

# Large Sample Asymptotics of the Pseudo-Marginal Method

Sebastian M Schmon, Arnaud Doucet, George Deligiannidis and Michael K Pitt



DEPARTMENT OF  
**STATISTICS**

## Introduction: Bayesian Inference and MCMC

- Observed data  $y_{1:T} := (y_1, \dots, y_T)$ ,  $T \geq 1$ , likelihood function  $p_\theta(y_{1:T})$  where  $\theta \in \Theta \subseteq \mathbb{R}^d$  and prior density  $p(\theta)$ .

- Bayesian inference relies on the posterior

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

- For complex models inference relies usually on Markov chain Monte Carlo techniques, i.e. one simulates an ergodic Markov chain  $\{\vartheta_i\}_{i \geq 1}$  with limiting distribution  $\pi_T(\theta)$ .

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

## Pseudo-Marginal Approach (I)

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

At iteration  $i$

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

- Compute an estimate  $\hat{p}_\theta(y)$  of  $p_\theta(y)$ .

- With probability

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

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

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

## Pseudo-Marginal Approach (II)

**Proposition:** If  $\hat{p}_\theta(y_{1:T})$  is a non-negative unbiased estimator of  $p_\theta(y_{1:T})$  then the pseudo-marginal MH admits  $\pi_T(\theta)$  as invariant density.

Let  $Z_T(\theta) = \log \{ \hat{p}_\theta(y_{1:T}) / p_\theta(y_{1:T}) \}$  be the error in log-likelihood estimator and introduce an auxiliary target density on  $\Theta \times \mathbb{R}$

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

where  $Z_T(\theta) \sim g_T(\cdot | \theta)$ ; e.g. importance sampling or particle filter.

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

$$\frac{\bar{\pi}_T(\vartheta, w)}{\bar{\pi}_T(\theta, z)} \frac{q(\theta | \vartheta) g_\theta(z)}{q(\vartheta | \theta) g_\theta(w)} = \frac{\hat{p}_\vartheta(y_{1:T}) p(\vartheta)}{\hat{p}_\theta(y_{1:T}) p(\theta)} \frac{q(\theta | \vartheta)}{q(\vartheta | \theta)}.$$

## Assumptions

**Assumption 1:** The posterior  $\{\pi_T(d\theta); T \geq 1\}$  concentrates in a "Bernstein-von Mises" sense:

$$\int \left| \pi_T(\theta) - \varphi(\theta; \hat{\theta}_T, \Sigma/T) \right| d\theta \xrightarrow{\mathbb{P}^Y} 0, \quad \hat{\theta}_T \xrightarrow{\mathbb{P}^Y} \bar{\theta},$$

with covariance matrix  $\Sigma$  and estimators  $\hat{\theta}_T$ .

**Assumption 2:** The proposal distributions  $\{q_T(\theta, d\theta'); T \geq 1\}$  are scaled with the data

$$\theta' = \theta + \frac{\xi}{\sqrt{T}}, \quad \xi \sim \nu(\cdot)$$

where  $\nu$  is a continuous probability density on  $\mathbb{R}^d$  with  $\mathbb{E}_\nu[\|\xi\|] < \infty$ .

**Assumption 3:** There exists an  $\varepsilon$ -ball  $B(\bar{\theta})$  around  $\bar{\theta}$  such that the distributions of the error in the log-likelihood estimator  $\{g_T(dz | \theta); T \geq 1\}$  satisfy

$$\sup_{\theta \in B(\bar{\theta})} d_{\text{BL}}(g_T(\cdot | \theta), \varphi(\cdot; -\sigma^2(\theta)/2, \sigma^2(\theta))) \xrightarrow{\mathbb{P}^Y} 0,$$

where  $d_{\text{BL}}$  denotes the bounded Lipschitz metric,  $\sigma: \Theta \rightarrow [0, \infty)$  is continuous at  $\bar{\theta}$  with  $0 < \sigma(\bar{\theta}) < \infty$ . An analogous result holds for  $\bar{g}_T(z | \theta) = \exp(z) g_T(z | \theta)$ , the distribution of this error at equilibrium, that is

$$\sup_{\theta \in B(\bar{\theta})} d_{\text{BL}}(\bar{g}_T(\cdot | \theta), \varphi(\cdot; \sigma^2(\theta)/2, \sigma^2(\theta))) \xrightarrow{\mathbb{P}^Y} 0.$$

## References and related work

Deligiannidis, G.; Doucet, A. & Pitt, M. K. *The Correlated Pseudo-Marginal Method*, Journal of the Royal Statistical Society Series B, to appear

Doucet, A.; Pitt, M. K.; Deligiannidis, G. & Kohn, R. *Efficient implementation of Markov chain Monte Carlo when using an unbiased likelihood estimator* Biometrika, 2015, 102 (2), 295-313

Pitt, M. K., dos Santos Silva, R., Giordani, P., & Kohn, R. *On some properties of Markov chain Monte Carlo simulation methods based on the particle filter* Journal of Econometrics, 2012, 171(2), 134-151.

Sherlock, C.; Thiery, A. H.; Roberts, G. O.; Rosenthal, J. S. & others *On the efficiency of pseudo-marginal random walk Metropolis algorithms* The Annals of Statistics, 2015, 43, 238-275

## Parameter Rescaling

Consider the pseudo-marginal chain  $\{(\vartheta_{T,k}, Z_{T,k}); k \geq 0\}$  started at  $(\vartheta_{T,0}, Z_{T,0}) \sim \bar{\pi}_T$ . Update  $(\vartheta_{T,k}, Z_{T,k}) \sim P_T(\vartheta_{T,k-1}, Z_{T,k-1}; \cdot)$  for  $k \geq 1$ . Let  $\chi_T = \{(\vartheta_{T,k}, Z_{T,k}); k \geq 0\}$  where  $\tilde{\vartheta}_{T,k} := \sqrt{T}(\vartheta_{T,k} - \theta_T)$  is the Markov chain arising from rescaling the parameter component of the pseudo-marginal chain. The transition kernel is

$$\tilde{P}_T(\tilde{\theta}, z; d\tilde{\theta}', dz') = \tilde{q}_T(\tilde{\theta}, d\tilde{\theta}') \tilde{g}_T(dz' | \tilde{\theta}') \tilde{\alpha}_T(\tilde{\theta}, z; \tilde{\theta}', z') + \tilde{\rho}_T(\tilde{\theta}, z) \delta_{(\tilde{\theta}, z)}(d\tilde{\theta}', dz'),$$

where

$$\tilde{\alpha}_T(\tilde{\theta}, z; \tilde{\theta}', z') = \min \left\{ 1, \frac{\tilde{\pi}_T(d\tilde{\theta}')}{\tilde{\pi}_T(d\tilde{\theta})} \frac{\tilde{q}_T(\tilde{\theta}', d\tilde{\theta})}{\tilde{q}_T(\tilde{\theta}, d\tilde{\theta}')} \exp(z' - z) \right\},$$

$\tilde{\rho}_T(\theta, z)$  is the corresponding rejection probability,  $\tilde{\pi}_T(\tilde{\theta}) := \pi_T(\hat{\theta}_T + \tilde{\theta}/\sqrt{T})/\sqrt{T}$ ,  $\tilde{q}_T(\tilde{\theta}, \tilde{\theta}') := q_T(\hat{\theta}_T + \tilde{\theta}/\sqrt{T}, \hat{\theta}_T + \tilde{\theta}'/\sqrt{T})/\sqrt{T}$  and  $\tilde{g}_T(z | \tilde{\theta}) := g_T(z | \hat{\theta}_T + \tilde{\theta}/\sqrt{T})$ .

## Theorem (Weak Convergence):

Under Assumptions 1, 2 and 3, the sequence of stationary Markov chains  $(\chi_T; T \geq 1)$  converges weakly in  $\mathbb{P}^Y$ -probability as  $T \rightarrow \infty$  to the law of a stationary Markov chain of initial distribution

$$\tilde{\pi}(d\tilde{\theta}, dz) := \varphi(d\tilde{\theta}; 0, \Sigma) \varphi(dz; \sigma^2/2, \sigma^2)$$

and transition kernel

$$\tilde{P}(\tilde{\theta}, z; d\tilde{\theta}', dz') = \tilde{q}(\tilde{\theta}, d\tilde{\theta}') \varphi(dz'; -\sigma^2/2, \sigma^2) \tilde{\alpha}(\tilde{\theta}, z; \tilde{\theta}', z') + \tilde{\rho}(\tilde{\theta}, z) \delta_{(\tilde{\theta}, z)}(d\tilde{\theta}', dz')$$

where  $\sigma := \sigma(\bar{\theta})$ ,

$$\tilde{\alpha}(\tilde{\theta}, z; \tilde{\theta}', z') = \min \left\{ 1, \frac{\varphi(\tilde{\theta}'; 0, \Sigma)}{\varphi(\tilde{\theta}; 0, \Sigma)} \frac{\tilde{q}(\tilde{\theta}', \tilde{\theta})}{\tilde{q}(\tilde{\theta}, \tilde{\theta}')} \exp(z' - z) \right\},$$

and  $\tilde{\rho}(\theta, z)$  is the corresponding rejection probability.

## Simulation Study: Tuning the Pseudo-Marginal algorithm

We optimize the performance of the limiting pseudo-marginal chain identified above as a proxy for the optimization of the original pseudo-marginal chain. We consider a Gaussian random walk proposal parameterized by  $\ell$

$$q(\theta, \theta') = \varphi\left(\theta'; \theta, \frac{\ell^2}{d} I_d\right).$$

Denoting  $\tau$  the integrated autocorrelation time, we minimize the computing time

$$\text{ct}(h, \tilde{P}_{\ell, \sigma}) = \frac{\tau(h, \tilde{P}_{\ell, \sigma})}{\sigma^2}$$

over a grid  $(\ell, \sigma) \in \{1.8, 1.9, \dots, 2.7\} \times \{1, 1.1, \dots, 2\}$ . We restrict attention here to the case where  $h(\theta, z) = \theta_1$ , the first component of  $\theta$ .

Dimension $d$	$\hat{\ell}_{\text{opt}}$	$\hat{\sigma}_{\text{opt}}$	$\text{ct}(\hat{\sigma}_{\text{opt}}, \hat{\ell}_{\text{opt}})$	$p_{\text{jump}}(\hat{\sigma}_{\text{opt}}, \hat{\ell}_{\text{opt}})$
$d = 1$	2.05 (0.25)	1.16 (0.07)	8.47	25.73%
$d = 2$	1.97 (0.14)	1.21 (0.06)	12.71	22.92%
$d = 3$	2.11 (0.07)	1.24 (0.05)	16.79	19.97%
$d = 5$	2.17 (0.12)	1.30 (0.05)	23.18	17.35%
$d = 10$	2.20 (0.08)	1.44 (0.05)	37.93	14.27%
$d = 15$	2.33 (0.08)	1.50 (0.00)	53.43	12.07%
$d = 20$	2.34 (0.10)	1.54 (0.05)	65.62	11.44%
$d = 30$	2.36 (0.11)	1.61 (0.03)	90.46	10.41%
$d = 50$	2.41 (0.10)	1.74 (0.05)	136.38	8.66%

**Table 1:** Optimal values for scaling  $\ell$  and noise  $\sigma$  and associated value of computing time and average acceptance probability.

## Real Data: Indonesian Preschool Children

We now consider a Bayesian logistic mixed effects model applied to a real data set with linear predictor

$$\eta_{t,j} = x_{t,j}^T \beta + U_t, \quad U_t \sim \mathcal{N}(0, \tau),$$

where  $U_t$  denotes the random intercept for children  $t = 1, \dots, T$  and  $\beta$  the regression parameters. The observations are assumed conditionally independent given the random effects and the likelihood of the population parameters is

$$p(y_{1:T} | \beta, \tau) = \prod_{t=1}^T \int \prod_{j=1}^J \frac{\exp(y_{t,j} \eta_{t,j})}{1 + \exp(\eta_{t,j})} \varphi(du_t, 0, \tau).$$

