Correlated Particle Metropolis

Michael Pitt King's College London Joint work George Deligiannidis & Arnaud Doucet (Oxford)

IMS, Singapore 2018

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

∃ ► < ∃ ►</p>

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

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

- Observations $\{Y_t\}_{t \ge 1}$ are conditionally independent given $\{X_t\}_{t \ge 0}$.
- Likelihood of $y_{1:T} = (y_1, ..., y_T)$ is

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

<ロ> <四> <四> <四> <三< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< =< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=< => <=<

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

A directed acyclic graph (DAG) of the problem:



イロト イヨト イヨト イヨト

- 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)

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

4 3 4 3 4

- 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, \ldots, x_t^N with associated discrete probability masses π_t^1, \ldots, π_t^N , drawn from the density $f(x_t|y_{1:t})$.

Bootrap Filter: GSS (1993, IEE)

We start at t = 0 with samples from $x_0^k \sim p(x_0)$. For t=1,...,T: We have samples $x_t^k \sim p(x_t|y_{1:t})$ for k = 1, ..., N. Solution For k = 1 : N, sample $\tilde{x}_{t+1}^k \sim f(x_{t+1}|x_t^k)$. For k = 1 : N, $\pi_{t+1}^k = \frac{g(y_{t+1}|\tilde{x}_{t+1}^k)}{\sum_{i=1}^N g(y_{t+1}|\tilde{x}_{t+1}^i)}$. 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;

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

• We need to estimate the function :

$$\widehat{p}(y_{1:T}|\theta) = \widehat{p}(y_1|\theta) \prod_{t=1}^{T-1} \widehat{p}(y_{t+1}|y_{1:t};\theta),$$
$$\widehat{p}(y_{t+1}|\theta;y_{1:t}) = \frac{1}{N} \sum_{i=1}^{N} p(y_{t+1}|\widetilde{x}_{t+1}^i).$$

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

・ 聞 と ・ 臣 と ・ 臣 と …

- 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.

• 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

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

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

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

3 1 4 3 1

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

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_θ (x_{1:T} | y_{1:T}); e.g. non-Gaussian state-space models.

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

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_θ (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

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

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_θ (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.

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

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_θ (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 $p(\theta, x_{1:T}|y_{1:T})$.

▲口> ▲圖> ▲注> ▲注> 三注

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.

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

ullet Introduce a target distribution on $\Theta\times {\mathbb U}$ of density

٠

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

Then unbiasedness yields

$$\int_{\mathbb{U}}\overline{\pi}\left(heta,u
ight)\mathsf{d}u=\pi\left(heta
ight)$$

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

Pseudo-Marginal Metropolis-Hastings algorithm

Can form an unbiased estimator, based on N particles p̂ (y; θ, U).
Set (θ⁽⁰⁾, U⁽⁰⁾) and iterate for j = 1, 2, ...
Sample ϑ ~ q(·|ϑ^(j-1)), U ~ mϑ(·) to obtain p̂ (y; ϑ, U).
Compute

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

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

• Standard non-linear model

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

$$Y_t = rac{1}{20}X_t^2 + W_t$$
, $W_t \stackrel{ ext{i.i.d.}}{\sim} \mathcal{N}\left(0, \sigma_W^2
ight)$.

- 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}, θ) is highly multimodal.
- We sample from p (θ| y_{1:T}) using a random walk pseudo-marginal MH where p (y_{1:T}; θ) 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.

- 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

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

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

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

• Aim: Minimize the "computational time"

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

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

• Aim: Minimize the "computational time"

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

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

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

- 4 回 ト 4 ヨ ト - 4 ヨ ト -

• The IACT is

$$\mathrm{IACT}\left(\textit{\textit{Q}},\textit{\textit{h}}\right) = 1 + 2\sum_{\tau=1}^{\infty} \mathsf{corr}_{\overline{\pi},\textit{\textit{Q}}}\left\{h\left(\theta_{0}\right), h\left(\theta_{\tau}\right)\right\}$$

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

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

- where $r\left(heta, artheta
 ight)$ is exact Metropolis ratio,
- and

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

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

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

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

イロト 不得下 イヨト イヨト

IID data, fixed θ (IS)

$$Z = \log \widehat{p}(y| heta) - \log p(y| heta) = \sum_{t=1}^{T} \log \left\{ 1 + rac{\gamma_t}{\sqrt{N}} \varepsilon_t
ight\}$$
 ,

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

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

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

- More detail (and more generality, PFs) in Bérard, Del Moral, Doucet (EJP, 2014).
- Need concentration (in θ as $T \longrightarrow \infty$), 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 $q(\vartheta|\theta) = \pi(\vartheta)$, $\sigma_{opt} = 0.92$ (Pitt et al., JoE 2012). The acf simply reduces to:

$$\phi_n(\theta, Q) = \int \Pr(R \mid z; \sigma)^n \pi_Z(z \mid \sigma) \, dz$$

$$IACT(\sigma) = \int \frac{1 + \Pr(R \mid z; \sigma)}{1 - \Pr(R \mid z; \sigma)} \pi_Z(z \mid \sigma) \, dz.$$

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

• Sherlock, Thiery, Roberts and Rosenthal, (Annals 2015) consider the (joint) optimisation for a limiting target in the important case of RWM. Again anytically 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.

- 4 同 6 4 日 6 4 日 6

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

$$egin{aligned} Q^{*}\left\{ \left(heta,z
ight),\left(dartheta,dw
ight)
ight\} &=q(artheta| heta)g(w)lpha_{\mathsf{EX}}\left(heta,artheta
ight)lpha_{\mathsf{Z}}\left(z,w
ight)dartheta dw\ &+\left\{1-arrho_{\mathsf{EX}}\left(heta
ight)arrho_{\mathsf{Z}}\left(z
ight)
ight\}\delta_{\left(heta,z
ight)}\left(dartheta,dw
ight), \end{aligned}$$

where we have a product acceptance criterion:

$$lpha_{\mathsf{EX}}\left(heta,artheta
ight)=\min\left\{1,r\left(heta,artheta
ight)
ight\}$$
, $lpha_{\mathsf{Z}}\left(z,w
ight)=\min\left\{1,\exp\left(w-z
ight)
ight\}$.

通 ト イヨ ト イヨ ト

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

$$\begin{aligned} Q^*\left\{\left(\theta,z\right),\left(d\vartheta,dw\right)\right\} &= q(\vartheta|\theta)g(w)\alpha_{\mathsf{EX}}\left(\theta,\vartheta\right)\alpha_{\mathsf{Z}}\left(z,w\right)d\vartheta dw \\ &+ \left\{1-\varrho_{\mathsf{EX}}\left(\theta\right)\varrho_{\mathsf{Z}}\left(z\right)\right\}\delta_{\left(\theta,z\right)}\left(d\vartheta,dw\right),\end{aligned}$$

where we have a product acceptance criterion:

 $\alpha_{\mathsf{EX}}\left(\theta,\vartheta\right) = \min\left\{1,r\left(\theta,\vartheta\right)\right\}, \ \ \alpha_{\mathsf{Z}}\left(z,w\right) = \min\left\{1,\exp\left(w-z\right)\right\}.$

• Peskun's theorem (1973) guarantees that $IF_{h}^{Q}(\sigma) \leq IF_{h}^{Q^{*}}(\sigma)$ so that $CT_{h}^{Q}(\sigma) \leq CT_{h}^{Q^{*}}(\sigma)$.

▲ロト ▲圖 ト ▲ ヨト ▲ ヨト 二 ヨー わえの

pseudo-proof: Jump chain

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

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

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

Simpler Bounds on the Relative Inefficiency

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

$$\frac{\mathit{IF}_{h}^{Q}\left(\sigma\right)}{\mathit{IF}_{h}^{\mathsf{EX}}} \leq \frac{\mathit{IF}_{h}^{Q^{*}}\left(\sigma\right)}{\mathit{IF}_{h}^{\mathsf{EX}}} \leq \frac{1}{2}\left(1 + \frac{1}{\mathit{IF}_{h}^{\mathsf{EX}}}\right)\pi_{\mathsf{Z}}^{\sigma}\left(1/\varrho_{\mathsf{Z}}^{\sigma}\right) - \frac{1}{\mathit{IF}_{h}^{\mathsf{EX}}}$$

and the bound is tight as $IF_h^{\mathsf{EX}} \to 1 \text{ or } \sigma \to 0$. • As $IF_{J,h/\varrho_{\mathsf{EX}}}^{\mathsf{EX}} \to \infty$,

$$\frac{\mathit{IF}_{h}^{Q^{*}}(\sigma)}{\mathit{IF}_{h}^{\mathsf{EX}}} \to \frac{1}{\pi^{\sigma}_{\mathsf{Z}}(\varrho^{\sigma}_{\mathsf{Z}})} = \frac{1}{\mathit{P}_{\mathsf{A}}(\sigma)}.$$

• Results used to minimize w.r.t σ upper bounds on $CT_{h}^{Q}(\sigma) = IF_{h}^{Q}(\sigma) / \sigma^{2}$.

・ロト ・聞 と ・ 臣 と ・ 臣 と … 臣

Bounds on Relative Computational Time



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

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

$$\begin{aligned} dv_{1}\left(t\right) &= -k_{1}\left\{v_{1}\left(t\right) - \mu_{1}\right\}dt + \sigma_{1}dW_{1}\left(t\right),\\ dv_{2}\left(t\right) &= -k_{2}v_{2}\left(t\right) + \left\{1 + \beta_{12}v_{2}\left(t\right)\right\}dW_{2}\left(t\right),\\ d\log P\left(t\right) &= \mu_{y}dt + \text{s-exp}\left[\left\{v_{1}\left(t\right) + \beta_{2}v_{2}\left(t\right)\right\}/2\right]dB\left(t\right), \end{aligned}$$

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

- 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.

イロト イヨト イヨト イヨト
- Daily returns $y = (y_1, ..., y_T)$ 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 $\overline{\theta}$.

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 $\overline{\theta}$ (left) and marginalized over $\pi(\theta)$.

(IMS, Singapore 2018)

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: 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 √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 √T; e.g. Lee, Econometric Theory, 1999.
- Problem: the ratio p_θ (y_{1:T}) / p_θ (y_{1:T}) is estimated by estimating independently the numerator and denominator in pseudo-marginal.

イロト イヨト イヨト イヨト

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 $p_{ heta}\left(y_{1:T}\right)$ and $p_{artheta}\left(y_{1:T}\right)$ by setting

$$\widehat{p}_{\vartheta}\left(y_{1:T}\right) = \widehat{p}_{\vartheta}\left(y_{1:T};V\right)$$

where

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

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

- < A > < B > < B >

- < A > < B > < B >

<u>At iteration i</u>

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

イロト イポト イヨト イヨト

<u>At iteration i</u>

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

<u>At iteration i</u>

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

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

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

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

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \widehat{p}_{\theta}(y; u)$ 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) .

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \widehat{p}_{\theta}(y; u)$ 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 $X_t^{\sigma_t(1)} \leq \ldots \leq X_t^{\sigma_t(N)}$.

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \widehat{p}_{\theta}(y; u)$ 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 $X_t^{\sigma_t(1)} \leq \ldots \leq X_t^{\sigma_t(N)}$.
- For $n \ge 2$, (Gerber & Chopin, 2015) use Hilbert curve; e.g. map particles to $[0,1]^n$ using $\psi : \mathbb{R}^n \to [0,1]^n$ (e.g. logistic) and $h : [0,1]^n \to [0,1]$ (pseudo-inverse Hilbert curve) and sort projected particles on [0,1].

・ロト ・聞 ト ・ 国 ト ・ 国 ト … 国

- If the likelihood is computed using standard SMC, $(\theta, u) \mapsto \widehat{p}_{\theta}(y; u)$ 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 $X_t^{\sigma_t(1)} \leq \ldots \leq X_t^{\sigma_t(N)}$.
- For $n \ge 2$, (Gerber & Chopin, 2015) use Hilbert curve; e.g. map particles to $[0,1]^n$ using $\psi : \mathbb{R}^n \to [0,1]^n$ (e.g. logistic) and $h: [0,1]^n \to [0,1]$ (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: $\exists \ \widehat{\theta}^T \xrightarrow{P} \overline{\theta}$ and Σ a p.d. matrix s.t.

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

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

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

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

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

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

where

$$\kappa^{2}(\theta) = 2\psi \mathbb{E}\left[\left\|\partial_{u}\mathcal{O}\left(U,Y;\theta\right)\right\|^{2}\right]$$

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

|山田 | 小田 | 小田 |

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

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

where

$$\kappa^{2}\left(\theta\right)=2\psi\mathbb{E}\left[\left\Vert \partial_{u}\mathcal{D}\left(U,Y;\theta
ight)
ight\Vert ^{2}
ight]$$

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

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

where

$$\kappa^{2}\left(\theta\right)=2\psi\mathbb{E}\left[\left\Vert \partial_{u}\mathcal{O}\left(U,Y;\theta
ight)
ight\Vert ^{2}
ight]$$

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

ヘロト 人間 とくほ とくほ とう

Let Θ_T := {ϑ_i^T}_{i≥0} be the stationary *non-Markovian* sequence of the correlated PM targetting p (θ| Y_{1:T}).

- Let Θ_T := {ϑ_i^T}_{i≥0} be the stationary *non-Markovian* sequence of the correlated PM targetting p (θ | Y_{1:T}).
- **Proposition** (Deligiannidis et al., 2016): The sequences $\{\Theta_T\}_{T\geq 1}$ converge weakly as $T \to \infty$ to a stationary Markov chain of invariant density $\phi\left(\tilde{\theta}; 0, \Sigma\right)$ and kernel given for $\tilde{\theta} \neq \tilde{\vartheta}$ by

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

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

- Let Θ_T := {ϑ_i^T}_{i≥0} be the stationary *non-Markovian* sequence of the correlated PM targetting p (θ | Y_{1:T}).
- **Proposition** (Deligiannidis et al., 2016): The sequences $\{\Theta_T\}_{T\geq 1}$ converge weakly as $T \to \infty$ to a stationary Markov chain of invariant density $\phi\left(\tilde{\theta}; 0, \Sigma\right)$ and kernel given for $\tilde{\theta} \neq \tilde{\vartheta}$ by

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

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

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

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

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

Image: Image:

- 4 3 6 4 3 6

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

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

• 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 $h(\vartheta) = \vartheta$

$$\mathbb{E}\left[\left.\vartheta\right|\left.U_{i}\right]=\underbrace{\widehat{\theta}_{T}+\frac{\Sigma}{T}\nabla_{\theta}\log\frac{\widehat{p}\left(\left.Y_{1:T}\right|\widehat{\theta}_{T},U_{i}\right)}{p\left(\left.Y_{1:T}\right|\widehat{\theta}_{T}\right)}}_{\Psi\left(U_{i}\right)}+O_{P}\left(\left.T^{-2}\right).$$

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

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

• 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 $h(\vartheta) = \vartheta$

$$\mathbb{E}\left[\left.\vartheta\right|\left.U_{i}\right]=\widehat{\theta}_{T}+\frac{\Sigma}{T}\nabla_{\theta}\log\frac{\widehat{p}\left(\left.Y_{1:T}\right|\widehat{\theta}_{T},U_{i}\right)}{p\left(\left.Y_{1:T}\right|\widehat{\theta}_{T}\right)}+O_{P}\left(T^{-2}\right).$$

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

- 4 回 ト 4 ヨ ト - 4 ヨ ト -

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

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

• 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 $h(\vartheta) = \vartheta$

$$\mathbb{E}\left[\left.\vartheta\right|\left.U_{i}\right]=\underbrace{\widehat{\theta}_{\mathcal{T}}+\frac{\Sigma}{\mathcal{T}}\nabla_{\theta}\log\frac{\widehat{p}\left(\left.Y_{1:\mathcal{T}}\right|\widehat{\theta}_{\mathcal{T}},U_{i}\right)}{p\left(\left.Y_{1:\mathcal{T}}\right|\widehat{\theta}_{\mathcal{T}}\right)}+O_{P}\left(\left.\mathcal{T}^{-2}\right).}_{\Psi\left(U_{i}\right)}$$

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

• This result suggests we need at least $\sqrt{T}/N = O(1)$.

Graphical Illustration



< A

Example: Gaussian Latent Variable Model

MH (<i>T</i> = 8192)		$IACT(\theta)$	
		15.6	
PM ($ ho=$ 0.0)			
Ν		$RIACT(\theta)$	$RCT(\theta)$
5000		2.2	11210
CPM ($ ho = 0.9963$)			
N	κ	$RIACT(\theta)$	$RCT(\theta)$
N 10	κ 3.1	RIACT(θ) 14.0	RCT(θ) 126.2
N 10 20	κ 3.1 2.2	RIACT(θ) 14.0 4.7	RCT(θ) 126.2 93.3
N 10 20 25	κ 3.1 2.2 2.0	RIACT(θ) 14.0 4.7 2.8	RCT(θ) 126.2 93.3 69.3
N 10 20 25 35	κ 3.1 2.2 2.0 1.7	RIACT(θ) 14.0 4.7 2.8 1.7	RCT(θ) 126.2 93.3 69.3 61.1

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

イロン イ理と イヨン イヨン

Example: Noisy Autoregressive Model

MH ($T = 16,000$)		$IACT(\theta)$	
		5.8	
PM ($ ho=$ 0.0)			
N		$RIACT(\theta)$	$RCT(\theta)$
2500		3.1	8427.0
$CPM(\rho=0.9965)$			
N	κ	$RIACT(\theta)$	$RCT(\theta)$
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.

44 / 48

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.

A B F A B F

- 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.

∃ ► < ∃ ►</p>

- 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 $O(T^2)$.

個 ト イヨト イヨト

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

▶ < ∃ ▶ < ∃ ▶</p>

- 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+ n/n+1}).

- 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.