

Correlated Particle Metropolis

Michael Pitt

King's College London

Joint work George Deligiannidis & Arnaud Doucet (Oxford)

Imperial College, May 2018

Organisation of the Talk

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

Sequential Monte Carlo Estimator

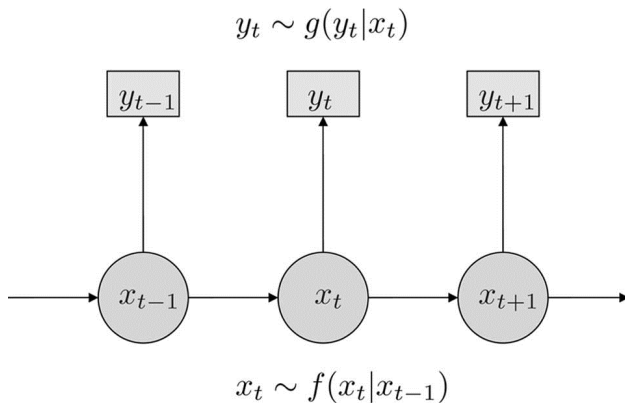
- Assume $\{X_t\}_{t \geq 1}$ is a latent Markov process, i.e. $X_1 \sim \mu_\theta(\cdot)$ and
$$X_{t+1} | (X_t = x) \sim f_\theta(\cdot | x), \quad Y_t | (X_t = x) \sim g_\theta(\cdot | x).$$
- Observations $\{Y_t\}_{t \geq 1}$ are conditionally independent given $\{X_t\}_{t \geq 0}$.
- Likelihood of $y_{1:T} = (y_1, \dots, y_T)$ is

$$p(y_{1:T}; \theta) = \int_{\mathbb{X}^{T+1}} p(x_{0:T}, y_{1:T}; \theta) dx_{0:T}.$$

Sequential Monte Carlo Estimator

- SMC provides an unbiased estimator of relative variance $\mathcal{O}(T/N)$ where N is the number of particles.
- Whatever being $N \geq 1$, the pseudo-marginal MH admits $\pi(\theta)$ as invariant distribution.

A directed acyclic graph (DAG) of the problem:



- Both $g(y_t|x_t)$, $f(x_t|x_{t-1})$ may be indexed by fixed parameters θ .
- Filtering density

$$p(x_t|y_{1:t}; \theta).$$

- As an important byproduct we also obtain the one step predictive density

$$p(y_t|y_{1:t-1}; \theta).$$

- Yields the likelihood (KF)

$$p(y_{1:T}|\theta) = p(y_1|\theta) \prod_{t=1}^{T-1} p(y_{t+1}|y_{1:t}; \theta).$$

Particle Filter Estimation

- Simulation based methods to perform filtering in nonlinear/non-Gaussian state space models.
- See Gordon, Salmond and Smith (1993) (GSS), Kitagawa (1996), Pitt and Shephard (1999) and reviewed by Doucet et al. (2000).
- We aim to have 'particles', x_t^1, \dots, x_t^N with associated discrete probability masses π_t^1, \dots, π_t^N , drawn from the density $f(x_t|y_{1:t})$.

Bootstrap Filter: GSS (1993, IEE)

We start at $t = 0$ with samples from $x_0^k \sim p(x_0)$. For $t=1, \dots, T$:
We have samples $x_t^k \sim p(x_t|y_{1:t})$ for $k = 1, \dots, N$.

① For $k = 1 : N$, sample $\tilde{x}_{t+1}^k \sim f(x_{t+1}|x_t^k)$.

② For $k = 1 : N$,

$$\pi_{t+1}^k = \frac{g(y_{t+1}|\tilde{x}_{t+1}^k)}{\sum_{j=1}^N g(y_{t+1}|\tilde{x}_{t+1}^j)}.$$

③ For $j = 1 : N$, sample $x_{t+1}^j \sim \sum_{k=1}^N \pi_{t+1}^k \delta(x_{t+1}^j - \tilde{x}_{t+1}^k)$.

Step 3 multinomial (or stratified) sampling (*from the mixture*).

This will yield an *approximate sample* the desired posterior density, $f(x_t|y_{1:t})$ as t varies.

SMC Likelihood Estimation

- Parameter estimation using likelihood function, via prediction decomposition given by;

$$\log L(\theta) = \log p(y_1, \dots, y_T | \theta) = \sum_{t=1}^T \log p(y_{t+1} | \theta; y_{1:t}).$$

- We need to estimate the function :

$$\hat{p}(y_{1:T} | \theta) = \hat{p}(y_1 | \theta) \prod_{t=1}^{T-1} \hat{p}(y_{t+1} | y_{1:t}; \theta),$$

$$\hat{p}(y_{t+1} | \theta; y_{1:t}) = \frac{1}{N} \sum_{i=1}^N p(y_{t+1} | \tilde{x}_{t+1}^i).$$

where $\tilde{x}_{t+1}^i \sim f(x_{t+1} | y_{1:t}; \theta)$, from step (2).

- Remarkably (just like in IS) the lik estimator $\hat{p}(y_{1:T} | \theta)$ is unbiased for $p(y_{1:T} | \theta)$ regardless of N , Del Moral (04).

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

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

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

Particle Metropolis: Intractable Likelihood Function

- In numerous scenarios, $p(y; \theta)$ cannot be evaluated pointwise; e.g.

$$p(y; \theta) = \int p(x, y; \theta) dx$$

where the integral cannot be evaluated.

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

$$p(\theta, x|y) = \frac{p(x, y; \theta) p(\theta)}{p(y)}$$

Standard MCMC Approaches

- Standard MCMC schemes target

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

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

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

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. non-Gaussian state-space models.

Standard MCMC Approaches

- Standard MCMC schemes target

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

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. non-Gaussian state-space models.
- **Problem 2:** Even when it is implementable, Gibbs can converge very slowly e.g. diffusions

Standard MCMC Approaches

- Standard MCMC schemes target

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

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. non-Gaussian state-space models.
- **Problem 2:** Even when it is implementable, Gibbs can converge very slowly e.g. diffusions
- **Problem 3:** It may only be possible to generate from $f(x_t | x_{t-1})$ not to evaluate it, e.g. DSGE models.

Standard MCMC Approaches

- Standard MCMC schemes target

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

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. non-Gaussian state-space models.
- **Problem 2:** Even when it is implementable, Gibbs can converge very slowly e.g. diffusions
- **Problem 3:** It may only be possible to generate from $f(x_t | x_{t-1})$ not to evaluate it, e.g. DSGE models.
- Pseudo-marginal methods mimic an algorithm targetting directly $\pi(\theta) = p(\theta | y_{1:T})$ instead of $p(\theta, x_{1:T} | y_{1:T})$.

Pseudo Metropolis (PM)

- Let $\hat{p}(y; \theta, U)$ be an *unbiased non-negative estimator of the likelihood* where $U \sim m_\theta(\cdot)$; i.e.

$$p(y; \theta) = \int_{\mathbf{U}} \hat{p}(y; \theta, u) m_\theta(u) du.$$

- Introduce a target distribution on $\Theta \times \mathbf{U}$ of density

$$\bar{\pi}(\theta, u) = \pi(\theta) \frac{\hat{p}(y; \theta, u)}{p(y; \theta)} m_\theta(u) = \frac{p(\theta) \hat{p}(y; \theta, u) m_\theta(u)}{p(y)}$$

- Then unbiasedness yields

$$\int_{\mathbf{U}} \bar{\pi}(\theta, u) du = \pi(\theta)$$

Any MCMC algorithm sampling from $\bar{\pi}(\theta, u)$ yields samples from $\pi(\theta)$.

Pseudo-Marginal Metropolis-Hastings algorithm

- Can form an **unbiased estimator**, based on N particles $\hat{p}(y; \theta, U)$.
- Set $(\vartheta^{(0)}, U^{(0)})$ and iterate for $j = 1, 2, \dots$

Sample $\vartheta \sim q(\cdot | \vartheta^{(j-1)})$, $U \sim m_{\vartheta}(\cdot)$ to obtain $\hat{p}(y; \vartheta, U)$.

- 1 Compute

$$\alpha = 1 \wedge \frac{\hat{p}(y; \vartheta, U)}{\hat{p}(y; \vartheta^{(j-1)}, U^{(j-1)})} \frac{p(\vartheta)}{p(\vartheta^{(j-1)})} \frac{q(\vartheta^{(j-1)} | \vartheta)}{q(\vartheta | \vartheta^{(j-1)})}$$

- 2 With proba α , set $(\vartheta^{(j)}, U^{(j)}) := (\vartheta, U)$ and stay where you are otherwise.

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 $p(y_{1:T}; \theta)$ is estimated using SMC with N particles.

A Nonlinear State-Space Model

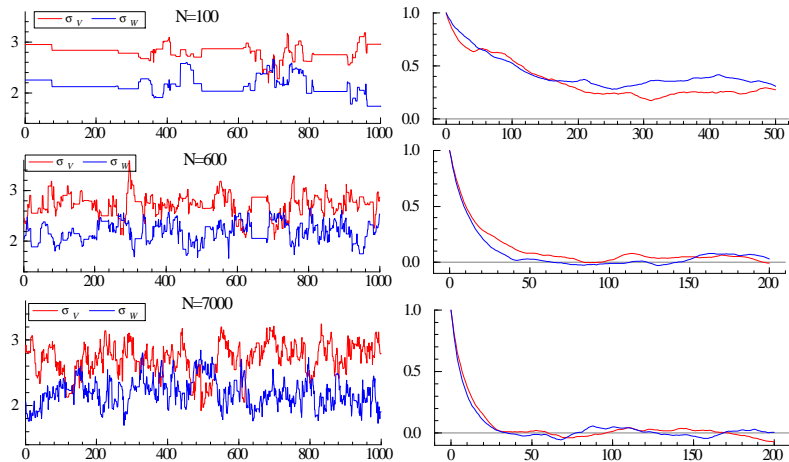


Figure: Autocorrelation of $\{\sigma_V^{(i)}\}$ and $\{\sigma_W^{(i)}\}$ of the MH sampler for various N .

How to Select the Number of Samples

- If N is too **small**, then the algorithm **mixes poorly** and will require many MCMC iterations.
- If N is too **large**, then each iteration is **expensive** due to estimating the likelihood.
- Equivalently we will examine σ^2 , the **variance** of the estimator of the log-likelihood to trade off these two concerns.
- **Simple Guideline:** We find the optimal value for σ around 1. *JoE* (2012, Kohn, Giordani, Silva), *Biometrika* (2015, Doucet, Deligiannidis, Kohn).

- Consider the error in the log-likelihood estimator

$$Z = \log \hat{p}(y; \theta, U) - \log p(y; \theta) \sim g_{\theta}(\cdot)$$

- In the (θ, Z) parameterization, the target is

$$\bar{\pi}(\theta, u) = \pi(\theta) \frac{\hat{p}(y; \theta, u)}{p(y; \theta)} m_{\theta}(u) \Rightarrow \bar{\pi}(\theta, z) = \pi(\theta) \exp(z) g_{\theta}(z).$$

How precise should the log-likelihood estimator be?

- **Aim:** Minimize the “computational time”

$$CT(Q, h) = IACT(Q, h) / \sigma^2$$

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

How precise should the log-likelihood estimator be?

- **Aim:** Minimize the “computational time”

$$CT(Q, h) = \text{IACT}(Q, h) / \sigma^2$$

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

- where

$$\text{IACT}(Q, h) = \text{Integrated Autocorrelation Time of } \{h(\vartheta_i)\}_{i \geq 1}$$

How precise should the log-likelihood estimator be?

- The IACT is

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

- where Q is the pseudo-marginal kernel with acceptance criterion

$$\min \{1, r(\theta, \vartheta) \exp(w - z)\}.$$

- where $r(\theta, \vartheta)$ is exact Metropolis ratio,
- and

$$z = \log\{\widehat{p}_\theta(y_{1:T}) / p_\theta(y_{1:T})\}, \quad w = \log\{\widehat{p}_\vartheta(y_{1:T}) / p_\vartheta(y_{1:T})\}.$$

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

$$\bar{\pi}(\theta, z) = \pi(\theta) \underbrace{\exp(z) g(z)}_{\pi_Z(z|\sigma)} = \pi(\theta) \mathcal{N}(z; \sigma^2/2; \sigma^2).$$

- Justified empirically and theoretically (a CLT and concentration).

Pseudo-proof:

IID data, fixed θ (IS)

$$Z = \log \hat{p}(y|\theta) - \log p(y|\theta) = \sum_{t=1}^T \log \left\{ 1 + \frac{\gamma_t}{\sqrt{N}} \varepsilon_t \right\},$$

where ε_t white noise, variance 1 and $N = \beta T$,

$$\simeq \sum_{t=1}^T \frac{\gamma_t}{\sqrt{\beta T}} \varepsilon_t - \frac{1}{2} \frac{\gamma_t^2}{\beta T} \varepsilon_t^2 \longrightarrow \mathcal{N} \left(-\frac{1}{2} \frac{\bar{\gamma}^2}{\beta}, \frac{\bar{\gamma}^2}{\beta} \right),$$

where $\bar{\gamma}^2 = T^{-1} \sum_{t=1}^T \gamma_t^2$.

- More detail (and more generality, PFs) in Bérard, Del Moral, Doucet (EJP, 2014).
- Need concentration (in θ as $T \rightarrow \infty$), unproven.

Simulation:

SSF (from JoE(2012, Kohn, Giordani, Silva))

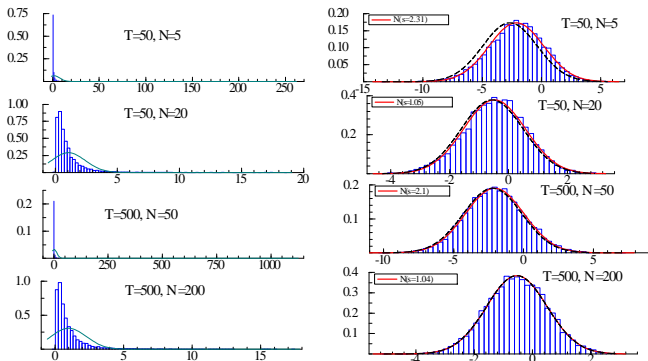


Figure: AR(1) plus noise model with fixed parameters. replications is 10,000. SIR likelihood estimator (divided by the true likelihood) on LEFT and for the error in the log of the SIR likelihood estimator on RIGHT. Both N and T vary

- When $q(\vartheta|\theta) = \pi(\vartheta)$, $\sigma_{\text{opt}} = 0.92$ (Pitt et al., JoE 2012). The acf simply reduces to:

$$\phi_n(\theta, Q) = \int \Pr(R | z; \sigma)^n \pi_Z(z|\sigma) dz$$
$$I_{ACT}(\sigma) = \int \frac{1 + \Pr(R | z; \sigma)}{1 - \Pr(R | z; \sigma)} \pi_Z(z|\sigma) dz.$$

$\Pr(R | z; \sigma)$ analytically available.

- Sherlock, Thiery, Roberts and Rosenthal, (Annals 2015) consider the (joint) optimisation for a limiting target in the important case of RWM. Again analytically available and exact. Similar results.

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 $\bar{\pi}(\theta, z)$ -reversible kernel

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

where we have a product acceptance criterion:

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

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 $\bar{\pi}(\theta, z)$ -reversible kernel

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

where we have a product acceptance criterion:

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

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

pseudo-proof: Jump chain

IACT: $\mathbb{E}_\pi[x] = 0$, on J.C.s Douc and Robert (Annals, 2011)

$$\tilde{\pi}(x) = \frac{\pi(x)p(x)}{P_A}$$

$$\begin{aligned} (IACT + 1)\mathbb{E}_\pi[x_o^2] &= 2 \sum_{t=0}^{\infty} \mathbb{E}_\pi[x_o x_t] \\ &= 2 \sum_{\tau=0}^{\infty} \mathbb{E}_{\pi, \tilde{\pi}}[x_o \tau \tilde{x}_\tau] \\ &= 2 \sum_{\tau=0}^{\infty} \mathbb{E}_{\pi, \tilde{\pi}} \left[x_o \frac{\tilde{x}_\tau}{p(\tilde{x}_\tau)} \right] \text{ because sojourns Geometri} \\ &= 2P_A \sum_{\tau=0}^{\infty} \mathbb{E}_{\tilde{\pi}} \left[\frac{\tilde{x}_o}{p(\tilde{x}_o)} \frac{\tilde{x}_\tau}{p(\tilde{x}_\tau)} \right] \text{ change of measure for} \\ &= 2P_A \left(\widetilde{IACT} + 1 \right) \mathbb{E}_{\tilde{\pi}} \left[\frac{\tilde{x}_o^2}{p(\tilde{x}_o)^2} \right]. \end{aligned}$$

Simpler Bounds on the Relative Inefficiency

- We obtain an explicit expression for $IF_h^{Q^*}(\sigma)$.
- If $IF_{h/\varrho_{EX}}^{\tilde{Q}^{EX}} \geq 1$, e.g. \tilde{Q}^{EX} is a positive kernel, then

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

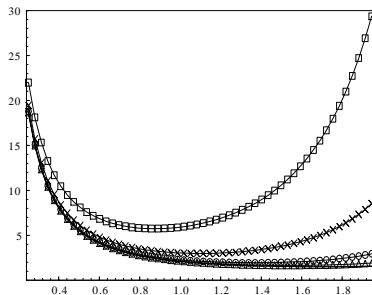
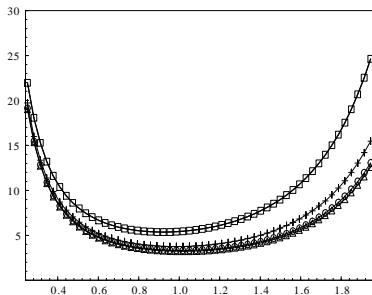
and the bound is tight as $IF_h^{EX} \rightarrow 1$ or $\sigma \rightarrow 0$.

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

$$\frac{IF_h^{Q^*}(\sigma)}{IF_h^{EX}} \rightarrow \frac{1}{\pi_Z^\sigma(\varrho_Z^\sigma)} = \frac{1}{P_A(\sigma)}.$$

- 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 Time



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

Application: Stochastic Volatility Model

- Chernov et al., *J. Econometrics* (2003) and Huang & Tauchen, *J. Financial Econometrics* (2005):

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

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

- Euler discretization of the volatilities $v_1(t)$ and $v_2(t)$ provides closed form expression for $Y(s) = \log P(\tau_{s+1}) - \log P(\tau_s)$.
- Straightforward to estimate the likelihood and simulate forward.

Application: Stochastic Volatility Model

- 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 $\bar{\theta}$.

Empirical vs Assumed Distributions of Z for SV model:

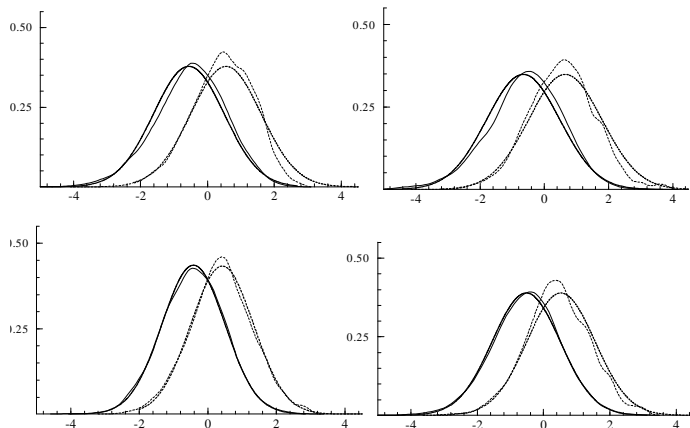


Figure: Empirical distributions (dashed) vs assumed Gaussians (solid) of Z . $T = 300$ and $T = 2700$ at $\bar{\theta}$ (left) and marginalized over $\pi(\theta)$.

Integrated Autocorrelation Time of Pseudo-Marginal MH

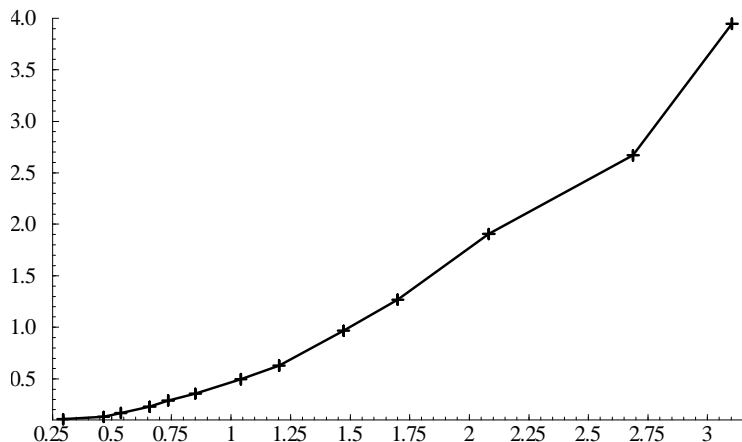


Figure: Average over the 9 parameter components of the log-integrated autocorrelation time of pseudo-marginal chain as a function of σ for $T = 300$.

Computational time for the SV model

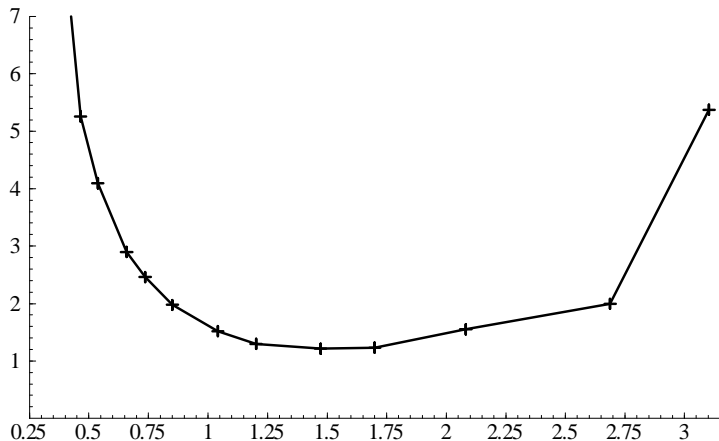


Figure: Computational time as a function of σ

- **Guideline:** Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2 - 1.3$ is a sweet spot.

Guideline and Discussion

- **Guideline:** Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2 - 1.3$ is a sweet spot.
- Pseudo-marginal scales in $\mathcal{O}(T^2)$ at each iteration as we require $N \propto T$.

- **Guideline:** Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2 - 1.3$ is a sweet spot.
- Pseudo-marginal scales in $\mathcal{O}(T^2)$ at each iteration as we require $N \propto T$.
- For i.i.d. data, simulated ML estimator is efficient as long as N increases at a rate faster than \sqrt{T} ; e.g. Lee, *Econometric Theory*, 1999.

- **Guideline:** Optimal σ depends on efficiency of the ideal MH algorithm but $\sigma \approx 1.2 - 1.3$ is a sweet spot.
- Pseudo-marginal scales in $\mathcal{O}(T^2)$ at each iteration as we require $N \propto T$.
- For i.i.d. data, simulated ML estimator is efficient as long as N increases at a rate faster than \sqrt{T} ; e.g. Lee, *Econometric Theory*, 1999.
- **Problem:** the ratio $p_{\theta}(y_{1:T}) / p_{\theta}(y_{1:T})$ is estimated by estimating independently the numerator and denominator in pseudo-marginal.

Correlated Pseudo-Marginal Algorithm

- Previously, we consider the likelihood estimator $\hat{p}_\theta(y_{1:T}; U)$ where $U \sim m_\theta(\cdot)$.
- Reparameterize the likelihood estimator so that $U \sim \mathcal{N}(0, I)$.
- 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

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

- In practice, ρ will be selected close to 1.
- The invariance of the PMH unaffected.
- Can be seen as a solution of a standard O-U over $[0, \delta]$, $\rho = \exp(-\delta)$.

Correlated Pseudo-Marginal Algorithm

- Correlated pseudo-marginal simulates a Markov chain $\{\theta_i, U_i\}_{i \geq 1}$ of limiting distribution $\bar{\pi}(\theta, u)$.

Correlated Pseudo-Marginal Algorithm

- Correlated pseudo-marginal simulates a Markov chain $\{\theta_i, U_i\}_{i \geq 1}$ of limiting distribution $\bar{\pi}(\theta, u)$.

Correlated Pseudo-Marginal Algorithm

- Correlated pseudo-marginal simulates a Markov chain $\{\vartheta_i, U_i\}_{i \geq 1}$ of limiting distribution $\bar{\pi}(\theta, u)$.

At iteration i

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

Correlated Pseudo-Marginal Algorithm

- Correlated pseudo-marginal simulates a Markov chain $\{\vartheta_i, U_i\}_{i \geq 1}$ of limiting distribution $\bar{\pi}(\theta, u)$.

At iteration i

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

Correlated Pseudo-Marginal Algorithm

- Correlated pseudo-marginal simulates a Markov chain $\{\vartheta_i, U_i\}_{i \geq 1}$ of limiting distribution $\bar{\pi}(\theta, u)$.

At iteration i

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

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

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

Likelihood estimation for state-space models

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \hat{p}_\theta(y; u)$ is not continuous and displays large variations for moderate N .

Likelihood estimation for state-space models

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \hat{p}_\theta(y; u)$ is not continuous and displays large variations for moderate N .
- Discontinuities arise from the resampling step: you can end up picking resampling particles very far from each other even with small variations in (θ, u) .

Likelihood estimation for state-space models

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \hat{p}_\theta(y; u)$ is not continuous and displays large variations for moderate N .
- Discontinuities arise from the resampling step: you can end up picking resampling particles very far from each other even with small variations in (θ, u) .
- For $n = 1$, (Pitt et al, 2012) propose sorting particles $X_t^{\sigma_t(1)} \leq \dots \leq X_t^{\sigma_t(N)}$.

Likelihood estimation for state-space models

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \hat{p}_\theta(y; u)$ is not continuous and displays large variations for moderate N .
- Discontinuities arise from the resampling step: you can end up picking resampling particles very far from each other even with small variations in (θ, u) .
- For $n = 1$, (Pitt et al, 2012) propose sorting particles
$$X_t^{\sigma_t(1)} \leq \dots \leq X_t^{\sigma_t(N)}.$$
- For $n \geq 2$, (Gerber & Chopin, 2015) use Hilbert curve; e.g. map particles to $[0, 1]^n$ using $\psi : \mathbb{R}^n \rightarrow [0, 1]^n$ (e.g. logistic) and $h : [0, 1]^n \rightarrow [0, 1]$ (pseudo-inverse Hilbert curve) and sort projected particles on $[0, 1]$.

Likelihood estimation for state-space models

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \hat{p}_\theta(y; u)$ is not continuous and displays large variations for moderate N .
- Discontinuities arise from the resampling step: you can end up picking resampling particles very far from each other even with small variations in (θ, u) .
- For $n = 1$, (Pitt et al, 2012) propose sorting particles
$$X_t^{\sigma_t(1)} \leq \dots \leq X_t^{\sigma_t(N)}.$$
- For $n \geq 2$, (Gerber & Chopin, 2015) use Hilbert curve; e.g. map particles to $[0, 1]^n$ using $\psi : \mathbb{R}^n \rightarrow [0, 1]^n$ (e.g. logistic) and $h : [0, 1]^n \rightarrow [0, 1]$ (pseudo-inverse Hilbert curve) and sort projected particles on $[0, 1]$.
- Alternative coupling ideas have been used to mitigate these fluctuations (Jacob et al., arXiv 2016).

- **Assumption 1** - *Asymptotic Normality*: $\exists \hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ a p.d. matrix s.t.

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

- **Assumption 1** - *Asymptotic Normality*: $\exists \hat{\theta}^T \xrightarrow{P} \bar{\theta}$ and Σ a p.d. matrix s.t.

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

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

Large sample analysis of the correlated PM - i.i.d. case

Proposition. Let $N \rightarrow \infty$ as $T \rightarrow \infty$ with $N = o(T)$. When $U \sim \bar{\pi}_T(\cdot|\theta)$ and $U' = \rho U + \sqrt{1 - \rho^2} \varepsilon$ with $\rho = \exp(-\psi \frac{N}{T})$ then as $T \rightarrow \infty$

$$\log \left\{ \frac{\hat{p}_{\theta+\xi/\sqrt{T}}(Y_{1:T}; U')}{\hat{p}_{\theta}(Y_{1:T}; U)} / \frac{p_{\theta+\xi/\sqrt{T}}(Y_{1:T})}{p_{\theta}(Y_{1:T})} \right\} \Bigg| \mathcal{Y}^T, \mathcal{U}^T \Rightarrow \mathcal{N}\left(-\frac{\kappa^2(\theta)}{2}, \kappa^2(\theta)\right)$$

where

$$\kappa^2(\theta) = 2\psi \mathbb{E} \left[\|\partial_u \omega(U, Y; \theta)\|^2 \right]$$

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

Large sample analysis of the correlated PM - i.i.d. case

Proposition. Let $N \rightarrow \infty$ as $T \rightarrow \infty$ with $N = o(T)$. When $U \sim \bar{\pi}_T(\cdot|\theta)$ and $U' = \rho U + \sqrt{1 - \rho^2} \varepsilon$ with $\rho = \exp(-\psi \frac{N}{T})$ then as $T \rightarrow \infty$

$$\log \left\{ \frac{\hat{p}_{\theta+\xi/\sqrt{T}}(Y_{1:T}; U')}{\hat{p}_{\theta}(Y_{1:T}; U)} / \frac{p_{\theta+\xi/\sqrt{T}}(Y_{1:T})}{p_{\theta}(Y_{1:T})} \right\} \Bigg| \mathcal{Y}^T, \mathcal{U}^T \Rightarrow \mathcal{N}\left(-\frac{\kappa^2(\theta)}{2}, \kappa^2(\theta)\right)$$

where

$$\kappa^2(\theta) = 2\psi \mathbb{E} \left[\|\partial_u \omega(U, Y; \theta)\|^2 \right]$$

- 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.

Large sample analysis of the correlated PM - i.i.d. case

Proposition. Let $N \rightarrow \infty$ as $T \rightarrow \infty$ with $N = o(T)$. When $U \sim \bar{\pi}_T(\cdot|\theta)$ and $U' = \rho U + \sqrt{1 - \rho^2} \varepsilon$ with $\rho = \exp(-\psi \frac{N}{T})$ then as $T \rightarrow \infty$

$$\log \left\{ \frac{\hat{p}_{\theta+\xi/\sqrt{T}}(Y_{1:T}; U')}{\hat{p}_{\theta}(Y_{1:T}; U)} / \frac{p_{\theta+\xi/\sqrt{T}}(Y_{1:T})}{p_{\theta}(Y_{1:T})} \right\} \Big| \mathcal{Y}^T, \mathcal{U}^T \Rightarrow \mathcal{N}\left(-\frac{\kappa^2(\theta)}{2}, \kappa^2(\theta)\right)$$

where

$$\kappa^2(\theta) = 2\psi \mathbb{E} \left[\|\partial_u \omega(U, Y; \theta)\|^2 \right]$$

- 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 $N \sim \log(T)$.

Large sample analysis of the correlated PM - i.i.d. case

- Let $\Theta_T := \{\theta_i^T\}_{i \geq 0}$ be the stationary *non-Markovian* sequence of the correlated PM targetting $p(\theta | Y_{1:T})$.

Large sample analysis of the correlated PM - i.i.d. case

- Let $\Theta_T := \{\vartheta_i^T\}_{i \geq 0}$ be the stationary *non-Markovian* sequence of the correlated PM targetting $p(\theta | Y_{1:T})$.
- **Proposition** (Deligiannidis et al., 2016): The sequences $\{\Theta_T\}_{T \geq 1}$ converge weakly as $T \rightarrow \infty$ to a stationary Markov chain of invariant density $\phi(\tilde{\theta}; 0, \Sigma)$ and kernel given for $\tilde{\theta} \neq \tilde{\theta}'$ by

$$\tilde{Q}(\tilde{\theta}, d\tilde{\theta}') = v(\tilde{\theta}' - \tilde{\theta}) \mathbb{E}_{R \sim \mathcal{N}(-\kappa^2/2, \kappa^2)} \left[\min \left\{ 1, \frac{\phi(\tilde{\theta}'; 0, \Sigma)}{\phi(\tilde{\theta}; 0, \Sigma)} \exp R \right\} \right] d\tilde{\theta}'$$

where $\kappa := \kappa(\bar{\theta})$.

Large sample analysis of the correlated PM - i.i.d. case

- Let $\Theta_T := \{\vartheta_i^T\}_{i \geq 0}$ be the stationary *non-Markovian* sequence of the correlated PM targetting $p(\theta | Y_{1:T})$.
- **Proposition** (Deligiannidis et al., 2016): The sequences $\{\Theta_T\}_{T \geq 1}$ converge weakly as $T \rightarrow \infty$ to a stationary Markov chain of invariant density $\phi(\tilde{\theta}; 0, \Sigma)$ and kernel given for $\tilde{\theta} \neq \tilde{\theta}'$ by

$$\tilde{Q}(\tilde{\theta}, d\tilde{\theta}') = v(\tilde{\theta}' - \tilde{\theta}) \mathbb{E}_{R \sim \mathcal{N}(-\kappa^2/2, \kappa^2)} \left[\min \left\{ 1, \frac{\phi(\tilde{\theta}'; 0, \Sigma)}{\phi(\tilde{\theta}; 0, \Sigma)} \exp R \right\} \right] d\tilde{\theta}'$$

where $\kappa := \kappa(\bar{\theta})$.

- It is tempting to use this result to provide guidelines on the optimization of CPM... but one has to be careful.

Decomposition in fast and slow components

- For a stationary CPM chain (ϑ_i, U_i) , decompose

$$h(\vartheta_i) = \underbrace{\mathbb{E}[h(\vartheta) | U_i]}_{\text{slow}} + \underbrace{h(\vartheta_i) - \mathbb{E}[h(\vartheta) | U_i]}_{\text{fast}}.$$

Decomposition in fast and slow components

- For a stationary CPM chain (ϑ_i, U_i) , decompose

$$h(\vartheta_i) = \underbrace{\mathbb{E}[h(\vartheta) | U_i]}_{\text{slow}} + \underbrace{h(\vartheta_i) - \mathbb{E}[h(\vartheta) | U_i]}_{\text{fast}}.$$

- U_i is proposed according to AR scheme with persistency $\approx 1 - \psi N/T$: when N grows to slowly with T , IACT of $h(\vartheta_i)$ is dominated by IACT of $\mathbb{E}[h(\vartheta) | U_i]$ where for scalar $h(\vartheta) = \vartheta$

$$\mathbb{E}[\vartheta | U_i] = \underbrace{\hat{\theta}_T + \frac{\Sigma}{T} \nabla_{\theta} \log \frac{\hat{p}(Y_{1:T} | \hat{\theta}_T, U_i)}{p(Y_{1:T} | \hat{\theta}_T)}}_{\Psi(U_i)} + O_P(T^{-2}).$$

Decomposition in fast and slow components

- For a stationary CPM chain (ϑ_i, U_i) , decompose

$$h(\vartheta_i) = \underbrace{\mathbb{E}[h(\vartheta) | U_i]}_{\text{slow}} + \underbrace{h(\vartheta_i) - \mathbb{E}[h(\vartheta) | U_i]}_{\text{fast}}.$$

- U_i is proposed according to AR scheme with persistency $\approx 1 - \psi N/T$: when N grows to slowly with T , IACT of $h(\vartheta_i)$ is dominated by IACT of $\mathbb{E}[h(\vartheta) | U_i]$ where for scalar $h(\vartheta) = \vartheta$

$$\mathbb{E}[\vartheta | U_i] = \underbrace{\hat{\theta}_T + \frac{\Sigma}{T} \nabla_{\theta} \log \frac{\hat{p}(Y_{1:T} | \hat{\theta}_T, U_i)}{p(Y_{1:T} | \hat{\theta}_T)}}_{\Psi(U_i)} + O_P(T^{-2}).$$

- Proposition.** Let $N \propto T^\alpha$ for $0 < \alpha < 1$ then $\text{IACT}(Q, \Psi) \gtrsim T^{1-2\alpha}$.

Decomposition in fast and slow components

- For a stationary CPM chain (ϑ_i, U_i) , decompose

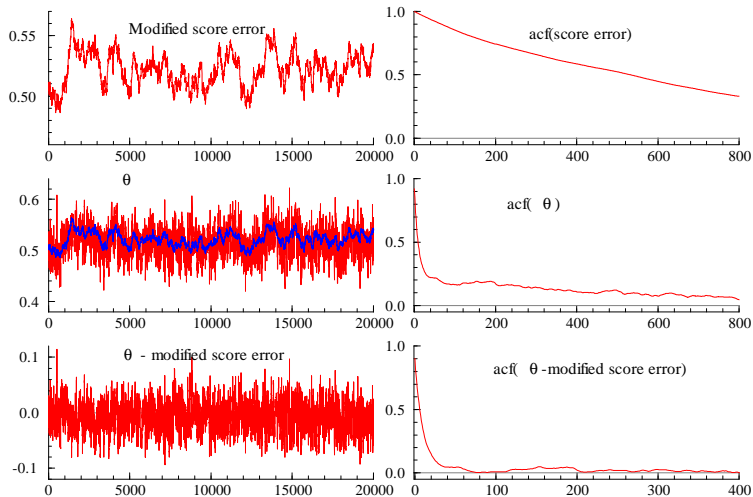
$$h(\vartheta_i) = \underbrace{\mathbb{E}[h(\vartheta) | U_i]}_{\text{slow}} + \underbrace{h(\vartheta_i) - \mathbb{E}[h(\vartheta) | U_i]}_{\text{fast}}.$$

- U_i is proposed according to AR scheme with persistency $\approx 1 - \psi N/T$: when N grows to slowly with T , IACT of $h(\vartheta_i)$ is dominated by IACT of $\mathbb{E}[h(\vartheta) | U_i]$ where for scalar $h(\vartheta) = \vartheta$

$$\mathbb{E}[\vartheta | U_i] = \underbrace{\hat{\theta}_T + \frac{\Sigma}{T} \nabla_{\theta} \log \frac{\hat{p}(Y_{1:T} | \hat{\theta}_T, U_i)}{p(Y_{1:T} | \hat{\theta}_T)}}_{\Psi(U_i)} + O_P(T^{-2}).$$

- **Proposition.** Let $N \propto T^\alpha$ for $0 < \alpha < 1$ then $\text{IACT}(Q, \Psi) \gtrsim T^{1-2\alpha}$.
- This result suggests we need at least $\sqrt{T}/N = O(1)$.

Graphical Illustration



Example: Gaussian Latent Variable Model

MH ($T = 8192$)		IACT(θ)	
		15.6	
PM ($\rho = 0.0$)			
N		RIACT(θ)	RCT(θ)
5000		2.2	11210
CPM ($\rho = 0.9963$)			
N	κ	RIACT(θ)	RCT(θ)
10	3.1	14.0	126.2
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

Here $\text{RIACT} = \text{IACT} / \text{IACT}_{MH}$ and $\text{RCT} = N \times \text{RIACT}$. Improvement by 180 fold.

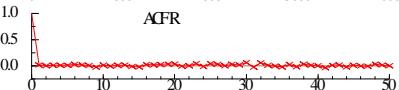
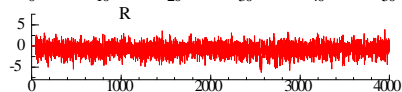
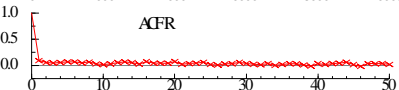
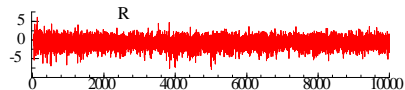
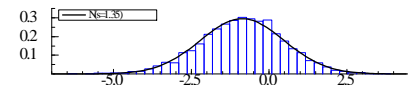
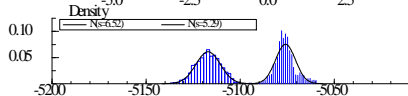
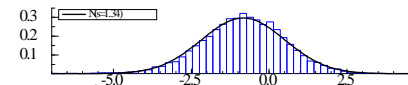
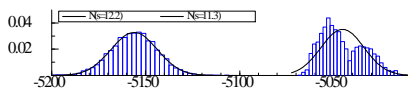
Example: Noisy Autoregressive Model

MH ($T = 16,000$)		IACT(θ)	
		5.8	
PM ($\rho = 0.0$)			
N		RIACT(θ)	RCT(θ)
2500		3.1	8427.0
CPM ($\rho = 0.9965$)			
N	κ	RIACT(θ)	RCT(θ)
6	6.7	43.8	262.8
10	3.3	8.7	86.7
16	1.9	6.0	85.8
22	1.3	3.9	85.6
35	0.8	2.4	85.0
40	0.7	2.4	94.8

Improvement by 100 fold.

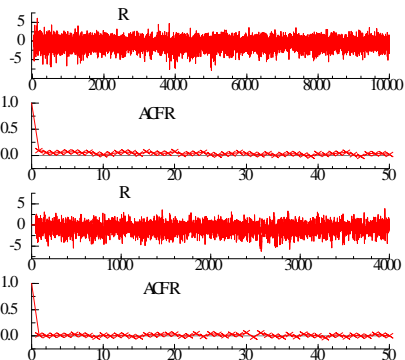
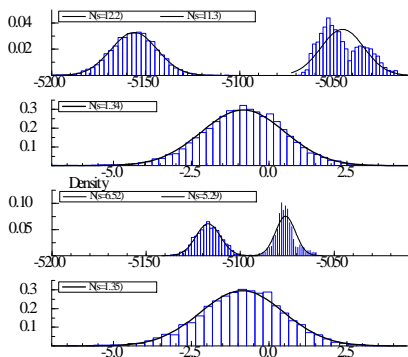
Heston Stochastic Volatility Model

- Inference for a discretized one-dimensional SDE, 40,000 latent variables. 100-fold gain compared to PM.



Heston Stochastic Volatility Model

- Inference for a discretized one-dimensional SDE, 40,000 latent variables. 100-fold gain compared to PM.
- Real data: 4,000 returns from the S&P 500 index from 15/08/1990 to 03/07/2006.



- CLT appears to hold only for N at least of order $T^{n/n+1}$ where n is state-dimension.

Higher-dimensional SSM

- CLT appears to hold only for N at least of order $T^{n/n+1}$ where n is state-dimension.
- Long-range effect is consequently limited.

Higher-dimensional SSM

- CLT appears to hold only for N at least of order $T^{n/n+1}$ where n is state-dimension.
- Long-range effect is consequently limited.
- Still significant gains over PM: over 50 fold for 2-d complex SV model, over 70-fold for 4-d model in (Jacob et al., 2016).

- Large sample analysis of pseudo-marginal algorithm provides useful guidelines, overall complexity is $O(T^2)$.

- Large sample analysis of pseudo-marginal algorithm provides useful guidelines, overall complexity is $O(T^2)$.
- Correlated pseudo-marginal can achieve very substantial improvement.

- Large sample analysis of pseudo-marginal algorithm provides useful guidelines, overall complexity is $O(T^2)$.
- Correlated pseudo-marginal can achieve very substantial improvement.
- In i.i.d. case, analysis shows $\sqrt{T}/N = O(1)$ is necessary and we conjecture it is sufficient leading to complexity $O(T^{3/2})$ vs $O(T^2)$.

- Large sample analysis of pseudo-marginal algorithm provides useful guidelines, overall complexity is $O(T^2)$.
- Correlated pseudo-marginal can achieve very substantial improvement.
- In i.i.d. case, analysis shows $\sqrt{T}/N = O(1)$ is necessary and we conjecture it is sufficient leading to complexity $O(T^{3/2})$ vs $O(T^2)$.
- Implementation for state-space models in state dimension $n > 1$ relies on non-standard particle scheme (e.g., Gerber & Chopin, 2015): our analysis does not capture these cases, experimental results suggest $O(T^{1+\frac{n}{n+1}})$.

Some References

- C. Andrieu, A.D. & R. Holenstein, “Particle Markov chain Monte Carlo methods”, *JRSS B*, 2000.
- J. Berard, P. Del Moral & A.D., “A Lognormal CLT for Particle Approximations of Normalizing Constants”, *Elec. J. Proba.*, 2014.
- G. Deligiannidis, A.D. and M.K. Pitt, “The Correlated Pseudo-marginal Method”, arXiv:1511.04992, 2015.
- 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.