frac{1}{IF_h^{EX}}$$

and the bound is tight as $IF_h^{EX} \rightarrow 1$ or $\sigma \rightarrow 0$.

Simpler Bounds on the Relative Inefficiency

- If $IF_{h/\varrho_{EX}}^{\tilde{Q}^{EX}} \geq 1$, e.g. \tilde{Q}^{EX} is a positive kernel, then

$$\frac{IF_h^Q(\sigma)}{IF_h^{EX}} \leq \frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \leq \frac{1}{2} \left(1 + \frac{1}{IF_h^{EX}} \right) \pi_{Z,\sigma}(1/\varrho_{Z,\sigma}) - \frac{1}{IF_h^{EX}}$$

and the bound is tight as $IF_h^{EX} \rightarrow 1$ or $\sigma \rightarrow 0$.

- As $IF_{J,h/\varrho_{EX}}^{EX} \rightarrow \infty$,

$$\frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \rightarrow \frac{1}{\pi_{Z,\sigma}(\varrho_{Z,\sigma})}.$$

Simpler Bounds on the Relative Inefficiency

- If $IF_{h/\varrho_{EX}}^{\tilde{Q}^{EX}} \geq 1$, e.g. \tilde{Q}^{EX} is a positive kernel, then

$$\frac{IF_h^Q(\sigma)}{IF_h^{EX}} \leq \frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \leq \frac{1}{2} \left(1 + \frac{1}{IF_h^{EX}} \right) \pi_{Z,\sigma}(1/\varrho_{Z,\sigma}) - \frac{1}{IF_h^{EX}}$$

and the bound is tight as $IF_h^{EX} \rightarrow 1$ or $\sigma \rightarrow 0$.

- As $IF_{J,h/\varrho_{EX}}^{EX} \rightarrow \infty$,

$$\frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \rightarrow \frac{1}{\pi_{Z,\sigma}(\varrho_{Z,\sigma})}.$$

- Results used to minimize w.r.t σ upper bounds on $CT_h^Q(\sigma) = IF_h^Q(\sigma) / \sigma^2$.

Bounds on Relative Computational Costs

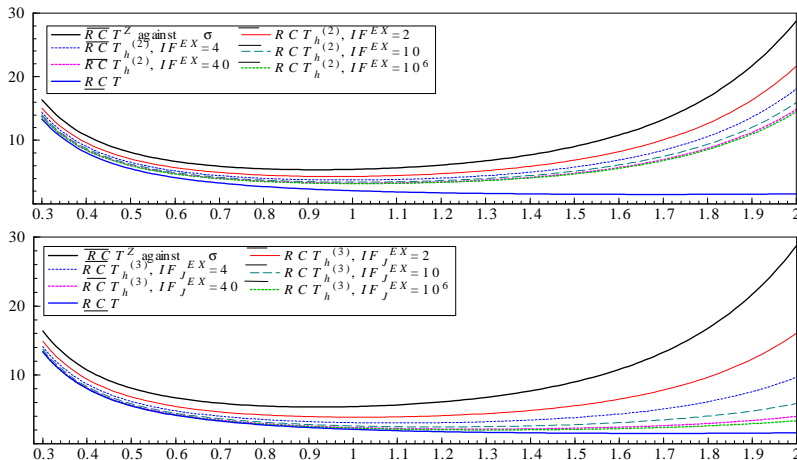


Figure: Bounds on $IF_h^Q(\sigma) / (\sigma^2 IF_h^{EX})$

Practical Guidelines

- For good proposals, select $\sigma \approx 1$ whereas for poor proposals, select $\sigma \approx 1.7$.

Practical Guidelines

- For good proposals, select $\sigma \approx 1$ whereas for poor proposals, select $\sigma \approx 1.7$.
- When you have no clue about the proposal efficiency,

Practical Guidelines

- For good proposals, select $\sigma \approx 1$ whereas for poor proposals, select $\sigma \approx 1.7$.
- When you have no clue about the proposal efficiency,
- ① If $\sigma_{\text{opt}} = 1$ and you pick $\sigma = 1.7$, computing time increases by $\approx 150\%$.

Practical Guidelines

- For good proposals, select $\sigma \approx 1$ whereas for poor proposals, select $\sigma \approx 1.7$.
 - When you have no clue about the proposal efficiency,
- 1 If $\sigma_{\text{opt}} = 1$ and you pick $\sigma = 1.7$, computing time increases by $\approx 150\%$.
 - 2 If $\sigma_{\text{opt}} = 1.7$ and you pick $\sigma = 1$, computing time increases by $\approx 50\%$.

Practical Guidelines

- For good proposals, select $\sigma \approx 1$ whereas for poor proposals, select $\sigma \approx 1.7$.
 - When you have no clue about the proposal efficiency,
- 1 If $\sigma_{\text{opt}} = 1$ and you pick $\sigma = 1.7$, computing time increases by $\approx 150\%$.
 - 2 If $\sigma_{\text{opt}} = 1.7$ and you pick $\sigma = 1$, computing time increases by $\approx 50\%$.
 - 3 If $\sigma_{\text{opt}} = 1$ or $\sigma_{\text{opt}} = 1.7$ and you pick $\sigma = 1.2$, computing time increases by $\approx 15\%$.

Example: Noisy Autoregressive Example

- Consider

$$X_t = \mu(1 - \phi) + \phi X_t + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\eta^2),$$

$$Y_t = X_t + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2),$$

where $\theta = (\phi, \mu, \sigma_\eta^2)$.

Example: Noisy Autoregressive Example

- Consider

$$X_t = \mu(1 - \phi) + \phi X_t + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\eta^2),$$

$$Y_t = X_t + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2),$$

where $\theta = (\phi, \mu, \sigma_\eta^2)$.

- Likelihood can be computed exactly using Kalman.

Example: Noisy Autoregressive Example

- Consider

$$X_t = \mu(1 - \phi) + \phi X_t + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\eta^2),$$

$$Y_t = X_t + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2),$$

where $\theta = (\phi, \mu, \sigma_\eta^2)$.

- Likelihood can be computed exactly using Kalman.
- Autoregressive Metropolis proposal of coefficient ρ for θ based on multivariate t-distribution.

Example: Noisy Autoregressive Example

- Consider

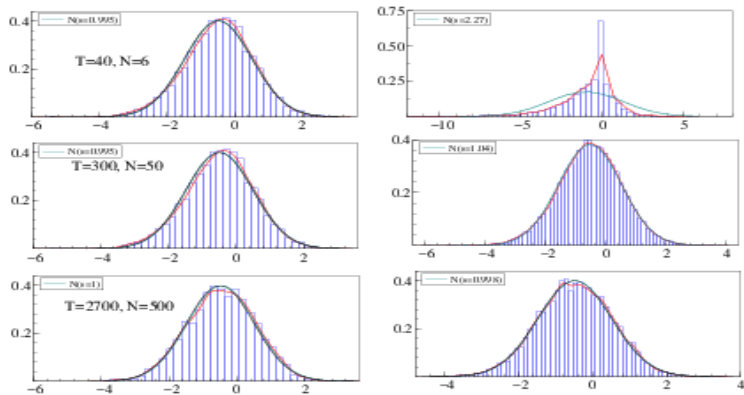
$$X_t = \mu(1 - \phi) + \phi X_t + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\eta^2),$$

$$Y_t = X_t + W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2),$$

where $\theta = (\phi, \mu, \sigma_\eta^2)$.

- Likelihood can be computed exactly using Kalman.
- Autoregressive Metropolis proposal of coefficient ρ for θ based on multivariate t-distribution.
- N is selected so as to obtain $\sigma(\bar{\theta}) \approx \text{constant}$ where $\bar{\theta}$ posterior mean.

Empirical vs Asymptotic Distribution of Log-Likelihood Estimator



Empirical distribution of Z at posterior mean (left) and marginalized over samples from $\pi q(\vartheta) = \int \pi(\theta) q(\theta, \vartheta) d\theta$.

Relative Inefficiency and Computing Time

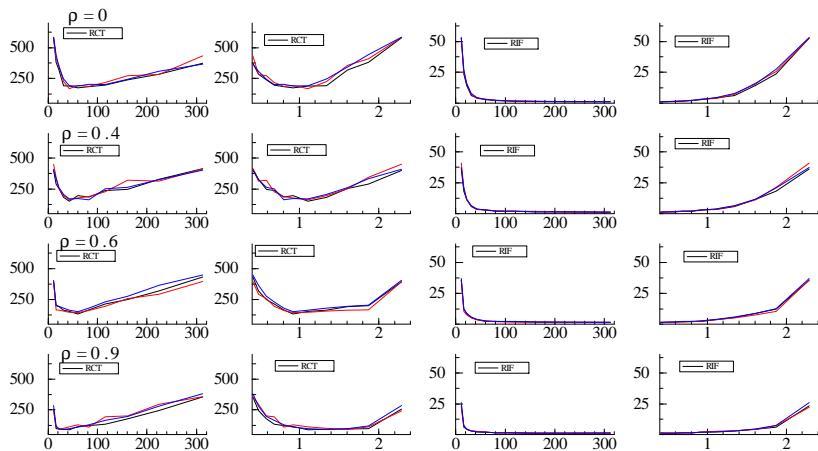


Figure: From left to right: RCT_h^Q vs N , RCT_h^Q vs $\sigma(\bar{\theta})$, RIF_h^Q against N and RIF_h^Q against $\sigma(\bar{\theta})$ for various values of ρ and different parameters.

- Simplified quantitative analysis of the particle MH algorithm.

- Simplified quantitative analysis of the particle MH algorithm.
- Particle MH scales roughly in $O(T^2)$.

- Simplified quantitative analysis of the particle MH algorithm.
- Particle MH scales roughly in $O(T^2)$.
- Particle Gibbs sampling displays better theoretical properties: scaling?