

Pseudo-marginal MCMC methods for inference in latent variable models

Arnaud Doucet

Department of Statistics, Oxford University

Joint work with George Deligiannidis (Oxford) & Mike Pitt (Kings)

07/07/2016

Organization of the talk

- Latent variable models

Organization of the talk

- Latent variable models
- The pseudo-marginal method

Organization of the talk

- Latent variable models
- The pseudo-marginal method
- Optimal tuning

Organization of the talk

- Latent variable models
- The pseudo-marginal method
- Optimal tuning
- The correlated pseudo-marginal method

Organization of the talk

- Latent variable models
- The pseudo-marginal method
- Optimal tuning
- The correlated pseudo-marginal method
- Illustrations

- Assume

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mu_\theta(\cdot), \quad Y_t | (X_t = x) \sim g_\theta(\cdot | x) \quad \text{for } t = 1, \dots, T$$

where $(X_t)_{t \geq 1}$ are latent variables and $(Y_t)_{t \geq 1}$ correspond to observations.

- Assume

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mu_\theta(\cdot), \quad Y_t | (X_t = x) \sim g_\theta(\cdot | x) \quad \text{for } t = 1, \dots, T$$

where $(X_t)_{t \geq 1}$ are latent variables and $(Y_t)_{t \geq 1}$ correspond to observations.

- The likelihood of $Y_{1:T} = y_{1:T}$ for parameter $\theta \in \mathbb{R}^d$ is

$$p_\theta(y_{1:T}) = \prod_{t=1}^T p_\theta(y_t), \quad \text{where } p_\theta(y_t) = \int \mu_\theta(x_t) g_\theta(y_t | x_t) dx_t.$$

- Assume

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mu_\theta(\cdot), \quad Y_t | (X_t = x) \sim g_\theta(\cdot | x) \quad \text{for } t = 1, \dots, T$$

where $(X_t)_{t \geq 1}$ are latent variables and $(Y_t)_{t \geq 1}$ correspond to observations.

- The likelihood of $Y_{1:T} = y_{1:T}$ for parameter $\theta \in \mathbb{R}^d$ is

$$p_\theta(y_{1:T}) = \prod_{t=1}^T p_\theta(y_t), \quad \text{where } p_\theta(y_t) = \int \mu_\theta(x_t) g_\theta(y_t | x_t) dx_t.$$

- In many scenarios, $p_\theta(y_{1:T})$ cannot be evaluated exactly.

Example: Multivariate Probit model

- Multivariate latent Gaussian variables

$$X_t = Z_t\beta + \varepsilon_t, \quad \varepsilon_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, R).$$

Example: Multivariate Probit model

- Multivariate latent Gaussian variables

$$X_t = Z_t\beta + \varepsilon_t, \quad \varepsilon_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, R).$$

- Multivariate binary observations

$$Y_{ti} = \mathbb{I}(X_{ti} \geq 0), \quad i = 1, \dots, n$$

Example: Multivariate Probit model

- Multivariate latent Gaussian variables

$$X_t = Z_t\beta + \varepsilon_t, \quad \varepsilon_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, R).$$

- Multivariate binary observations

$$Y_{ti} = \mathbb{I}(X_{ti} \geq 0), \quad i = 1, \dots, n$$

- Likelihood of (β, R) is the product of T integrals of n -dimensional truncated multivariate normals.

State-Space Models

- Assume $\{X_t\}_{t \geq 1}$ is a latent Markov process, i.e. $X_1 \sim \mu_\theta(\cdot)$ and $X_{t+1} | (X_t = x) \sim f_\theta(\cdot | x)$, $Y_t | (X_t = x) \sim g_\theta(\cdot | x)$.

State-Space Models

- Assume $\{X_t\}_{t \geq 1}$ is a latent Markov process, i.e. $X_1 \sim \mu_\theta(\cdot)$ and

$$X_{t+1} | (X_t = x) \sim f_\theta(\cdot | x), \quad Y_t | (X_t = x) \sim g_\theta(\cdot | x).$$

- The likelihood of observations $Y_{1:T} = y_{1:T}$ is

$$p_\theta(y_{1:T}) = \int p_\theta(x_{1:T}, y_{1:T}) dx_{1:T}$$

where

$$p_\theta(x_{1:T}, y_{1:T}) = \mu_\theta(x_1) g_\theta(y_1 | x_1) \prod_{t=2}^T f_\theta(x_t | x_{t-1}) g_\theta(y_t | x_t).$$

State-Space Models

- Assume $\{X_t\}_{t \geq 1}$ is a latent Markov process, i.e. $X_1 \sim \mu_\theta(\cdot)$ and $X_{t+1} | (X_t = x) \sim f_\theta(\cdot | x)$, $Y_t | (X_t = x) \sim g_\theta(\cdot | x)$.

- The likelihood of observations $Y_{1:T} = y_{1:T}$ is

$$p_\theta(y_{1:T}) = \int p_\theta(x_{1:T}, y_{1:T}) dx_{1:T}$$

where

$$p_\theta(x_{1:T}, y_{1:T}) = \mu_\theta(x_1) g_\theta(y_1 | x_1) \prod_{t=2}^T f_\theta(x_t | x_{t-1}) g_\theta(y_t | x_t).$$

- State-space models are ubiquitous in time series analysis but inference is difficult as $p_\theta(y_{1:T})$ is intractable for non-linear/non-Gaussian models.

- Two species X_s^1 (prey) and X_s^2 (predator)

$$\Pr \left(X_{s+ds}^1 = x_s^1 + 1, X_{s+ds}^2 = x_s^2 \mid x_s^1, x_s^2 \right) = \alpha x_s^1 ds + o(ds),$$

$$\Pr \left(X_{s+ds}^1 = x_s^1 - 1, X_{s+ds}^2 = x_s^2 + 1 \mid x_s^1, x_s^2 \right) = \beta x_s^1 x_s^2 ds + o(ds),$$

$$\Pr \left(X_{s+ds}^1 = x_s^1, X_{s+ds}^2 = x_s^2 - 1 \mid x_s^1, x_s^2 \right) = \gamma x_s^2 ds + o(ds),$$

observed at discrete times

$$Y_t = X_{\Delta t}^1 + W_t \text{ with } W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

- Two species X_s^1 (prey) and X_s^2 (predator)

$$\Pr (X_{s+ds}^1 = x_s^1 + 1, X_{s+ds}^2 = x_s^2 \mid x_s^1, x_s^2) = \alpha x_s^1 ds + o(ds),$$

$$\Pr (X_{s+ds}^1 = x_s^1 - 1, X_{s+ds}^2 = x_s^2 + 1 \mid x_s^1, x_s^2) = \beta x_s^1 x_s^2 ds + o(ds),$$

$$\Pr (X_{s+ds}^1 = x_s^1, X_{s+ds}^2 = x_s^2 - 1 \mid x_s^1, x_s^2) = \gamma x_s^2 ds + o(ds),$$

observed at discrete times

$$Y_t = X_{\Delta t}^1 + W_t \text{ with } W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

- Kinetic rate constants $\theta = (\alpha, \beta, \gamma)$.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.
- **Epidemiology**: disease dynamic models.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.
- **Epidemiology**: disease dynamic models.
- **Ecology**: population dynamics.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.
- **Epidemiology**: disease dynamic models.
- **Ecology**: population dynamics.
- **Environmentrics**: phytoplankton-zooplankton model, paleoclimate reconstruction.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.
- **Epidemiology**: disease dynamic models.
- **Ecology**: population dynamics.
- **Environmentrics**: phytoplankton-zooplankton model, paleoclimate reconstruction.
- **Macroeconomics**: dynamic generalized stochastic equilibrium.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.
- **Epidemiology**: disease dynamic models.
- **Ecology**: population dynamics.
- **Environmentrics**: phytoplankton-zooplankton model, paleoclimate reconstruction.
- **Macroeconomics**: dynamic generalized stochastic equilibrium.
- **Signal Processing**: target tracking.

- State-space models are ubiquitous: 16,700 hits on Google Scholar since January 2015.
- **Econometrics**: stochastic volatility models.
- **Epidemiology**: disease dynamic models.
- **Ecology**: population dynamics.
- **Environmentrics**: phytoplankton-zooplankton model, paleoclimate reconstruction.
- **Macroeconomics**: dynamic generalized stochastic equilibrium.
- **Signal Processing**: target tracking.
- **Systems biology**: stochastic kinetic models.

Bayesian Inference for Latent Variable Models

- Prior distribution of density $p(\theta)$.

Bayesian Inference for Latent Variable Models

- Prior distribution of density $p(\theta)$.
- Likelihood function $p_{\theta}(y_{1:T})$.

Bayesian Inference for Latent Variable Models

- Prior distribution of density $p(\theta)$.
- Likelihood function $p_{\theta}(y_{1:T})$.
- Bayesian inference relies on the posterior

$$\pi(\theta) = p(\theta | y_{1:T}) = \frac{p_{\theta}(y_{1:T}) p(\theta)}{\int_{\Theta} p_{\theta'}(y_{1:T}) p(\theta') d\theta'}.$$

Bayesian Inference for Latent Variable Models

- Prior distribution of density $p(\theta)$.
- Likelihood function $p_{\theta}(y_{1:T})$.
- Bayesian inference relies on the posterior

$$\pi(\theta) = p(\theta | y_{1:T}) = \frac{p_{\theta}(y_{1:T}) p(\theta)}{\int_{\Theta} p_{\theta'}(y_{1:T}) p(\theta') d\theta'}.$$

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

Standard MCMC Approaches

- Standard MCMC schemes target $p(\theta, x_{1:T} | y_{1:T})$ where

$$p(\theta, x_{1:T} | y_{1:T}) \propto p(\theta) p_{\theta}(x_{1:T}, y_{1:T})$$

using Gibbs type strategy; i.e. sample alternately $X_{1:T} \sim p_{\theta}(\cdot | y_{1:T})$ and $\theta \sim p(\cdot | y_{1:T}, X_{1:T})$.

Standard MCMC Approaches

- Standard MCMC schemes target $p(\theta, x_{1:T} | y_{1:T})$ where

$$p(\theta, x_{1:T} | y_{1:T}) \propto p(\theta) p_{\theta}(x_{1:T}, y_{1:T})$$

using Gibbs type strategy; i.e. sample alternately $X_{1:T} \sim p_{\theta}(\cdot | y_{1:T})$ and $\theta \sim p(\cdot | y_{1:T}, X_{1:T})$.

- **Problem 1:** it can be difficult to sample $p_{\theta}(x_{1:T} | y_{1:T})$; e.g. state-space models.

Standard MCMC Approaches

- Standard MCMC schemes target $p(\theta, x_{1:T} | y_{1:T})$ where

$$p(\theta, x_{1:T} | y_{1:T}) \propto p(\theta) p_{\theta}(x_{1:T}, y_{1:T})$$

using Gibbs type strategy; i.e. sample alternately $X_{1:T} \sim p_{\theta}(\cdot | y_{1:T})$ and $\theta \sim p(\cdot | y_{1:T}, X_{1:T})$.

- **Problem 1:** it can be difficult to sample $p_{\theta}(x_{1:T} | y_{1:T})$; e.g. state-space models.
- **Problem 2:** Even when it is implementable, Gibbs can converge very slowly.

Standard MCMC Approaches

- Standard MCMC schemes target $p(\theta, x_{1:T} | y_{1:T})$ where

$$p(\theta, x_{1:T} | y_{1:T}) \propto p(\theta) p_{\theta}(x_{1:T}, y_{1:T})$$

using Gibbs type strategy; i.e. sample alternately $X_{1:T} \sim p_{\theta}(\cdot | y_{1:T})$ and $\theta \sim p(\cdot | y_{1:T}, X_{1:T})$.

- **Problem 1:** it can be difficult to sample $p_{\theta}(x_{1:T} | y_{1:T})$; e.g. state-space models.
- **Problem 2:** Even when it is implementable, Gibbs can converge very slowly.
- Pseudo-marginal methods mimick an algorithm targetting directly $p(\theta | y_{1:T})$ instead of $p(\theta, x_{1:T} | y_{1:T})$.

Ideal Marginal Metropolis-Hastings algorithm

- Metropolis–Hastings (MH) algorithm simulates an ergodic Markov chain $\{\theta_i\}_{i \geq 1}$ of limiting distribution $\pi(\theta)$.

Ideal Marginal Metropolis-Hastings algorithm

- Metropolis–Hastings (MH) algorithm simulates an ergodic Markov chain $\{\theta_i\}_{i \geq 1}$ of limiting distribution $\pi(\theta)$.

Ideal Marginal Metropolis-Hastings algorithm

- Metropolis–Hastings (MH) algorithm simulates an ergodic Markov chain $\{\vartheta_i\}_{i \geq 1}$ of limiting distribution $\pi(\theta)$.

At iteration i

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

Ideal Marginal Metropolis-Hastings algorithm

- Metropolis–Hastings (MH) algorithm simulates an ergodic Markov chain $\{\vartheta_i\}_{i \geq 1}$ of limiting distribution $\pi(\theta)$.

At iteration i

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

$$\min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} \right\} = \min \left\{ 1, \frac{p_{\vartheta}(y_{1:T}) p(\vartheta)}{p_{\vartheta_{i-1}}(y_{1:T}) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} \right\},$$

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

Ideal Marginal Metropolis-Hastings algorithm

- Metropolis–Hastings (MH) algorithm simulates an ergodic Markov chain $\{\vartheta_i\}_{i \geq 1}$ of limiting distribution $\pi(\theta)$.

At iteration i

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

$$\min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} \right\} = \min \left\{ 1, \frac{p_{\vartheta}(y_{1:T}) p(\vartheta)}{p_{\vartheta_{i-1}}(y_{1:T}) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} \right\},$$

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

- **Problem:** MH cannot be implemented if $p_{\vartheta}(y_{1:T})$ cannot be evaluated.

Pseudo-Marginal Metropolis–Hastings algorithm

- **“Idea”**: Replace $p_\theta(y_{1:T})$ by an estimate $\hat{p}_\theta(y_{1:T})$ in MH.

Pseudo-Marginal Metropolis–Hastings algorithm

- **“Idea”**: Replace $p_\theta(y_{1:T})$ by an estimate $\hat{p}_\theta(y_{1:T})$ in MH.

Pseudo-Marginal Metropolis–Hastings algorithm

- **“Idea”**: Replace $p_{\vartheta}(y_{1:T})$ by an estimate $\hat{p}_{\vartheta}(y_{1:T})$ in MH.

At iteration i

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

Pseudo-Marginal Metropolis–Hastings algorithm

- **“Idea”**: Replace $p_{\vartheta}(y_{1:T})$ by an estimate $\hat{p}_{\vartheta}(y_{1:T})$ in MH.

At iteration i

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

Pseudo-Marginal Metropolis–Hastings algorithm

- **“Idea”**: Replace $p_\theta (y_{1:T})$ by an estimate $\hat{p}_\theta (y_{1:T})$ in MH.

At iteration i

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

$$\min \left\{ 1, \underbrace{\frac{p_\vartheta (y_{1:T})}{p_{\vartheta_{i-1}} (y_{1:T})} \frac{p (\vartheta)}{p (\vartheta_{i-1})} \frac{q (\vartheta_{i-1} | \vartheta)}{q (\vartheta | \vartheta_{i-1})}}_{\text{exact MH ratio}} \times \underbrace{\frac{\hat{p}_\vartheta (y_{1:T}) / p_\vartheta (y_{1:T})}{\hat{p}_{\vartheta_{i-1}} (y_{1:T}) / p_{\vartheta_{i-1}} (y_{1:T})}}_{\text{noise}} \right\}$$
$$= \min \left\{ 1, \frac{\hat{p}_\vartheta (y_{1:T}) p (\vartheta)}{\hat{p}_{\vartheta_{i-1}} (y_{1:T}) p (\vartheta_{i-1})} \frac{q (\vartheta_{i-1} | \vartheta)}{q (\vartheta | \vartheta_{i-1})} \right\}$$

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

Key Result

- **Proposition** (Lin, Liu & Sloan, 2000; Andrieu & Roberts, 2009): If $\hat{p}_\theta (y_{1:T})$ is a non-negative unbiased estimator of $p_\theta (y_{1:T})$ then the pseudo-marginal MH kernel admits $\pi (\theta)$ as invariant density.

Key Result

- **Proposition** (Lin, Liu & Sloan, 2000; Andrieu & Roberts, 2009): If $\hat{p}_\theta(y_{1:T})$ is a non-negative unbiased estimator of $p_\theta(y_{1:T})$ then the pseudo-marginal MH kernel admits $\pi(\theta)$ as invariant density.
- Let U be the r.v. such that $\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; U)$ and $\mathbb{E}[\hat{p}_\theta(y_{1:T}; U)] = p_\theta(y_{1:T})$ when $U \sim m(\cdot)$.

Key Result

- **Proposition** (Lin, Liu & Sloan, 2000; Andrieu & Roberts, 2009): If $\hat{p}_\theta(y_{1:T})$ is a non-negative unbiased estimator of $p_\theta(y_{1:T})$ then the pseudo-marginal MH kernel admits $\pi(\theta)$ as invariant density.
- Let U be the r.v. such that $\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; U)$ and $\mathbb{E}[\hat{p}_\theta(y_{1:T}; U)] = p_\theta(y_{1:T})$ when $U \sim m(\cdot)$.
- Consider the auxiliary target density on $\Theta \times \mathcal{U}$

$$\bar{\pi}(\theta, u) = \pi(\theta) \underbrace{\frac{\hat{p}_\theta(y_{1:T}; u)}{p_\theta(y_{1:T})}}_{\int(\cdot)du=1} m(u)$$

Key Result

- **Proposition** (Lin, Liu & Sloan, 2000; Andrieu & Roberts, 2009): If $\hat{p}_\theta(y_{1:T})$ is a non-negative unbiased estimator of $p_\theta(y_{1:T})$ then the pseudo-marginal MH kernel admits $\pi(\theta)$ as invariant density.
- Let U be the r.v. such that $\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; U)$ and $\mathbb{E}[\hat{p}_\theta(y_{1:T}; U)] = p_\theta(y_{1:T})$ when $U \sim m(\cdot)$.
- Consider the auxiliary target density on $\Theta \times \mathcal{U}$

$$\bar{\pi}(\theta, u) = \pi(\theta) \underbrace{\frac{\hat{p}_\theta(y_{1:T}; u)}{p_\theta(y_{1:T})} m(u)}_{\int(\cdot)du=1}$$

- Pseudo-marginal MH is a standard MH with target $\bar{\pi}(\theta, u)$ and proposal $q(\vartheta|\theta) m(v)$ as

$$\frac{\bar{\pi}(\vartheta, v) q(\theta|\vartheta) m(u)}{\bar{\pi}(\theta, u) q(\vartheta|\theta) m(v)} = \frac{\hat{p}_\theta(y_{1:T}; v) p(\vartheta) q(\theta|\vartheta)}{\hat{p}_\theta(y_{1:T}; u) p(\theta) q(\vartheta|\theta)}$$

Importance Sampling Estimator

- For latent variable models, one has

$$p_{\theta}(y_t) = \int \mu_{\theta}(x_t) g_{\theta}(y_t | x_t) dx_t.$$

Importance Sampling Estimator

- For latent variable models, one has

$$p_{\theta}(y_t) = \int \mu_{\theta}(x_t) g_{\theta}(y_t | x_t) dx_t.$$

- An non-negative unbiased estimator is given by

$$\hat{p}_{\theta}(y_{1:T}) = \prod_{t=1}^T \hat{p}_{\theta}(y_t) = \prod_{t=1}^T \left\{ \frac{1}{N} \sum_{k=1}^N g_{\theta}(y_t | x_t^k) \right\}, \quad x_t^k \stackrel{\text{i.i.d.}}{\sim} \mu_{\theta},$$

i.e.

$$m(u) = \prod_{t=1}^T \prod_{k=1}^N \mu_{\theta}(x_t^k).$$

Importance Sampling Estimator

- For latent variable models, one has

$$p_{\theta}(y_t) = \int \mu_{\theta}(x_t) g_{\theta}(y_t | x_t) dx_t.$$

- An non-negative unbiased estimator is given by

$$\hat{p}_{\theta}(y_{1:T}) = \prod_{t=1}^T \hat{p}_{\theta}(y_t) = \prod_{t=1}^T \left\{ \frac{1}{N} \sum_{k=1}^N g_{\theta}(y_t | x_t^k) \right\}, \quad x_t^k \stackrel{\text{i.i.d.}}{\sim} \mu_{\theta},$$

i.e.

$$m(u) = \prod_{t=1}^T \prod_{k=1}^N \mu_{\theta}(x_t^k).$$

- Computational complexity is $O(NT)$.

Particle Filter Estimator

- For state-space models, previous approach provides an estimator whose relative variance scales typically exponentially with T .

Particle Filter Estimator

- For state-space models, previous approach provides an estimator whose relative variance scales typically exponentially with T .
- An alternative is to use particle filter where

$$\begin{aligned}\widehat{p}_\theta(y_{1:T}) &= \widehat{p}_\theta(y_1) \prod_{t=2}^T \widehat{p}_\theta(y_t | y_{1:t-1}) \\ &= \prod_{t=1}^T \left\{ \frac{1}{N} \sum_{k=1}^N g_\theta(y_t | X_n^k) \right\}\end{aligned}$$

where

$$m(u) = \prod_{k=1}^N \mu_\theta(x_1^k) \prod_{t=2}^T \left\{ \prod_{k=1}^N w_t^{a_{t-1}^k} f(x_t^k | x_{t-1}^{a_{t-1}^k}) \right\}$$

with $a_{t-1}^k \in \{1, \dots, N\}$, $w_t^j \propto g_\theta(y_t | X_t^j)$, $\sum_j w_t^j = 1$.

Particle Filter Estimator

- For state-space models, previous approach provides an estimator whose relative variance scales typically exponentially with T .
- An alternative is to use particle filter where

$$\begin{aligned}\widehat{p}_\theta(y_{1:T}) &= \widehat{p}_\theta(y_1) \prod_{t=2}^T \widehat{p}_\theta(y_t | y_{1:t-1}) \\ &= \prod_{t=1}^T \left\{ \frac{1}{N} \sum_{k=1}^N g_\theta(y_t | X_n^k) \right\}\end{aligned}$$

where

$$m(u) = \prod_{k=1}^N \mu_\theta(x_1^k) \prod_{t=2}^T \left\{ \prod_{k=1}^N w_t^{a_{t-1}^k} f(x_t^k | x_{t-1}^{a_{t-1}^k}) \right\}$$

with $a_{t-1}^k \in \{1, \dots, N\}$, $w_t^j \propto g_\theta(y_t | X_t^j)$, $\sum_j w_t^j = 1$.

- Computational complexity is $O(NT)$.

Particle Filter Estimator

- For state-space models, previous approach provides an estimator whose relative variance scales typically exponentially with T .
- An alternative is to use particle filter where

$$\begin{aligned}\hat{p}_\theta(y_{1:T}) &= \hat{p}_\theta(y_1) \prod_{t=2}^T \hat{p}_\theta(y_t | y_{1:t-1}) \\ &= \prod_{t=1}^T \left\{ \frac{1}{N} \sum_{k=1}^N g_\theta(y_t | X_n^k) \right\}\end{aligned}$$

where

$$m(u) = \prod_{k=1}^N \mu_\theta(x_1^k) \prod_{t=2}^T \left\{ \prod_{k=1}^N w_t^{a_{t-1}^k} f(x_t^k | x_{t-1}^{a_{t-1}^k}) \right\}$$

with $a_{t-1}^k \in \{1, \dots, N\}$, $w_t^j \propto g_\theta(y_t | X_t^j)$, $\sum_j w_t^j = 1$.

- Computational complexity is $O(NT)$.
- The estimator $\hat{p}_\theta(y_{1:T})$ of $p_\theta(y_{1:T})$ is unbiased and its relative variance is bounded uniformly over T if $N \propto T$ (Cerou, Del Moral &

Pseudo-Marginal Metropolis–Hastings algorithm

Pseudo-Marginal Metropolis–Hastings algorithm

At iteration i

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

Pseudo-Marginal Metropolis–Hastings algorithm

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- Use particle filter to compute an estimate $\hat{p}_\vartheta(y_{1:T})$ of $p_\vartheta(y_{1:T})$.

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$.
- Use particle filter to compute an estimate $\hat{p}_\vartheta(y_{1:T})$ of $p_\vartheta(y_{1:T})$.
- With probability

$$\min\left\{1, \frac{\hat{p}_\vartheta(y_{1:T}) p(\vartheta)}{\hat{p}_{\vartheta_{i-1}}(y_{1:T}) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})}\right\}$$

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

Empirical performance: Stochastic kinetic model

- Two species X_s^1 (prey) and X_s^2 (predator)

$$\Pr \left(X_{s+ds}^1 = x_s^1 + 1, X_{s+ds}^2 = x_s^2 \mid x_s^1, x_s^2 \right) = \alpha x_s^1 ds + o(ds),$$

$$\Pr \left(X_{s+ds}^1 = x_s^1 - 1, X_{s+ds}^2 = x_s^2 + 1 \mid x_s^1, x_s^2 \right) = \beta x_s^1 x_s^2 ds + o(ds),$$

$$\Pr \left(X_{s+ds}^1 = x_s^1, X_{s+ds}^2 = x_s^2 - 1 \mid x_s^1, x_s^2 \right) = \gamma x_s^2 ds + o(ds),$$

observed at discrete times

$$Y_t = X_{\Delta t}^1 + W_t \text{ with } W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

Empirical performance: Stochastic kinetic model

- Two species X_s^1 (prey) and X_s^2 (predator)

$$\Pr \left(X_{s+ds}^1 = x_s^1 + 1, X_{s+ds}^2 = x_s^2 \mid x_s^1, x_s^2 \right) = \alpha x_s^1 ds + o(ds),$$

$$\Pr \left(X_{s+ds}^1 = x_s^1 - 1, X_{s+ds}^2 = x_s^2 + 1 \mid x_s^1, x_s^2 \right) = \beta x_s^1 x_s^2 ds + o(ds),$$

$$\Pr \left(X_{s+ds}^1 = x_s^1, X_{s+ds}^2 = x_s^2 - 1 \mid x_s^1, x_s^2 \right) = \gamma x_s^2 ds + o(ds),$$

observed at discrete times

$$Y_t = X_{\Delta t}^1 + W_t \text{ with } W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

- We are interested in the kinetic rate constants $\theta = (\alpha, \beta, \gamma)$ a priori distributed as (Boys et al., 2008; Kunsch, 2011)

$$\alpha \sim \mathcal{G}(1, 10), \quad \beta \sim \mathcal{G}(1, 0.25), \quad \gamma \sim \mathcal{G}(1, 7.5).$$

Empirical performance: Stochastic kinetic model

- Two species X_s^1 (prey) and X_s^2 (predator)

$$\begin{aligned}\Pr\left(X_{s+ds}^1 = x_s^1 + 1, X_{s+ds}^2 = x_s^2 \mid x_s^1, x_s^2\right) &= \alpha x_s^1 ds + o(ds), \\ \Pr\left(X_{s+ds}^1 = x_s^1 - 1, X_{s+ds}^2 = x_s^2 + 1 \mid x_s^1, x_s^2\right) &= \beta x_s^1 x_s^2 ds + o(ds), \\ \Pr\left(X_{s+ds}^1 = x_s^1, X_{s+ds}^2 = x_s^2 - 1 \mid x_s^1, x_s^2\right) &= \gamma x_s^2 ds + o(ds),\end{aligned}$$

observed at discrete times

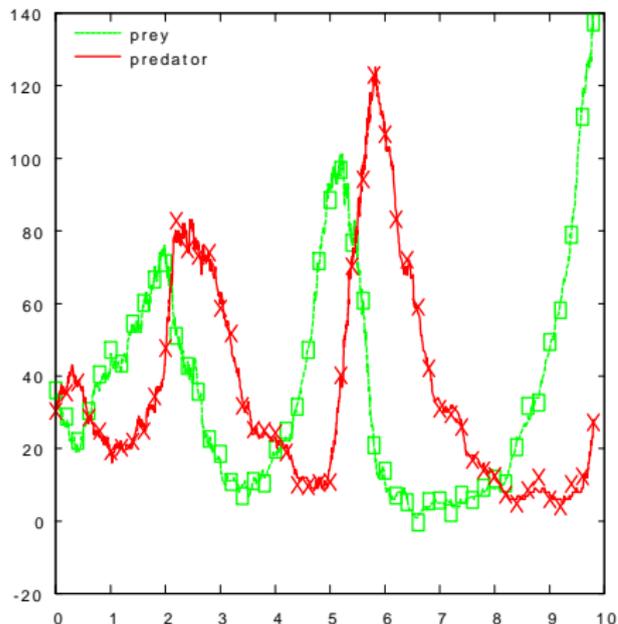
$$Y_t = X_{\Delta t}^1 + W_t \text{ with } W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

- We are interested in the kinetic rate constants $\theta = (\alpha, \beta, \gamma)$ a priori distributed as (Boys et al., 2008; Kunsch, 2011)

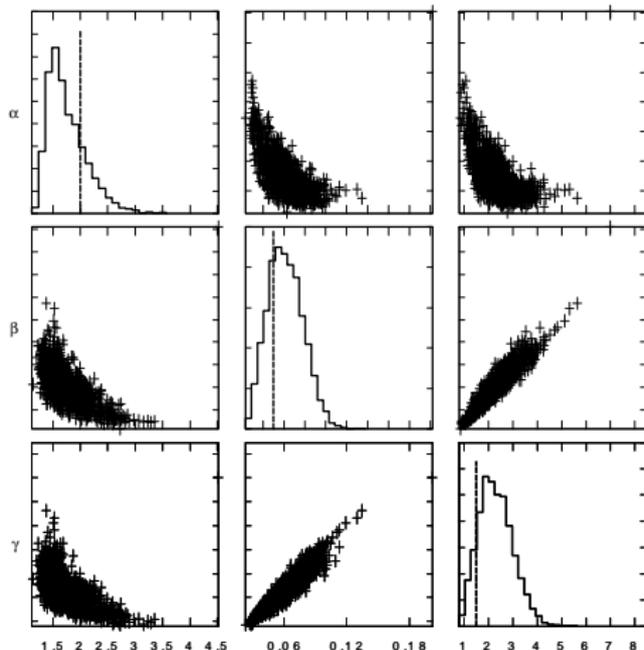
$$\alpha \sim \mathcal{G}(1, 10), \quad \beta \sim \mathcal{G}(1, 0.25), \quad \gamma \sim \mathcal{G}(1, 7.5).$$

- Pseudo-marginal MH with RW proposal, likelihood is approximated using particle filter.

Empirical performance: Stochastic kinetic model

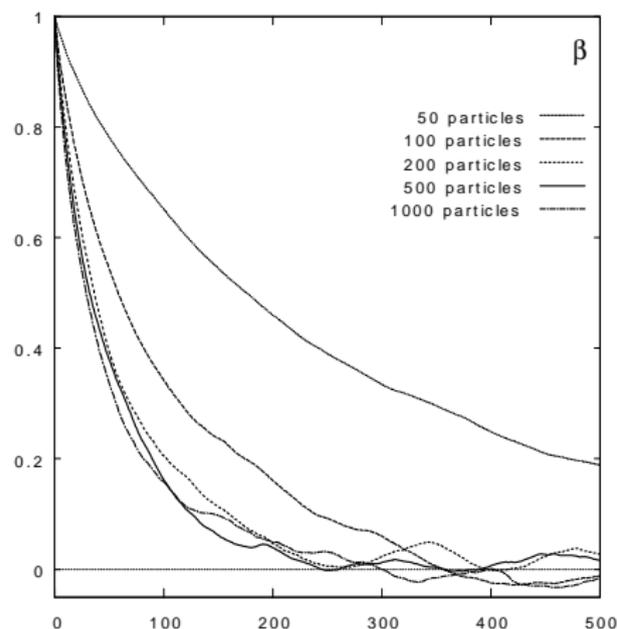
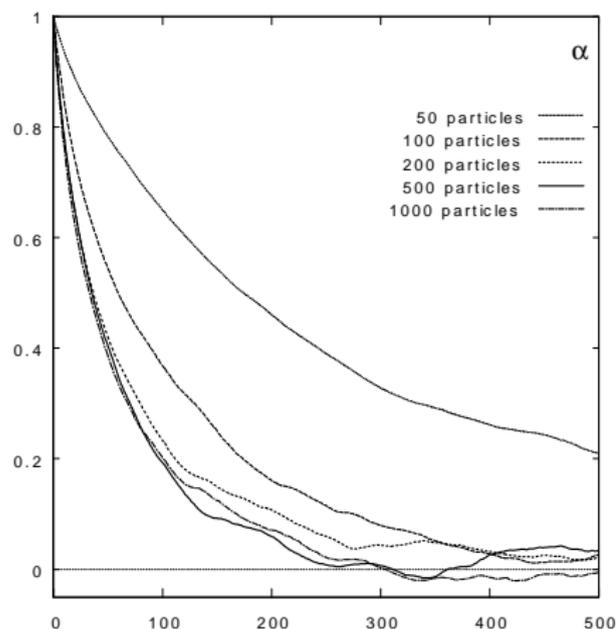


Simulated data



Estimated posteriors

Empirical performance: Stochastic kinetic model



Autocorrelation of α (left) and β (right) for the PM sampler for various N .

- Huang & Tauchen, *J. Financial Econometrics* (2005):

$$dv_1(s) = -k_1 \{v_1(s) - \mu_1\} ds + \sigma_1 dW_1(s),$$

$$dv_2(s) = -k_2 v_2(s) ds + \{1 + \beta_{12} v_2(s)\} dW_2(s),$$

$$d \log P(s) = \mu_y ds + s \cdot \exp \left[\frac{v_1(s) + \beta_2 v_2(s)}{2} \right] dB(s),$$

with $\phi_1 = \text{corr}\{B(s), W_1(s)\}$ and $\phi_2 = \text{corr}\{B(s), W_2(s)\}$.

- Huang & Tauchen, *J. Financial Econometrics* (2005):

$$\begin{aligned}dv_1(s) &= -k_1 \{v_1(s) - \mu_1\} ds + \sigma_1 dW_1(s), \\dv_2(s) &= -k_2 v_2(s) ds + \{1 + \beta_{12} v_2(s)\} dW_2(s), \\d \log P(s) &= \mu_y ds + s \cdot \exp \left[\frac{v_1(s) + \beta_2 v_2(s)}{2} \right] dB(s),\end{aligned}$$

with $\phi_1 = \text{corr}\{B(s), W_1(s)\}$ and $\phi_2 = \text{corr}\{B(s), W_2(s)\}$.

- Euler discretization of the volatilities $v_1(s)$ and $v_2(s)$ provides closed form expression for $Y_t = \log P(\Delta t) - \log P(\Delta(t-1))$.

- Huang & Tauchen, *J. Financial Econometrics* (2005):

$$dv_1(s) = -k_1 \{v_1(s) - \mu_1\} ds + \sigma_1 dW_1(s),$$

$$dv_2(s) = -k_2 v_2(s) ds + \{1 + \beta_{12} v_2(s)\} dW_2(s),$$

$$d \log P(s) = \mu_y ds + s \cdot \exp \left[\frac{v_1(s) + \beta_2 v_2(s)}{2} \right] dB(s),$$

with $\phi_1 = \text{corr}\{B(s), W_1(s)\}$ and $\phi_2 = \text{corr}\{B(s), W_2(s)\}$.

- Euler discretization of the volatilities $v_1(s)$ and $v_2(s)$ provides closed form expression for $Y_t = \log P(\Delta t) - \log P(\Delta(t-1))$.
- Daily returns $y = (y_1, \dots, y_T)$ of the S&P 500 index.

- Huang & Tauchen, *J. Financial Econometrics* (2005):

$$dv_1(s) = -k_1 \{v_1(s) - \mu_1\} ds + \sigma_1 dW_1(s),$$

$$dv_2(s) = -k_2 v_2(s) ds + \{1 + \beta_{12} v_2(s)\} dW_2(s),$$

$$d \log P(s) = \mu_y ds + s \cdot \exp[\{v_1(s) + \beta_2 v_2(s)\} / 2] dB(s),$$

with $\phi_1 = \text{corr}\{B(s), W_1(s)\}$ and $\phi_2 = \text{corr}\{B(s), W_2(s)\}$.

- Euler discretization of the volatilities $v_1(s)$ and $v_2(s)$ provides closed form expression for $Y_t = \log P(\Delta t) - \log P(\Delta(t-1))$.
- Daily returns $y = (y_1, \dots, y_T)$ of the S&P 500 index.
- Bayesian Inference on $\theta = (k_1, \mu_1, \sigma_1, k_2, \beta_{12}, \beta_2, \mu_y, \phi_1, \phi_2)$.

- Huang & Tauchen, *J. Financial Econometrics* (2005):

$$\begin{aligned}dv_1(s) &= -k_1 \{v_1(s) - \mu_1\} ds + \sigma_1 dW_1(s), \\dv_2(s) &= -k_2 v_2(s) ds + \{1 + \beta_{12} v_2(s)\} dW_2(s), \\d \log P(s) &= \mu_y ds + s \cdot \exp[\{v_1(s) + \beta_2 v_2(s)\} / 2] dB(s),\end{aligned}$$

with $\phi_1 = \text{corr}\{B(s), W_1(s)\}$ and $\phi_2 = \text{corr}\{B(s), W_2(s)\}$.

- Euler discretization of the volatilities $v_1(s)$ and $v_2(s)$ provides closed form expression for $Y_t = \log P(\Delta t) - \log P(\Delta(t-1))$.
- Daily returns $y = (y_1, \dots, y_T)$ of the S&P 500 index.
- Bayesian Inference on $\theta = (k_1, \mu_1, \sigma_1, k_2, \beta_{12}, \beta_2, \mu_y, \phi_1, \phi_2)$.
- Performance of the pseudo-marginal for RW proposal w.r.t σ , standard deviation of $\log \hat{p}_\theta(y)$ at posterior mean $\bar{\theta}$.

Integrated Autocorrelation Time of Pseudo-Marginal MH

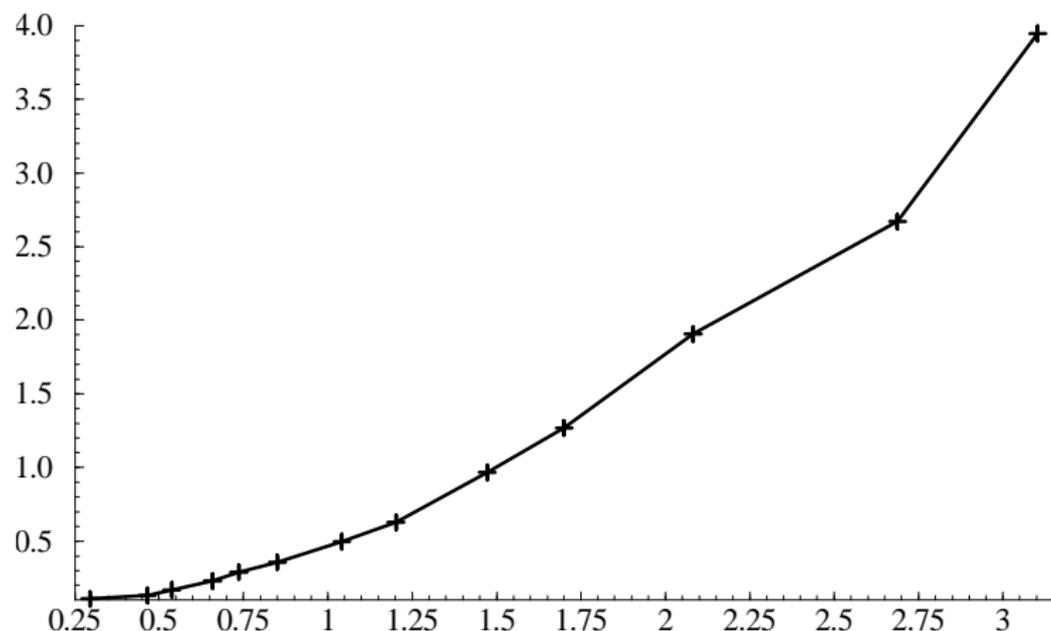


Figure: Average over the 9 parameter components of the log-integrated autocorrelation time of pseudo-marginal chain as a function of σ for $T = 300$.

How precise should the log-likelihood estimator be?

- **Aim:** Minimize the computational time

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

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

IF_h^Q = Integrated Autocorrelation Time of PM average

How precise should the log-likelihood estimator be?

- **Aim:** Minimize the computational time

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

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

IF_h^Q = Integrated Autocorrelation Time of PM average

- Call the IACT the *inefficiency*

$$IF_h^Q = 1 + 2 \sum_{\tau=1}^{\infty} \text{corr}_{\bar{\pi}, Q} \{h(\theta_0), h(\theta_\tau)\}$$

where Q is the pseudo-marginal kernel given for $(\theta, z) \neq (\vartheta, w)$ by

$$Q \{(\theta, z), (d\vartheta, dw)\} = q(\vartheta|\theta) g_\vartheta(w) \min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} \exp(w - z) \right\} d\vartheta dw,$$

where

$$z = \log \{ \hat{p}_\theta(y_{1:T}) / p_\theta(y_{1:T}) \},$$

$$w = \log \{ \hat{p}_\vartheta(y_{1:T}) / p_\vartheta(y_{1:T}) \}.$$

Computational time for the SV model

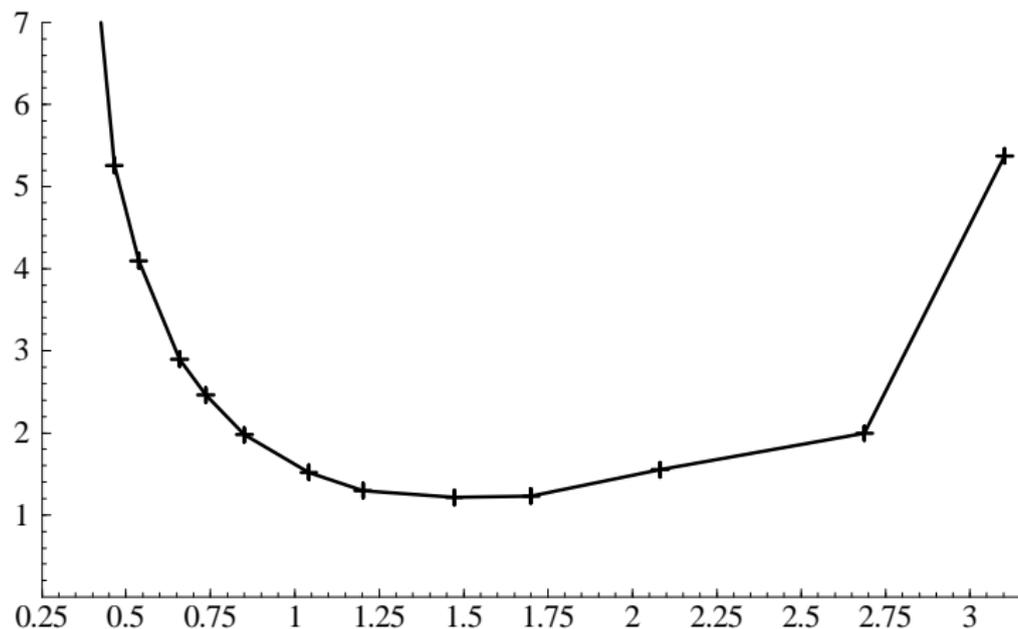


Figure: Computational time as a function of σ

Analysis in the large data regime

- Standard asymptotic study of MCMC relies on $d \rightarrow \infty$ and independence assumption on the target, interested here in fixed d , large T .

Analysis in the large data regime

- Standard asymptotic study of MCMC relies on $d \rightarrow \infty$ and independence assumption on the target, interested here in fixed d , large T .
- **Assumption 1** - Asymptotic Normality: We have

$$\int \left| p(\theta | Y_{1:T}) - \phi(\theta; \hat{\theta}^T, \Sigma/T) \right| d\theta \xrightarrow{P} 0,$$

where $\hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ is a p.d. matrix.

Analysis in the large data regime

- Standard asymptotic study of MCMC relies on $d \rightarrow \infty$ and independence assumption on the target, interested here in fixed d , large T .
- **Assumption 1** - Asymptotic Normality: We have

$$\int \left| p(\theta | Y_{1:T}) - \phi(\theta; \hat{\theta}^T, \Sigma / T) \right| d\theta \xrightarrow{P} 0,$$

where $\hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ is a p.d. matrix.

- **Assumption 2** - CLT: For any θ in a neighbourhood of $\bar{\theta}$,

$$\log \frac{\hat{p}_\theta(Y_{1:T})}{p_\theta(Y_{1:T})} \Big| \mathcal{Y}^T \Rightarrow \mathcal{N}(-\sigma^2(\theta) / 2, \sigma^2(\theta))$$

in probability and $\sigma^2(\cdot)$ continuous at $\bar{\theta}$.

Analysis in the large data regime

- Standard asymptotic study of MCMC relies on $d \rightarrow \infty$ and independence assumption on the target, interested here in fixed d , large T .

- **Assumption 1** - Asymptotic Normality: We have

$$\int \left| p(\theta | Y_{1:T}) - \phi(\theta; \hat{\theta}^T, \Sigma / T) \right| d\theta \xrightarrow{P} 0,$$

where $\hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ is a p.d. matrix.

- **Assumption 2** - CLT: For any θ in a neighbourhood of $\bar{\theta}$,

$$\log \frac{\hat{p}_\theta(Y_{1:T})}{p_\theta(Y_{1:T})} \Big| \mathcal{Y}^T \Rightarrow \mathcal{N}(-\sigma^2(\theta) / 2, \sigma^2(\theta))$$

in probability and $\sigma^2(\cdot)$ continuous at $\bar{\theta}$.

- **Assumption 3** - Proposal: $\vartheta = \theta + \varepsilon / \sqrt{T}$ where $\varepsilon \sim v(\cdot)$ with $v(\varepsilon) = v(-\varepsilon)$.

- Assumption 1 holds if for example Bernstein-von Mises holds (in correctly specified/misspecified scenarios).

Analysis in the large data regime

- Assumption 1 holds if for example Bernstein-von Mises holds (in correctly specified/misspecified scenarios).
- Assumption 2 has been shown to hold under regularity assumptions if $N \propto T$ (Berard et al, 2014, Deligiannidis et al, 2015).

Analysis in the large data regime

- Assumption 1 holds if for example Bernstein-von Mises holds (in correctly specified/misspecified scenarios).
- Assumption 2 has been shown to hold under regularity assumptions if $N \propto T$ (Berard et al, 2014, Deligiannidis et al, 2015).
- Assumption 3 can be easily enforced.

Weak convergence

- Let $\{\theta_i^T, Z_i^T := \log \widehat{p}_{\theta_i^T}(Y_{1:T}) / p_{\theta_i^T}(Y_{1:T})\}_{i \geq 0}$ the stationary PM Markov chain of invariant density $p(\theta | Y_{1:T}) \exp(z) g_\theta^T(z)$.

Weak convergence

- Let $\{\vartheta_i^T, Z_i^T := \log \widehat{p}_{\vartheta_i^T}(Y_{1:T}) / p_{\vartheta_i^T}(Y_{1:T})\}_{i \geq 0}$ the stationary PM Markov chain of invariant density $p(\theta | Y_{1:T}) \exp(z) g_\theta^T(z)$.
- **Proposition** (Schmon et al, 2016): The F.D.D. of the rescaled sequence $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \widehat{\theta}_T), Z_i^T\}_{i \geq 0}$ converge weakly as $T \rightarrow \infty$ to those of a stationary Markov chain of invariant density $\phi(\tilde{\theta}; 0, \Sigma) \phi(z; -\sigma^2(\bar{\theta})/2, \sigma^2(\bar{\theta}))$ and kernel given by

$$\begin{aligned} \tilde{Q}\{(\tilde{\theta}, z), (d\tilde{\vartheta}, dw)\} &= v(\tilde{\vartheta} - \tilde{\theta}) \phi(w; -\sigma^2(\bar{\theta})/2, \sigma^2(\bar{\theta})) \\ &\quad \times \min \left\{ 1, \frac{\phi(\tilde{\vartheta}; 0, \Sigma)}{\phi(\tilde{\theta}; 0, \Sigma)} \exp(w - z) \right\} d\tilde{\vartheta} dw \end{aligned}$$

for $(\tilde{\theta}, z) \neq (\tilde{\vartheta}, w)$.

Weak convergence

- These results suggests that a simplified analysis of the PM chain can be performed by looking at

$$\begin{aligned}\widehat{Q}\{(\theta, z), (d\vartheta, dw)\} &= q(\vartheta|\theta) \phi(w; -\sigma^2/2, \sigma^2) \\ &\times \min\left\{1, \frac{\pi(\vartheta)}{\pi(\theta)} \exp(w - z)\right\} d\vartheta dw,\end{aligned}$$

where $\sigma^2 = \sigma^2(\bar{\theta})$.

- These results suggests that a simplified analysis of the PM chain can be performed by looking at

$$\begin{aligned}\widehat{Q}\{(\theta, z), (d\vartheta, dw)\} &= q(\vartheta|\theta)\phi(w; -\sigma^2/2, \sigma^2) \\ &\quad \times \min\left\{1, \frac{\pi(\vartheta)}{\pi(\theta)} \exp(w - z)\right\} d\vartheta dw,\end{aligned}$$

where $\sigma^2 = \sigma^2(\bar{\theta})$.

- It would be more satisfactory to show that

$$\left|IF_h^Q - IF_h^{\widehat{Q}}\right| \rightarrow 0$$

as $T \rightarrow \infty$. The analysis relies on (Andrieu & Vihola, 2015) and is much more involved.

Empirical vs Assumed Distributions for SV model

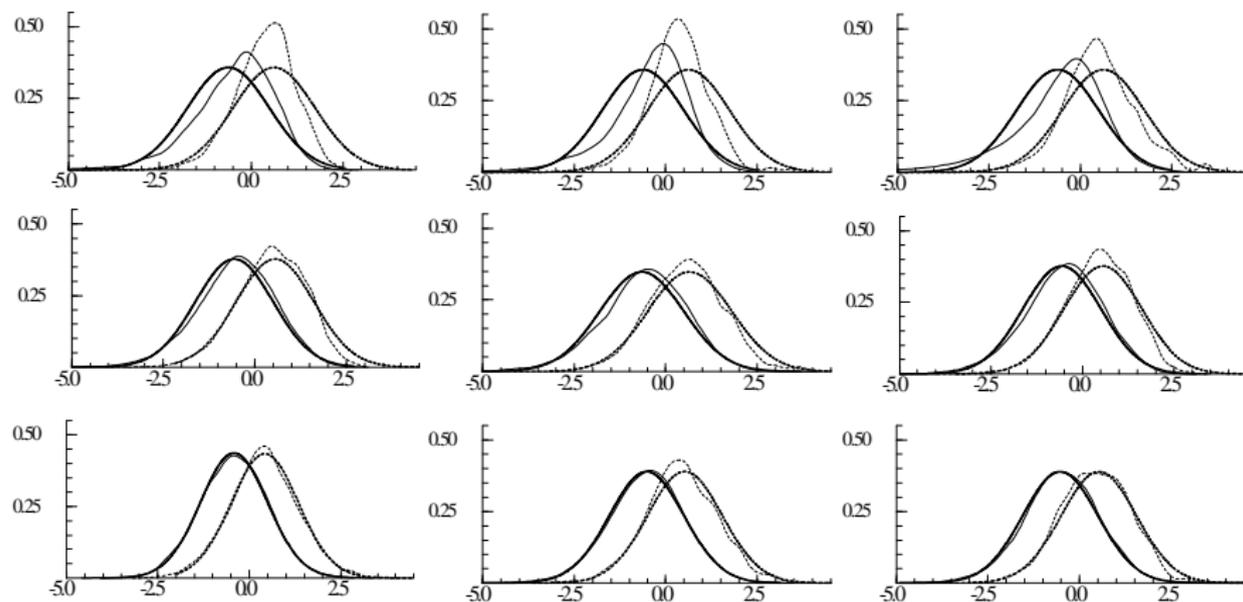


Figure: Empirical distributions (dashed) vs assumed Gaussians (solid) of Z at $\bar{\theta}$ (left) and marginalized over samples from $\pi(\theta)$ (center) and $\int \pi(d\theta) q(\theta | \bar{\theta})$ (right) for $T = 40$, $T = 300$ and $T = 2700$.

- **Aim:** Minimize the computational cost

$$CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2.$$

- **Aim:** Minimize the computational cost

$$CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2.$$

- **Special cases:**

- **Aim:** Minimize the computational cost

$$CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2.$$

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

- **Aim:** Minimize the computational cost

$$CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2.$$

- **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 $d \rightarrow \infty$, diffusion limit suggests $\sigma_{\text{opt}} = 1.81$ (Sherlock et al., 2015).

Sketch of the Analysis

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

Sketch of the Analysis

- For general proposals and targets, direct minimization of $CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2$ impossible so minimize an upper bound over it.
- Theoretical study relies on $\bar{\pi}$ -invariant kernel Q^* given for $(\theta, z) \neq (\vartheta, w)$ by

$$q(\vartheta|\theta)\phi(w; -\sigma^2/2, \sigma^2) \min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} \right\} \min \{1, \exp(w - z)\} d\vartheta dw,$$

instead of

$$q(\vartheta|\theta)\phi(w; -\sigma^2/2, \sigma^2) \min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} \exp(w - z) \right\} d\vartheta dw.$$

Sketch of the Analysis

- For general proposals and targets, direct minimization of $CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2$ impossible so minimize an upper bound over it.
- Theoretical study relies on $\bar{\pi}$ -invariant kernel Q^* given for $(\theta, z) \neq (\vartheta, w)$ by

$$q(\vartheta|\theta)\phi(w; -\sigma^2/2, \sigma^2) \min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} \right\} \min \{1, \exp(w - z)\} d\vartheta dw,$$

instead of

$$q(\vartheta|\theta)\phi(w; -\sigma^2/2, \sigma^2) \min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} \exp(w - z) \right\} d\vartheta dw.$$

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

Main Theoretical Result

- **Proposition:** If $IF_h^{Q^*}(\sigma) < \infty$ then $IF_h^{\hat{Q}}(\sigma) \leq IF_h^{Q^*}(\sigma)$ and

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

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

Main Theoretical Result

- **Proposition:** If $IF_h^{Q^*}(\sigma) < \infty$ then $IF_h^{\hat{Q}}(\sigma) \leq IF_h^{Q^*}(\sigma)$ and

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

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

- \tilde{Q}^{EX} and \tilde{Q}_σ^Z correspond to the jump kernels associated to Q^{EX} and Q_σ^Z , $\varrho_{\text{EX}}(\theta)$ and $\varrho_Z^\sigma(z)$ are acceptance proba of Q^{EX} and Q_σ^Z .

Main Theoretical Result

- **Proposition:** If $IF_h^{Q^*}(\sigma) < \infty$ then $IF_h^{\hat{Q}}(\sigma) \leq IF_h^{Q^*}(\sigma)$ and

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

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

- \tilde{Q}^{EX} and \tilde{Q}_σ^Z correspond to the jump kernels associated to Q^{EX} and Q_σ^Z , $\varrho_{\text{EX}}(\theta)$ and $\varrho_Z^\sigma(z)$ are acceptance proba of Q^{EX} and Q_σ^Z .
- This identity allows us to “decouple” the influence of the parameter and noise components on $IF_h^{Q^*}(\sigma)$.

Simpler Bounds on the Relative Inefficiency

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

$$\frac{IF_h^{\hat{Q}}(\sigma)}{IF_h^{EX}} \leq \frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \leq \frac{1}{2}(1 + 1/IF_h^{EX})\pi_Z^\sigma(1/q_Z^\sigma) - \frac{1}{IF_h^{EX}}$$

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

$$\frac{IF_h^{\hat{Q}}(\sigma)}{IF_h^{EX}} \leq \frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \leq \frac{1}{2}(1 + 1/IF_h^{EX})\pi_Z^\sigma(1/\varrho_Z^\sigma) - \frac{1}{IF_h^{EX}}$$

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

$$\frac{IF_h^{\hat{Q}}(\sigma)}{IF_h^{EX}} \leq \frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \leq \frac{1}{2}(1 + 1/IF_h^{EX})\pi_Z^\sigma(1/\varrho_Z^\sigma) - \frac{1}{IF_h^{EX}}$$

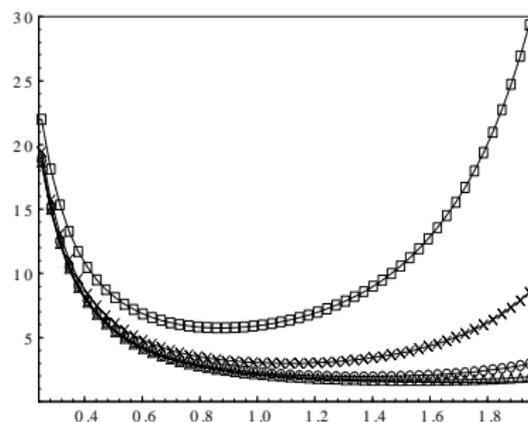
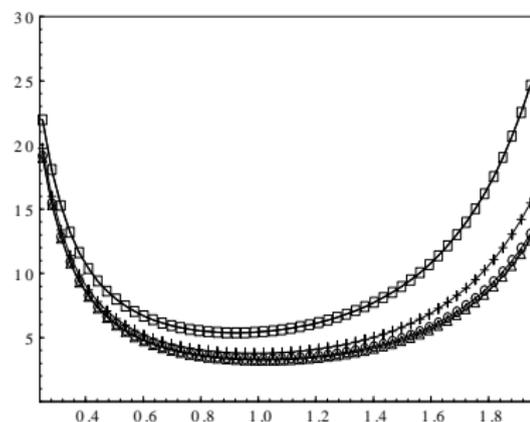
and the bound is tight as $IF_h^{EX} \rightarrow 1$ or $\sigma \rightarrow 0$.

- 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 $CT_h^{\hat{Q}}(\sigma) = IF_h^{\hat{Q}}(\sigma) / \sigma^2$.

Bounds on Relative Computational Time



Left: upper bound on $CT_h^{Q^*}(\sigma) / IF_h^{EX}$ as a function of σ for $IF_h^{EX} = 1$ (square), 4 (crosses), 20 (circles), 80 (triangles). Right: upper bounds on $CT_h^{Q^*}(\sigma) / IF_{J,h}^{EX}$ as a function of σ for $IF_{J,h}^{EX} / \rho_{EX} = 1$ for $IF_{J,h}^{EX} / \rho_{EX} = 1, 4, 20, 80$ and lower bound (solid line).

Practical Guidelines

- For good proposals, select $\sigma \approx 1.0$ whereas for poor proposals, select $\sigma \approx 1.7$.

Practical Guidelines

- For good proposals, select $\sigma \approx 1.0$ whereas for poor proposals, select $\sigma \approx 1.7$.
- When you have no clue about the proposal efficiency,

Practical Guidelines

- For good proposals, select $\sigma \approx 1.0$ whereas for poor proposals, select $\sigma \approx 1.7$.
- When you have no clue about the proposal efficiency,
- ① If $\sigma_{\text{opt}} = 1.0$ and you pick $\sigma = 1.7$, computing time increases by $\approx 150\%$.

Practical Guidelines

- For good proposals, select $\sigma \approx 1.0$ whereas for poor proposals, select $\sigma \approx 1.7$.
- When you have no clue about the proposal efficiency,
 - 1 If $\sigma_{\text{opt}} = 1.0$ and you pick $\sigma = 1.7$, computing time increases by $\approx 150\%$.
 - 2 If $\sigma_{\text{opt}} = 1.7$ and you pick $\sigma = 1.0$, computing time increases by $\approx 50\%$.

Practical Guidelines

- For good proposals, select $\sigma \approx 1.0$ whereas for poor proposals, select $\sigma \approx 1.7$.
- When you have no clue about the proposal efficiency,
 - 1 If $\sigma_{\text{opt}} = 1.0$ and you pick $\sigma = 1.7$, computing time increases by $\approx 150\%$.
 - 2 If $\sigma_{\text{opt}} = 1.7$ and you pick $\sigma = 1.0$, computing time increases by $\approx 50\%$.
 - 3 If $\sigma_{\text{opt}} = 1.0$ or $\sigma_{\text{opt}} = 1.7$ and you pick $\sigma = 1.2 - 1.3$, 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.

Relative Inefficiency and Computing Time

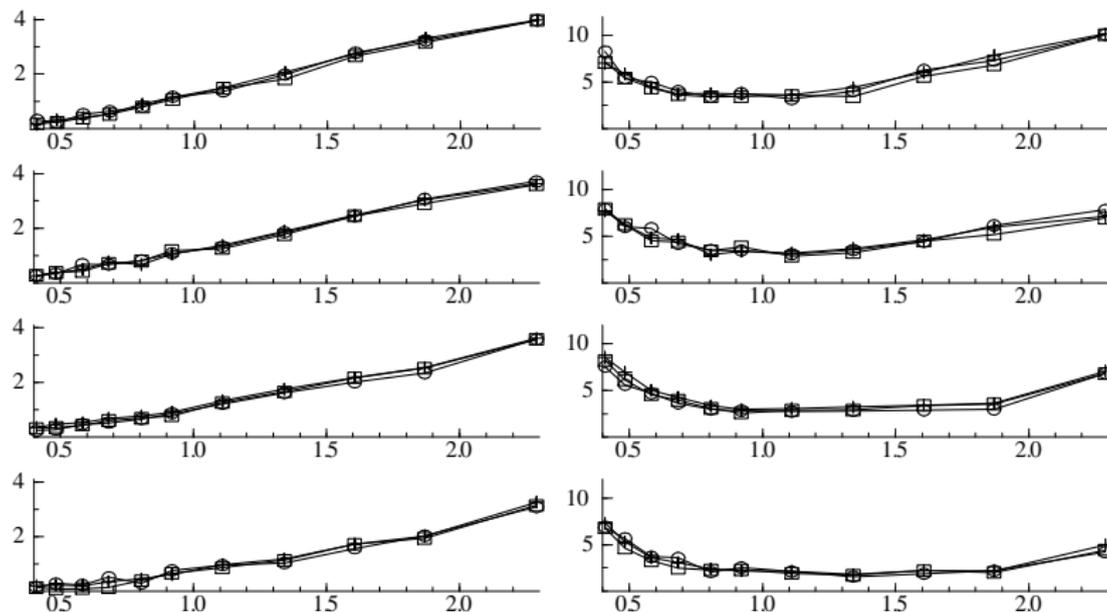


Figure: From left to right: RCT_h^Q vs N , RCT_h^Q vs $\sigma(\bar{\theta})$, RIF_h^Q against N and RIF_h^Q against $\sigma(\bar{\theta})$ for various values of ρ and different parameters.

- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.

- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.
- Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2$ is a sweet spot.

- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.
- Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2$ is a sweet spot.
- Pseudo-marginal MH scales in $\mathcal{O}(T^2)$ as we require $N \propto T$, while simulated likelihood scales in $\mathcal{O}(T^{3/2})$, i.e. $N \propto \sqrt{T}$.

- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.
- Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2$ is a sweet spot.
- Pseudo-marginal MH scales in $\mathcal{O}(T^2)$ as we require $N \propto T$, while simulated likelihood scales in $\mathcal{O}(T^{3/2})$, i.e. $N \propto \sqrt{T}$.
- However, pseudo-marginal MH much more generally applicable than simulated likelihood.

The Correlated Pseudo-Marginal Algorithm

- Reparameterize the likelihood estimator $\hat{p}_\theta(y_{1:T})$ as a function of normal variates $U \sim \mathcal{N}(0, I)$

$$\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; U)$$

The Correlated Pseudo-Marginal Algorithm

- Reparameterize the likelihood estimator $\hat{p}_\theta(y_{1:T})$ as a function of normal variates $U \sim \mathcal{N}(0, I)$

$$\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; U)$$

- Correlate estimators of $p_\theta(y_{1:T})$ and $p_\theta(y_{1:T})$ by setting

$$\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; V)$$

where

$$V = \rho U + \sqrt{1 - \rho^2} \varepsilon, \varepsilon \sim \mathcal{N}(0, I)$$

for $\rho \in (-1, 1)$.

The Correlated Pseudo-Marginal Algorithm

- Reparameterize the likelihood estimator $\hat{p}_\theta(y_{1:T})$ as a function of normal variates $U \sim \mathcal{N}(0, I)$

$$\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; U)$$

- Correlate estimators of $p_\theta(y_{1:T})$ and $p_\theta(y_{1:T})$ by setting

$$\hat{p}_\theta(y_{1:T}) = \hat{p}_\theta(y_{1:T}; V)$$

where

$$V = \rho U + \sqrt{1 - \rho^2} \varepsilon, \varepsilon \sim \mathcal{N}(0, I)$$

for $\rho \in (-1, 1)$.

- In practice, ρ will be select close to 1.

Correlated Pseudo-Marginal Metropolis–Hastings algorithm

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$ and $V = \rho U_{i-1} + \sqrt{1 - \rho^2} \varepsilon$, $\varepsilon \sim \mathcal{N}(0, I)$.

Correlated Pseudo-Marginal Metropolis–Hastings algorithm

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$ and $V = \rho U_{i-1} + \sqrt{1 - \rho^2} \varepsilon$, $\varepsilon \sim \mathcal{N}(0, I)$.
- Compute the estimate $\hat{p}_\vartheta(y_{1:T}; V)$ of $p_\vartheta(y_{1:T})$.

At iteration i

- Sample $\vartheta \sim q(\cdot | \vartheta_{i-1})$ and $V = \rho U_{i-1} + \sqrt{1 - \rho^2} \varepsilon$, $\varepsilon \sim \mathcal{N}(0, I)$.
- Compute the estimate $\hat{p}_\vartheta(y_{1:T}; V)$ of $p_\vartheta(y_{1:T})$.
- With probability

$$\min\left\{1, \frac{\hat{p}_\vartheta(y_{1:T}; V)}{\hat{p}_{\vartheta_{i-1}}(y_{1:T}; U_{i-1})} \frac{p(\vartheta)}{p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})}\right\}$$

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

Proposition. Let $N = N(T) \rightarrow \infty$ as $T \rightarrow \infty$ with $N = o(T)$. When $U \sim \bar{\pi}(\cdot|\theta)$ and $V = \rho_T U + \sqrt{1 - \rho_T^2} \varepsilon$ with $\rho_T = \exp(-\psi \frac{N}{T})$ then as $T \rightarrow \infty$

$$\log \left\{ \frac{\hat{p}_{\theta+\xi/\sqrt{T}}(y_{1:T}; V)}{\hat{p}_{\theta}(y_{1:T}; U)} / \frac{p_{\theta+\xi/\sqrt{T}}(y_{1:T})}{p_{\theta}(y_{1:T})} \right\} \Big| \mathcal{Y}^T, \mathcal{U}^T \Rightarrow \mathcal{N} \left(-\frac{\kappa^2(\theta)}{2}, \kappa^2(\theta) \right).$$

- This CLT is conditional on the observation sequence and the current auxiliary variables.

Proposition. Let $N = N(T) \rightarrow \infty$ as $T \rightarrow \infty$ with $N = o(T)$. When $U \sim \bar{\pi}(\cdot|\theta)$ and $V = \rho_T U + \sqrt{1 - \rho_T^2} \varepsilon$ with $\rho_T = \exp(-\psi \frac{N}{T})$ then as $T \rightarrow \infty$

$$\log \left\{ \frac{\hat{p}_{\theta+\xi/\sqrt{T}}(y_{1:T}; V)}{\hat{p}_{\theta}(y_{1:T}; U)} / \frac{p_{\theta+\xi/\sqrt{T}}(y_{1:T})}{p_{\theta}(y_{1:T})} \right\} \Big| \mathcal{Y}^T, \mathcal{U}^T \Rightarrow \mathcal{N} \left(-\frac{\kappa^2(\theta)}{2}, \kappa^2(\theta) \right).$$

- This CLT is conditional on the observation sequence and the current auxiliary variables.
- Asymptotically the distribution of the log-ratio decouples from the current location of the Markov chain.

Proposition. Let $N = N(T) \rightarrow \infty$ as $T \rightarrow \infty$ with $N = o(T)$. When $U \sim \bar{\pi}(\cdot|\theta)$ and $V = \rho_T U + \sqrt{1 - \rho_T^2} \varepsilon$ with $\rho_T = \exp(-\psi \frac{N}{T})$ then as $T \rightarrow \infty$

$$\log \left\{ \frac{\hat{p}_{\theta+\xi/\sqrt{T}}(y_{1:T}; V)}{\hat{p}_{\theta}(y_{1:T}; U)} / \frac{p_{\theta+\xi/\sqrt{T}}(y_{1:T})}{p_{\theta}(y_{1:T})} \right\} \Big| \mathcal{Y}^T, \mathcal{U}^T \Rightarrow \mathcal{N} \left(-\frac{\kappa^2(\theta)}{2}, \kappa^2(\theta) \right).$$

- This CLT is conditional on the observation sequence and the current auxiliary variables.
- Asymptotically the distribution of the log-ratio decouples from the current location of the Markov chain.
- The asymptotic variance is $O(1)$ even for $N \sim \log(T)$.

- **Assumption 1** - Asymptotic Normality: We have

$$\int \left| p(\theta | Y_{1:T}) - \phi(\theta; \hat{\theta}^T, \Sigma/T) \right| d\theta \xrightarrow{P} 0,$$

where $\hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ is a p.d. matrix.

- **Assumption 1** - Asymptotic Normality: We have

$$\int \left| p(\theta | Y_{1:T}) - \phi(\theta; \hat{\theta}^T, \Sigma/T) \right| d\theta \xrightarrow{P} 0,$$

where $\hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ is a p.d. matrix.

- **Assumption 2** - Proposal: $\vartheta = \theta + \xi/\sqrt{T}$ where $\varepsilon \sim v(\cdot)$ with $v(\xi) = v(-\xi)$.

- **Assumption 1** - Asymptotic Normality: We have

$$\int \left| p(\theta | Y_{1:T}) - \phi(\theta; \hat{\theta}^T, \Sigma/T) \right| d\theta \xrightarrow{P} 0,$$

where $\hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ is a p.d. matrix.

- **Assumption 2** - Proposal: $\vartheta = \theta + \xi/\sqrt{T}$ where $\varepsilon \sim v(\cdot)$ with $v(\xi) = v(-\xi)$.
- **Assumption 3** - For any θ in a neighbourhood of $\bar{\theta}$, the conditional CLT holds and $\kappa^2(\cdot)$ is continuous at $\bar{\theta}$.

Weak convergence

- Let $\{\vartheta_i^T\}_{i \geq 0}$ the stationary non-Markovian sequence of the correlated PM of invariant density $p(\theta | Y_{1:T})$.

Weak convergence

- Let $\{\vartheta_i^T\}_{i \geq 0}$ the stationary non-Markovian sequence of the correlated PM of invariant density $p(\theta | Y_{1:T})$.
- **Proposition** (Deligiannidis et al., 2016): The F.D.D. of the rescaled sequence $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \hat{\theta}_T)\}_{i \geq 0}$ converge weakly as $T \rightarrow \infty$ to those of a stationary Markov chain of invariant density $\phi(\tilde{\theta}; 0, \Sigma)$ and kernel given for $\tilde{\theta} \neq \tilde{\vartheta}$ by

$$\tilde{Q}(\tilde{\theta}, d\tilde{\vartheta}) = v(\tilde{\vartheta} - \tilde{\theta}) \mathbb{E}_R \left[\min \left\{ 1, \frac{\phi(\tilde{\vartheta}; 0, \Sigma)}{\phi(\tilde{\theta}; 0, \Sigma)} R \right\} \right] d\tilde{\vartheta}$$

where $R \sim \mathcal{N}(-\kappa^2(\bar{\theta})/2, \kappa^2(\bar{\theta}))$.

Weak convergence

- Let $\{\vartheta_i^T\}_{i \geq 0}$ the stationary non-Markovian sequence of the correlated PM of invariant density $p(\theta | Y_{1:T})$.
- **Proposition** (Deligiannidis et al., 2016): The F.D.D. of the rescaled sequence $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \hat{\theta}_T)\}_{i \geq 0}$ converge weakly as $T \rightarrow \infty$ to those of a stationary Markov chain of invariant density $\phi(\tilde{\vartheta}; 0, \Sigma)$ and kernel given for $\tilde{\theta} \neq \tilde{\vartheta}$ by

$$\tilde{Q}(\tilde{\theta}, d\tilde{\vartheta}) = v(\tilde{\vartheta} - \tilde{\theta}) \mathbb{E}_R \left[\min \left\{ 1, \frac{\phi(\tilde{\vartheta}; 0, \Sigma)}{\phi(\tilde{\theta}; 0, \Sigma)} R \right\} \right] d\tilde{\vartheta}$$

where $R \sim \mathcal{N}(-\kappa^2(\bar{\theta})/2, \kappa^2(\bar{\theta}))$.

- These results suggests that a simplified analysis of the CPM chain can be performed by looking at

$$\hat{Q}(\theta, d\vartheta) = q(\vartheta | \theta) \mathbb{E}_R \left[\min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} R \right\} \right] d\vartheta$$

where $R \sim \mathcal{N}(-\kappa^2/2, \kappa^2)$.

Breakdown

- An analysis based on this limiting kernel shows that one should select $\kappa^2 \approx 4.5$ to optimize the performance of the algorithm at fixed computational complexity.

Breakdown

- An analysis based on this limiting kernel shows that one should select $\kappa^2 \approx 4.5$ to optimize the performance of the algorithm at fixed computational complexity.
- Too good to be true? Can I really pick N arbitrarily?

Breakdown

- An analysis based on this limiting kernel shows that one should select $\kappa^2 \approx 4.5$ to optimize the performance of the algorithm at fixed computational complexity.
- Too good to be true? Can I really pick N arbitrarily?
- Weak convergence does NOT show that $\left| IF_h^Q - IF_h^{\hat{Q}} \right| \rightarrow 0$.

Breakdown

- An analysis based on this limiting kernel shows that one should select $\kappa^2 \approx 4.5$ to optimize the performance of the algorithm at fixed computational complexity.
- Too good to be true? Can I really pick N arbitrarily?
- Weak convergence does NOT show that $\left| IF_h^Q - IF_h^{\hat{Q}} \right| \rightarrow 0$.
- Informally, we have for $h(\theta) = \theta$

$$\text{Cov}(\theta_0, \theta_\tau) \approx \underbrace{\mathbb{E}(\mathbf{C}(\theta_0, \theta_\tau | U_0, U_\tau))}_{\text{fast}} + \underbrace{\mathbf{C}(\mathbb{E}(\theta_0 | U_0), \mathbb{E}(\theta_\tau | U_\tau))}_{\text{slow}}$$

where $\mathbb{E}(\theta_0 | U_0) \approx \hat{\theta}^T + \Sigma / T \nabla_\theta \log \hat{p}_\theta(y_{1:T}; U) / p_\theta(y_{1:T})|_{\hat{\theta}^T}$ and $IF_h^Q \rightarrow \infty$ if $N / \sqrt{T} \rightarrow 0$.

Breakdown

- An analysis based on this limiting kernel shows that one should select $\kappa^2 \approx 4.5$ to optimize the performance of the algorithm at fixed computational complexity.
- Too good to be true? Can I really pick N arbitrarily?
- Weak convergence does NOT show that $\left| IF_h^Q - IF_h^{\hat{Q}} \right| \rightarrow 0$.
- Informally, we have for $h(\theta) = \theta$

$$\text{Cov}(\theta_0, \theta_\tau) \approx \underbrace{\mathbb{E}(\mathbf{C}(\theta_0, \theta_\tau | U_0, U_\tau))}_{\text{fast}} + \underbrace{\mathbf{C}(\mathbb{E}(\theta_0 | U_0), \mathbb{E}(\theta_\tau | U_\tau))}_{\text{slow}}$$

where $\mathbb{E}(\theta_0 | U_0) \approx \hat{\theta}^T + \Sigma / T \nabla_{\theta} \log \hat{p}_{\theta}(y_{1:T}; U) / p_{\theta}(y_{1:T})|_{\hat{\theta}^T}$ and $IF_h^Q \rightarrow \infty$ if $N / \sqrt{T} \rightarrow 0$.

- To ensure IF_h^Q , we need at least $N \propto \sqrt{T}$ and we conjecture it is sufficient.

Example: Gaussian Latent Variable Model

- Consider the toy model

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1), \quad Y_t | X_t \sim \mathcal{N}(X_t, \sigma^2).$$

Example: Gaussian Latent Variable Model

- Consider the toy model

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1), \quad Y_t | X_t \sim \mathcal{N}(X_t, \sigma^2).$$

- The likelihood can be computed exactly, allowing to implement the “exact” MH algorithm.

Example: Gaussian Latent Variable Model

- Consider the toy model

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1), \quad Y_t | X_t \sim \mathcal{N}(X_t, \sigma^2).$$

- The likelihood can be computed exactly, allowing to implement the “exact” MH algorithm.
- The likelihood estimator is based on importance sampling.

Example: Gaussian Latent Variable Model

- Consider the toy model

$$X_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1), \quad Y_t | X_t \sim \mathcal{N}(X_t, \sigma^2).$$

- The likelihood can be computed exactly, allowing to implement the “exact” MH algorithm.
- The likelihood estimator is based on importance sampling.
- Integrated Autocorrelation Time is referred to as the Inefficiency IF .

Example: Gaussian Latent Variable Model

MH ($T = 8192$)		IF(θ)	
		15.6	
PM ($\rho = 0.0$)			
N		RIF(θ)	RCT(θ)
5000		2.2	11210
CPM ($\rho = 0.9963$)			
N	κ	RIF(θ)	RCT(θ)
9	3.1	14.0	126.2
12	2.7	8.3	99.7
20	2.2	4.7	93.3
25	2.0	2.8	69.3
35	1.7	1.7	61.1
56	1.3	1.6	87.0
80	1.1	1.1	89.0
120	0.9	0.9	113.5

Here $RIF = IF/IF_{MH}$ and $RCT = N \times RIF$.

- In i.i.d. case, very substantial improvement over the PM algorithm can be achieved by introducing a correlation scheme.

- In i.i.d. case, very substantial improvement over the PM algorithm can be achieved by introducing a correlation scheme.
- Analysis suggests that complexity is $O\left(T\sqrt{T}\right)$ vs $O\left(T^2\right)$.

- In i.i.d. case, very substantial improvement over the PM algorithm can be achieved by introducing a correlation scheme.
- Analysis suggests that complexity is $O\left(T\sqrt{T}\right)$ vs $O\left(T^2\right)$.
- In state-space models, implementation relies on non-standard particle filter scheme (Hilbert sorting): our analysis does not hold experimentally for state dimension > 1 and theoretically and but still substantial gains.

- In i.i.d. case, very substantial improvement over the PM algorithm can be achieved by introducing a correlation scheme.
- Analysis suggests that complexity is $O\left(T\sqrt{T}\right)$ vs $O\left(T^2\right)$.
- In state-space models, implementation relies on non-standard particle filter scheme (Hilbert sorting): our analysis does not hold experimentally for state dimension > 1 and theoretically and but still substantial gains.
- Novel pseudo-marginal scheme using Conditional Sequential Monte Carlo (Andrieu, A.D., Yildirim, 2016) appears to suggest $O(T)$ is feasible.

Experimental results using conditional SMC

	Novel c-SMC PM		Standard PM	
	σ_v^2	σ_w^2	σ_v^2	σ_w^2
$T = 1000$	17.7	23.5	71.2	59.2
$T = 2000$	17.5	23.7	759.0	757.9
$T = 5000$	17.6	23.7	5808.6	5663.5
$T = 10000$	17.6	23.6	7368.1	7176.9

Estimated IACT on a nonlinear state-space model for $N = 200$ for novel c-SMC PM algorithm and $N = 2000$ for standard PM algorithm

Some References

- C. Andrieu, A.D. & R. Holenstein, “Particle Markov chain Monte Carlo Methods”, *JRSS B*, 2010.
- C. Andrieu & G.O. Roberts, “The Pseudo-Marginal Algorithm for Bayesian Computation”, *Ann. Stat.*, 2009.
- J. Berard, P. Del Moral & A.D., “A Lognormal CLT for Particle Approximations of Normalizing Constants”, *Electronic J. Proba.*, 2014.
- A.D., M.K. Pitt, G. Deligiannidis and R. Kohn, “Efficient Implementation of Markov Chain Monte Carlo when Using an Unbiased Likelihood Estimator”, *Biometrika*, 2015.
- L. Lin, K. Lin & J. Sloan, “A Noisy Monte Carlo Algorithm”, *Phys. Rev. D*, 2000.
- M.K. Pitt, R. Silva, P. Giordani & R. Kohn, “On Some Properties of MCMC Simulation Methods Based of the Particle Filter”, *J. Econometrics*, 2012.
- C. Sherlock, A. Thiery, G.O. Roberts & J.S. Rosenthal, “On the Efficiency of the RW Pseudo-Marginal MH”, *Ann. Stat.*, 2015.