

# 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)

MCQMC, 19/08/2016

# 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

# Latent Variable Models

- Assume  $(Y_t)_{t \geq 1}$  are i.i.d. random variables such that

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

where  $(X_t)_{t \geq 1}$  are latent variables.

# Latent Variable Models

- Assume  $(Y_t)_{t \geq 1}$  are i.i.d. random variables such that

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

where  $(X_t)_{t \geq 1}$  are latent variables.

- The likelihood of  $\theta \in \Theta \subseteq \mathbb{R}^d$  associated to observations  $Y_{1:T} = y_{1:T}$  is

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



- Assume  $(Y_t)_{t \geq 1}$  are i.i.d. random variables such that

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

where  $(X_t)_{t \geq 1}$  are latent variables.

- The likelihood of  $\theta \in \Theta \subseteq \mathbb{R}^d$  associated to observations  $Y_{1:T} = y_{1:T}$  is

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

- In many scenarios,  $p_\theta(y_{1:T})$  cannot be evaluated exactly; e.g. multivariate probit model.

# Latent Variable Models

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

# Latent Variable 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  $\theta \in \mathbb{R}^d$  associated to  $Y_{1:T} = y_{1:T}$  is

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

where

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

# Latent Variable 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  $\theta \in \mathbb{R}^d$  associated to  $Y_{1:T} = y_{1:T}$  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 a very popular class of time series models but inference is difficult as  $p_\theta(y_{1:T})$  is intractable for non-linear/non-Gaussian models.

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

# Bayesian Inference

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

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

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

- 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

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

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

# Standard MCMC Approaches

- Standard MCMC schemes target

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

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 mimic an algorithm targetting directly  $p(\theta | y_{1:T})$  instead of  $p(\theta, x_{1:T} | y_{1:T})$ .

# Ideal Marginal Metropolis-Hastings

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

# Ideal Marginal Metropolis-Hastings

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

# Ideal Marginal Metropolis-Hastings

- 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

- 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

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

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



# Ideal Marginal Metropolis-Hastings

- 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

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

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

- **Idea:** Replace  $p_\theta(y_{1:T})$  by a non-negative unbiased estimate  $\hat{p}_\theta(y_{1:T})$  in MH.

# Pseudo-Marginal Metropolis–Hastings

- **Idea:** Replace  $p_\theta(y_{1:T})$  by a non-negative unbiased estimate  $\hat{p}_\theta(y_{1:T})$  in MH.

# Pseudo-Marginal Metropolis–Hastings

- **Idea:** Replace  $p_{\vartheta}(y_{1:T})$  by a non-negative unbiased estimate  $\hat{p}_{\vartheta}(y_{1:T})$  in MH.

At iteration  $i$

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

# Pseudo-Marginal Metropolis–Hastings

- **Idea:** Replace  $p_{\vartheta}(y_{1:T})$  by a non-negative unbiased 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

- **Idea:** Replace  $p_{\vartheta}(y_{1:T})$  by a non-negative unbiased 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  $p_{\vartheta}(y_{1:T})$ .
- With probability

$$\begin{aligned} & 1 \wedge \frac{\hat{p}_{\vartheta}(y_{1:T}) p(\vartheta)}{\hat{p}_{\vartheta_{i-1}}(y_{1:T}) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} \\ = & 1 \wedge \frac{p_{\vartheta}(y_{1:T}) p(\vartheta)}{p_{\vartheta_{i-1}}(y_{1:T}) p(\vartheta_{i-1})} \frac{q(\vartheta_{i-1} | \vartheta)}{q(\vartheta | \vartheta_{i-1})} \frac{\hat{p}_{\vartheta}(y_{1:T}) / p_{\vartheta}(y_{1:T})}{\hat{p}_{\vartheta_{i-1}}(y_{1:T}) / p_{\vartheta_{i-1}}(y_{1:T})}, \end{aligned}$$

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

- **Proposition** (Lin, Liu & Sloan, *Phys. Rev. D*, 2000): 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.

# Pseudo-Marginal Metropolis–Hastings

- **Proposition** (Lin, Liu & Sloan, *Phys. Rev. D*, 2000): 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  $\mathcal{U}$ -valued 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  $U \sim m_\theta(\cdot)$ .



# Pseudo-Marginal Metropolis–Hastings

- **Proposition** (Lin, Liu & Sloan, *Phys. Rev. D*, 2000): 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  $\mathcal{U}$ -valued 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  $U \sim m_\theta(\cdot)$ .
- Pseudo-marginal MH is a standard MH on  $\Theta \times \mathcal{U}$  with target  $\bar{\pi}(\theta, u)$  and proposal  $q(\vartheta | \theta) m_\vartheta(v)$  where

$$\bar{\pi}(\theta, u) = \pi(\theta) \frac{\hat{p}_\theta(y_{1:T}; u)}{p_\theta(y_{1:T})} m_\theta(u)$$

satisfies

$$\int \bar{\pi}(\theta, u) du = \pi(\theta).$$

# Likelihood Estimators

- For latent variable models, one has

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

- For latent variable models, one has

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

- A 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.  $U = (X_1^1, \dots, X_1^N, \dots, X_T^1, \dots, X_T^N)$  with

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

- For latent variable models, one has

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

- A 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.  $U = (X_1^1, \dots, X_1^N, \dots, X_T^1, \dots, X_T^N)$  with

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

- Computational complexity is  $O(NT)$ .

- For state-space models, an alternative is to use particle filters where

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

where

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

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

- For state-space models, an alternative is to use particle filters where

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

where

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

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

- The estimator is unbiased (Del Moral, 1998), relative variance is bounded uniformly over  $T$  if  $N \propto T$  (Cerou, Del Moral & Guyader, 2011), SQMC could also be used (Gerber & Chopin, *JRSS B* 2015).

- For state-space models, an alternative is to use particle filters where

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

where

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

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

- The estimator is unbiased (Del Moral, 1998), relative variance is bounded uniformly over  $T$  if  $N \propto T$  (Cerou, Del Moral & Guyader, 2011), SQMC could also be used (Gerber & Chopin, *JRSS B* 2015).
- Computational complexity is  $O(NT)$ .

# Examples of recent applications

- Measuring the impact of Ebola control measures in Sierra Leone, *PNAS*, 2016.



# Examples of recent applications

- Measuring the impact of Ebola control measures in Sierra Leone, *PNAS*, 2016.
- Dynamic prediction pools: an investigation of financial frictions and forecasting performance, *J. Econometrics*, 2016.

# Examples of recent applications

- Measuring the impact of Ebola control measures in Sierra Leone, *PNAS*, 2016.
- Dynamic prediction pools: an investigation of financial frictions and forecasting performance, *J. Econometrics*, 2016.
- Capturing the dynamics of pathogens with many strains, *J. Mathematical Biology*, 2016.

# Examples of recent applications

- Measuring the impact of Ebola control measures in Sierra Leone, *PNAS*, 2016.
- Dynamic prediction pools: an investigation of financial frictions and forecasting performance, *J. Econometrics*, 2016.
- Capturing the dynamics of pathogens with many strains, *J. Mathematical Biology*, 2016.
- Monte Carlo estimation of stage structured development from cohort data, *Ecology*, 2016.

# Examples of recent applications

- Measuring the impact of Ebola control measures in Sierra Leone, *PNAS*, 2016.
- Dynamic prediction pools: an investigation of financial frictions and forecasting performance, *J. Econometrics*, 2016.
- Capturing the dynamics of pathogens with many strains, *J. Mathematical Biology*, 2016.
- Monte Carlo estimation of stage structured development from cohort data, *Ecology*, 2016.
- Bayesian analysis of entry games: A simulated likelihood approach without simulation errors, *Global Economic Review*, 2016.

# Examples of recent applications

- Measuring the impact of Ebola control measures in Sierra Leone, *PNAS*, 2016.
- Dynamic prediction pools: an investigation of financial frictions and forecasting performance, *J. Econometrics*, 2016.
- Capturing the dynamics of pathogens with many strains, *J. Mathematical Biology*, 2016.
- Monte Carlo estimation of stage structured development from cohort data, *Ecology*, 2016.
- Bayesian analysis of entry games: A simulated likelihood approach without simulation errors, *Global Economic Review*, 2016.
- Controlling procedural modeling programs. *SIGGRAPH*, 2015.

# Empirical performance: Stochastic volatility model

- SV model with leverage (Huang & Tauchen, 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  $\phi_2 = \text{corr}\{B(s), W_2(s)\}$ .

# Empirical performance: Stochastic volatility model

- SV model with leverage (Huang & Tauchen, 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  $Y_t = \log P(\Delta t) - \log P(\Delta(t-1))$ .

# Empirical performance: Stochastic volatility model

- SV model with leverage (Huang & Tauchen, 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.



- SV model with leverage (Huang & Tauchen, 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  $\theta = (k_1, \mu_1, \sigma_1, k_2, \beta_{12}, \beta_2, \mu_y, \phi_1, \phi_2)$ .

# Empirical performance: Stochastic volatility model

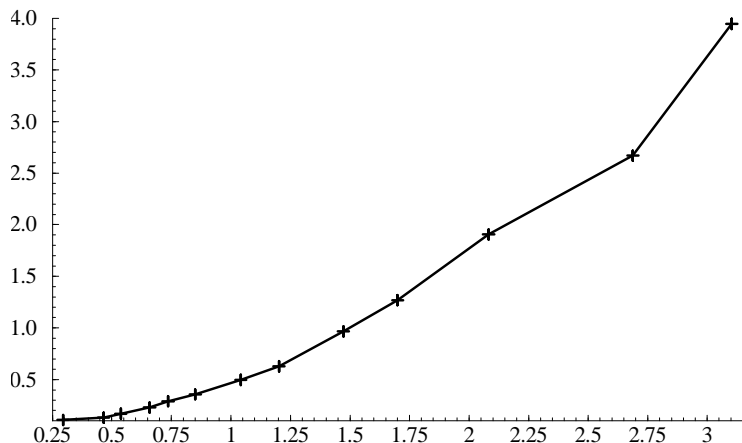
- SV model with leverage (Huang & Tauchen, 2005):

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

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 MH for RW proposal w.r.t  $\sigma$ , 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  $\sigma$  for  $T = 300$ .

# 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  $N$ , where

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

# How precise should the log-likelihood estimator be?

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

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

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

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

- The IACT is

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

where, for  $(\theta, z) \neq (\vartheta, w)$ , the PM kernel  $Q$  is

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

with

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

# Computational time for the SV model

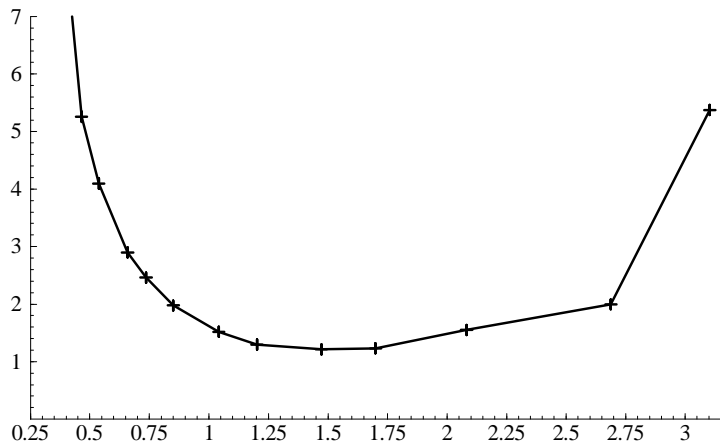


Figure: Computational time as a function of  $\sigma$

# Large sample analysis of the pseudo-marginal kernel

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

# Large sample analysis of the pseudo-marginal kernel

- 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  $\Sigma$  is a p.d. matrix.



# Large sample analysis of the pseudo-marginal kernel

- 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  $\Sigma$  is a p.d. matrix.

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

# Large sample analysis of the pseudo-marginal kernel

- 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  $\Sigma$  is a p.d. matrix.

- **Assumption 2** - *Proposal*:  $\vartheta = \theta + \varepsilon/\sqrt{T}$  where  $\varepsilon \sim v(\cdot)$  with  $v(\varepsilon) = v(-\varepsilon)$ .
- **Proposition** (Berard, Del Moral & D., *EJP*, 2014): Under regularity conditions, we have if  $N \propto T$

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

# Large sample analysis of the pseudo-marginal kernel

- 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 targetting  $\pi^T(\theta, z) = p(\theta | Y_{1:T}) \exp(z) g_\theta^T(z)$ .

# Large sample analysis of the pseudo-marginal kernel

- 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 targetting  $\pi^T(\theta, z) = p(\theta | Y_{1:T}) \exp(z) g_\theta^T(z)$ .
- **Proposition:** The F.D.D. of  $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \widehat{\theta}_T), Z_i^T\}_{i \geq 0}$  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

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

# Large sample analysis of the pseudo-marginal kernel

- 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 targetting  $\pi^T(\theta, z) = p(\theta | Y_{1:T}) \exp(z) g_\theta^T(z)$ .
- **Proposition:** The F.D.D. of  $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \widehat{\theta}_T), Z_i^T\}_{i \geq 0}$  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

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

- Let  $\sigma^2 = \sigma^2(\bar{\theta})$ , it suggests a simplified analysis based on

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

- **Aim:** Minimize the computational cost

$$\text{CT}(\hat{Q}_\sigma, h) = \frac{\text{IACT}(\hat{Q}_\sigma, h)}{\sigma^2}.$$

- **Aim:** Minimize the computational cost

$$\text{CT}(\hat{Q}_\sigma, h) = \frac{\text{IACT}(\hat{Q}_\sigma, h)}{\sigma^2}.$$

- **Special cases:**

- **Aim:** Minimize the computational cost

$$\text{CT}(\hat{Q}_\sigma, h) = \frac{\text{IACT}(\hat{Q}_\sigma, h)}{\sigma^2}.$$

- **Special cases:**
- ① When  $q(\vartheta|\theta) = p(\vartheta|y)$ ,  $\sigma_{\text{opt}} = 0.9$  (Pitt et al., 2012).



- **Aim:** Minimize the computational cost

$$\text{CT}(\hat{Q}_\sigma, h) = \frac{\text{IACT}(\hat{Q}_\sigma, h)}{\sigma^2}.$$

- **Special cases:**
  - 1 When  $q(\vartheta|\theta) = p(\vartheta|y)$ ,  $\sigma_{\text{opt}} = 0.9$  (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  $\sigma_{\text{opt}} = 1.8$  (Sherlock et al., 2015).

# Sketch of the Analysis

- For general proposals and targets, direct minimization of  $\text{CT}(\hat{Q}_\sigma, h)$  impossible so minimize an upper bound over it.

# Sketch of the Analysis

- For general proposals and targets, direct minimization of  $\text{CT}(\widehat{Q}_\sigma, h)$  impossible so minimize an upper bound over it.
- Theoretical study relies on  $\bar{\pi}$ -invariant kernel  $Q_\sigma^*$  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  $\text{CT}(\widehat{Q}_\sigma, h)$  impossible so minimize an upper bound over it.
- Theoretical study relies on  $\bar{\pi}$ -invariant kernel  $Q_\sigma^*$  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

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

- We have  $\text{IACT}(\widehat{Q}_\sigma, h) \leq \text{IACT}(Q_\sigma^*, h)$  by Peskun (1973) so  $\text{CT}(\widehat{Q}_\sigma, h) \leq \text{CT}(Q_\sigma^*, h)$  and exact expression of  $\text{IACT}(Q_\sigma^*, h)$  is available.

# 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  $\approx 15\%$ .

- Many sharp qualitative results for the pseudo-marginal MH algorithm have been obtained (Andrieu & Roberts, AoS 2009; Andrieu & Vihola, AAP 2015).

- Many sharp qualitative results for the pseudo-marginal MH algorithm have been obtained (Andrieu & Roberts, AoS 2009; Andrieu & Vihola, AAP 2015).
- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.

- Many sharp qualitative results for the pseudo-marginal MH algorithm have been obtained (Andrieu & Roberts, AoS 2009; Andrieu & Vihola, AAP 2015).
- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.
- Optimal  $\sigma$  depends on efficiency of the ideal MH algorithm but  $\sigma \approx 1.2 - 1.3$  is a sweet spot.

- Many sharp qualitative results for the pseudo-marginal MH algorithm have been obtained (Andrieu & Roberts, AoS 2009; Andrieu & Vihola, AAP 2015).
- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.
- Optimal  $\sigma$  depends on efficiency of the ideal MH algorithm but  $\sigma \approx 1.2 - 1.3$  is a sweet spot.
- Pseudo-marginal MH scales in  $\mathcal{O}(T^2)$  at each iteration as we require  $N \propto T$ .

- Many sharp qualitative results for the pseudo-marginal MH algorithm have been obtained (Andrieu & Roberts, AoS 2009; Andrieu & Vihola, AAP 2015).
- Simplified quantitative analysis of the pseudo-marginal MH algorithm, useful in large data regime.
- Optimal  $\sigma$  depends on efficiency of the ideal MH algorithm but  $\sigma \approx 1.2 - 1.3$  is a sweet spot.
- Pseudo-marginal MH scales in  $\mathcal{O}(T^2)$  at each iteration as we require  $N \propto T$ .
- **Problem:** the ratio  $p_\theta(y_{1:T}) / p_\theta(y_{1:T})$  is estimated by estimating independently the numerator and denominator.

# 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,  $\rho$  will be selected close to 1.

# Correlated Pseudo-Marginal Metropolis–Hastings algorithm

At iteration  $i$

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

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  $\vartheta_i = \vartheta_{i-1}$ ,  $U_i = U_{i-1}$ .

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

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

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

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

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

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

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

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

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



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

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

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

- 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.** The F.D.D. of  $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \hat{\theta}_T)\}_{i \geq 0}$  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

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

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

- 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.** The F.D.D. of  $\{\tilde{\vartheta}_i^T = \sqrt{T}(\vartheta_i^T - \hat{\theta}_T)\}_{i \geq 0}$  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

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

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

$$\hat{Q}(\theta, d\vartheta) = q(\vartheta | \theta) \mathbb{E}_{R \sim \mathcal{N}(-\kappa^2(\tilde{\theta})/2, \kappa^2(\tilde{\theta}))} \left[ \min \left\{ 1, \frac{\pi(\vartheta)}{\pi(\theta)} R \right\} \right] d\vartheta$$

and  $\kappa \approx 2.1$  minimizes computational time.

# Limitations of Weak Convergence

- Previous result does not imply convergence of IACT.

# Limitations of Weak Convergence

- Previous result does not imply convergence of IACT.
- In the CPM context, the auxiliary variables  $U^T$  evolves according to an AR scheme with persistency  $\approx 1 - N_T / T$ .

# Limitations of Weak Convergence

- Previous result does not imply convergence of IACT.
- In the CPM context, the auxiliary variables  $U^T$  evolves according to an AR scheme with persistency  $\approx 1 - N_T / T$ .
- When  $N_T$  grows too slowly with  $T$ , correlations decay too slowly leading to increasing IACT.

# Limitations of Weak Convergence

- Previous result does not imply convergence of IACT.
- In the CPM context, the auxiliary variables  $U^T$  evolves according to an AR scheme with persistency  $\approx 1 - N_T / T$ .
- When  $N_T$  grows too slowly with  $T$ , correlations decay too slowly leading to increasing IACT.
- For the i.i.d. latent variable case,  $\text{IACT}(Q, h)$  is driven by the IACT of  $\{\Psi_T(U_i^T) := \nabla_{\theta} \log \hat{p}(Y_{1:T}; U_i^T)\}_{i \geq 1}$ .

# Limitations of Weak Convergence

- Previous result does not imply convergence of IACT.
- In the CPM context, the auxiliary variables  $U^T$  evolves according to an AR scheme with persistency  $\approx 1 - N_T / T$ .
- When  $N_T$  grows too slowly with  $T$ , correlations decay too slowly leading to increasing IACT.
- For the i.i.d. latent variable case,  $\text{IACT}(Q, h)$  is driven by the IACT of  $\{\Psi_T(U_i^T) := \nabla_{\theta} \log \hat{p}(Y_{1:T}; U_i^T)\}_{i \geq 1}$ .
- **Proposition.** Let  $N_T \propto T^{\alpha}$  for  $0 < \alpha < 1$  then  $\text{IACT}(Q, h) \gtrsim T^{1-2\alpha}$ .



# Limitations of Weak Convergence

- Previous result does not imply convergence of IACT.
- In the CPM context, the auxiliary variables  $U^T$  evolves according to an AR scheme with persistency  $\approx 1 - N_T / T$ .
- When  $N_T$  grows too slowly with  $T$ , correlations decay too slowly leading to increasing IACT.
- For the i.i.d. latent variable case,  $\text{IACT}(Q, h)$  is driven by the IACT of  $\{\Psi_T(U_i^T) := \nabla_{\theta} \log \hat{p}(Y_{1:T}; U_i^T)\}_{i \geq 1}$ .
- **Proposition.** Let  $N_T \propto T^{\alpha}$  for  $0 < \alpha < 1$  then  $\text{IACT}(Q, h) \gtrsim T^{1-2\alpha}$ .
- This result suggests we need at least  $\sqrt{T} / N_T = O(1)$ .

# Example: Gaussian Latent Variable Model

MH ( $T = 8192$ )		IACT( $\theta$ )	
		15.6	
PM ( $\rho = 0.0$ )			
$N$		RIACT( $\theta$ )	RCT( $\theta$ )
5000		2.2	11210
CPM ( $\rho = 0.9963$ )			
$N$	$\kappa$	RIACT( $\theta$ )	RCT( $\theta$ )
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( $\theta$ )	
		5.8	
PM ( $\rho = 0.0$ )			
$N$		RIACT( $\theta$ )	RCT( $\theta$ )
2500		3.1	8427.0
CPM ( $\rho = 0.9965$ )			
$N$	$\kappa$	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.

- 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.
- Analysis shows  $\sqrt{T}/N_T = O(1)$  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.
- Analysis shows  $\sqrt{T}/N_T = O(1)$  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 (Gerber & Chopin, 2015): our analysis does not capture these cases, experimental results suggest  $O(T^{1+\frac{n}{n+1}})$ .

## Some References

- C. Andrieu & M Vihola, “Convergence properties of pseudo-marginal MCMC”, *Ann. Applied Proba.*, 2015.
- 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.
- C. Sherlock, A. Thiery, G.O. Roberts & J.S. Rosenthal, “On the Efficiency of the RW Pseudo-Marginal MH”, *Ann. Stat.*, 2015.