

# Particle Metropolis–Hastings using gradient and Hessian information

Johan Dahlin · Fredrik Lindsten · Thomas B. Schön

Received: 15 March 2014 / Accepted: 28 August 2014 / Published online: 20 September 2014  
© Springer Science+Business Media New York 2014

**Abstract** Particle Metropolis–Hastings (PMH) allows for Bayesian parameter inference in nonlinear state space models by combining Markov chain Monte Carlo (MCMC) and particle filtering. The latter is used to estimate the intractable likelihood. In its original formulation, PMH makes use of a marginal MCMC proposal for the parameters, typically a Gaussian random walk. However, this can lead to a poor exploration of the parameter space and an inefficient use of the generated particles. We propose a number of alternative versions of PMH that incorporate gradient and Hessian information about the posterior into the proposal. This information is more or less obtained as a byproduct of the likelihood estimation. Indeed, we show how to estimate the required information using a fixed-lag particle smoother, with a computational cost growing linearly in the number of particles. We conclude that the proposed methods can: (i) decrease the length of the burn-in phase, (ii) increase the mixing of the Markov chain at the stationary phase, and (iii) make the proposal distribution scale invariant which simplifies tuning.

**Keywords** Sequential Monte Carlo · Particle Markov chain Monte Carlo · Manifold MALA · Fixed-lag particle smoothing · Parameter inference

J. Dahlin (✉)  
Department of Electrical Engineering, Linköping University,  
Linköping, Sweden  
e-mail: johan.dahlin@liu.se

F. Lindsten  
Department of Engineering, University of Cambridge,  
Cambridge, UK  
e-mail: fredrik.lindsten@eng.cam.ac.uk

T. B. Schön  
Department of Information Technology, Uppsala University,  
Uppsala, Sweden  
e-mail: thomas.schon@it.uu.se

## 1 Introduction

We are interested in Bayesian parameter inference in nonlinear state space models (SSM) of the form

$$x_t|x_{t-1} \sim f_\theta(x_t|x_{t-1}), \quad y_t|x_t \sim g_\theta(y_t|x_t), \quad (1)$$

where the latent states and the measurements are denoted by  $\mathbf{x} = x_{0:T} \triangleq \{x_t\}_{t=0}^T$  and  $\mathbf{y} = y_{1:T} \triangleq \{y_t\}_{t=1}^T$ , respectively. Here,  $f_\theta(\cdot)$  and  $g_\theta(\cdot)$  denote the transition and observation kernels, respectively, parametrised by the unknown static parameter vector  $\theta \in \Theta \subset \mathbb{R}^d$ . The initial state is distributed according to some distribution  $\mu(x_0)$  which, for notational simplicity, is assumed to be independent of  $\theta$ .

The aim of Bayesian parameter inference (in SSMs) is to compute the *parameter posterior distribution*

$$p(\theta|\mathbf{y}) = \frac{p_\theta(\mathbf{y})p(\theta)}{p(\mathbf{y})}, \quad (2)$$

where  $p(\theta)$  denotes the prior of  $\theta$  and  $p_\theta(\mathbf{y})$  denotes the likelihood, which for an SSM can be expressed as

$$p_\theta(\mathbf{y}) = p_\theta(y_1) \prod_{t=2}^T p_\theta(y_t|y_{1:t-1}). \quad (3)$$

The one-step ahead predictor  $p_\theta(y_t|y_{1:t-1})$ , and thus also the likelihood function, is in general not analytically tractable. However, unbiased estimators of the likelihood can be constructed using sequential Monte Carlo (SMC) (Doucet and Johansen 2011; Del Moral 2004) and these can be used as *plug-in estimators*. This is especially useful in the Metropolis–Hastings (MH) algorithm that can be used for estimating the parameter posterior in (2).

This combination of MH and SMC is known as the particle Metropolis–Hastings (PMH) algorithm (Andrieu et al. 2010).

The MH acceptance probability depends on the intractable likelihood, which in PMH is estimated using SMC (see Sect. 2). Despite the apparent approximation, this results in an algorithm that targets the correct posterior distribution (Andrieu et al. 2010). The original PMH algorithm makes use of a marginal proposal for  $\theta$ , i.e. only the current parameter is used when proposing a new parameter. The theoretical properties of the marginal PMH algorithm have been analysed in Andrieu and Vihola (2011), Pitt et al. (2012), Doucet et al. (2012) and it has been applied for a number of interesting applications in, e.g., economics, social network analysis and ecology (Flury and Shephard 2011; Everitt 2012; Golightly and Wilkinson 2011).

In this paper, we show that information such as the gradient and the Hessian about the posterior can be included in the construction of the PMH proposal. This idea is first suggested by Doucet et al. (2011) in the discussions following Girolami and Calderhead (2011). In two previous proceedings, we have applied and extended this idea with gradient information (Dahlin et al. 2013) and also using Hessian information (Dahlin et al. 2014). The present article builds upon and extends this preliminary work. A PMH method using gradient information similar to Dahlin et al. (2013) has recently been proposed by Nemeth and Fearnhead (2014).

In the context of MH sampling, it has been recognised that the gradient and Hessian can be used to construct efficient proposal distributions. In the Metropolis adjusted Langevin algorithm (MALA) (Roberts and Stramer 2003), a drift term is added to the proposal in the direction of the gradient, which intuitively guides the Markov chain to regions of high posterior probability. In the manifold MALA (mMALA) (Girolami and Calderhead 2011), the Hessian (or some other appropriate metric tensor) is also included to scale the proposal to take the curvature of the log-posterior into account. Drawing parallels with the optimisation literature, mMALA shares some properties with Newton-type optimisation algorithms (where MALA is more similar to a steepest ascent method). In particular, scaling the proposal with the Hessian can considerably simplify the tedious tuning of the method since it removes the need for running costly pilot runs, which are commonly used to tune the covariance matrices of the random walk MH and the MALA.

In our problem, i.e. for inference in a nonlinear SSM (1), the gradient and Hessian cannot be computed analytically. However, in analogue with the intractable likelihood, these quantities can be estimated using SMC algorithms, see e.g. Poyiadjis et al. (2011); Doucet et al. (2013). This provides us with the tools necessary to construct PMH algorithms in the flavour of the MALA and the mMALA, resulting in the two methods proposed in this paper, PMH1 and PMH2, respectively. In particular, we make use of a fixed-lag (FL) particle smoother (Kitagawa and Sato 2001) to estimate the gradient and Hessian. The motivation for this is that this smoother

only makes use of the weighted particles computed by the particle filter. Consequently, we obtain this information as a *byproduct* of the likelihood computation in the PMH algorithm. This results in only a small computational overhead for the proposed methods when compared to the marginal method.

Finally, we provide numerical experiments to illustrate the benefits of using the gradient and Hessian and the accuracy of the FL smoother. We demonstrate some interesting properties of the proposed algorithms, in particular that they enjoy (i) a shorter burn-in compared with the marginal algorithm, (ii) a better mixing of the Markov chain in the stationary phase, and (iii) a simplified tuning of the step length(s), especially when the target distribution is non-isotropic.

## 2 Particle Metropolis–Hastings

In this section, we review the PMH algorithm and show how the random variables used to compute the likelihood estimator can be incorporated in the proposal construction. We also outline the idea of how this can be used to construct the proposed PMH1 and PMH2 algorithms.

### 2.1 MH sampling with unbiased likelihoods

The MH algorithm (see, e.g., Robert and Casella (2004)) is a member of the MCMC family for sampling from a target distribution  $\pi(\theta)$  by simulating a carefully constructed Markov chain on  $\Theta$ . The chain is constructed in such a way that it admits the target as its unique stationary distribution.

The algorithm consists of two steps: (i) a new parameter  $\theta''$  is sampled from a proposal distribution  $q(\theta''|\theta')$  given the current state  $\theta'$  and (ii) the current parameter is changed to  $\theta''$  with probability  $\alpha(\theta', \theta'')$ , otherwise the chain remains at the current state. The acceptance probability is given by

$$\alpha(\theta', \theta'') = 1 \wedge \frac{\pi(\theta'') q(\theta'|\theta'')}{\pi(\theta') q(\theta''|\theta')}, \quad (4)$$

where we use the notation  $a \wedge b \triangleq \min\{a, b\}$ .

In this paper, we have the parameter posterior distribution (2) as the target distribution, i.e.  $\pi(\theta) = p(\theta|\mathbf{y})$ . This implies that the acceptance probability (4) will depend explicitly on the intractable likelihood  $p_\theta(\mathbf{y})$ , preventing direct application of the MH algorithm to this problem. However, this difficulty can be circumvented by using a *pseudo-marginal* approach (Beaumont 2003; Andrieu and Roberts 2009).

Assume that there exists an unbiased, non-negative estimator of the likelihood  $\hat{p}_\theta(\mathbf{y}|u)$ . We introduce explicitly the random variable  $u \in \mathbf{U}$  used to construct this estimator, and we let  $m_\theta(u)$  denote the probability density of  $u$  on  $\mathbf{U}$ .

The pseudo-marginal method is then a standard MH algorithm operating in a non-standard extended space  $\Theta \times \mathbf{U}$ , with the *extended target*

$$\pi(\theta, u|\mathbf{y}) = \frac{\widehat{p}_\theta(\mathbf{y}|u)m_\theta(u)p(\theta)}{p(\mathbf{y})} = \frac{\widehat{p}_\theta(\mathbf{y}|u)m_\theta(u)p(\theta|\mathbf{y})}{p_\theta(\mathbf{y})},$$

and proposal distribution  $m_{\theta'}(u'')q(\theta''|\theta')$ .

Since the likelihood estimator is unbiased,  $\mathbb{E}_{u|\theta}[\widehat{p}_\theta(\mathbf{y}|u)] = p_\theta(\mathbf{y})$ , it follows that the extended target admits  $p(\theta|\mathbf{y})$  as a marginal. Hence, by simulating from the extended target  $\pi(\theta, u|\mathbf{y})$  we obtain samples from the original target distribution  $p(\theta|\mathbf{y})$  as a byproduct.

If the likelihood is estimated by using SMC (see Sect. 3) we obtain the PMH algorithm. The random variable  $u$  then corresponds to all the weighted particles generated by the SMC algorithm. However, these random variables carry useful information, not only about the likelihood, but also about the geometry of the posterior distribution. We suggest to incorporate this information into the proposal construction. With  $(\theta', u')$  being the current state of the Markov chain we simulate  $\theta'' \sim q(\cdot|\theta', u')$  and  $u'' \sim m_{\theta'}(\cdot)$ , using some proposal  $q$  (see Sect. 2.2).

It follows that the (standard) MH acceptance probability for the extended target is given by

$$\begin{aligned} \alpha(\theta'', u'', \theta', u') &= 1 \wedge \frac{\widehat{p}_{\theta''}(\mathbf{y}|u'')m_{\theta''}(u'')p(\theta'')}{\widehat{p}_{\theta'}(\mathbf{y}|u')m_{\theta'}(u')p(\theta')} \frac{m_{\theta'}(u')q(\theta'|\theta'', u'')}{m_{\theta''}(u'')q(\theta''|\theta', u')} \\ &= 1 \wedge \frac{\widehat{p}_{\theta''}(\mathbf{y}|u'')p(\theta'')}{\widehat{p}_{\theta'}(\mathbf{y}|u')p(\theta')} \frac{q(\theta'|\theta'', u'')}{q(\theta''|\theta', u')}. \end{aligned} \tag{5}$$

Note that  $q(\theta''|\theta', u')$  may depend on the auxiliary variable  $u'$  in a (formally) arbitrary way. In particular, in Sect. 3 we propose a construction making use of *biased* estimates of the gradient and Hessian of the log-posterior. Nevertheless, expression (5) still defines a correct MH acceptance probability for the extended target, ensuring the validity of our approach. Note also that the aforementioned proposal construction opens up for a wide range of adapted proposals, possibly different from the ones considered in this work.

### 2.2 Constructing PMH1 and PMH2

We now turn to the construction of a proposal that makes use of the gradient and Hessian of the log-posterior. Following Robert and Casella (2004), we do this by a Laplace approximation of the parameter posterior around the current state  $\theta'$ . Hence, consider a second order Taylor expansion of  $\log p(\theta''|\mathbf{y})$  at  $\theta'$ :

$$\begin{aligned} \log p(\theta''|\mathbf{y}) &\approx \log p(\theta'|\mathbf{y}) + (\theta'' - \theta')^\top \left[ \nabla \log p(\theta|\mathbf{y}) \right]_{\theta=\theta'} \\ &\quad + \frac{1}{2} (\theta'' - \theta')^\top \left[ \nabla^2 \log p(\theta|\mathbf{y}) \right]_{\theta=\theta'} (\theta'' - \theta'). \end{aligned}$$

Taking the exponential of both sides and completing the square, we obtain

$$p(\theta''|\mathbf{y}) \approx \mathbf{N}\left(\theta''; \theta' + \mathbf{l}_T^{-1}(\theta')\mathbf{S}_T(\theta'), \mathbf{l}_T^{-1}(\theta')\right),$$

where we have introduced  $\mathbf{S}_T(\theta') = \nabla \log p(\theta|\mathbf{y})|_{\theta=\theta'}$  and  $\mathbf{l}_T(\theta') = -\nabla^2 \log p(\theta|\mathbf{y})|_{\theta=\theta'}$ , for the gradient and the negative Hessian of the log-posterior, respectively. Here, we assume for now that the negative Hessian is positive definite; see Sect. 3.5 for further discussion on this matter.

As pointed out above, these quantities cannot be computed in closed form, but they can be estimated from the random variable  $u'$  (see Sect. 3). This suggests three different versions of the PMH algorithm, each resulting from a specific choice of the proposal:

$$q(\theta''|\theta', u') = \begin{cases} \mathbf{N}(\theta', \Gamma), & \text{[PMH0]} \\ \mathbf{N}\left(\theta' + \frac{1}{2}\Gamma\widehat{\mathbf{S}}_T(\theta'|u'), \Gamma\right), & \text{[PMH1]} \\ \mathbf{N}\left(\theta' + \widehat{\mathbf{G}}(\theta'|u'), \widehat{\mathbf{H}}(\theta'|u')\right). & \text{[PMH2]} \end{cases} \tag{6}$$

Here, we use the notation  $\widehat{\mathbf{G}}(\theta|u) = \frac{1}{2}\Gamma\widehat{\mathbf{l}}_T^{-1}(\theta|u)\widehat{\mathbf{S}}_T(\theta|u)$  and  $\widehat{\mathbf{H}}(\theta|u) = \Gamma\widehat{\mathbf{l}}_T^{-1}(\theta|u)$  for the natural gradient and scaled inverse Hessian, respectively. Furthermore,  $\Gamma$  denotes a scaling matrix that controls the step lengths of the proposal. For PMH0 and PMH1,  $\Gamma$  can be chosen as the inverse of an estimate of the posterior covariance matrix. However, computing this estimate typically requires costly and tedious trial runs. For PMH2, the curvature of the problem is captured by the Hessian matrix, i.e. a single step length can be used which can significantly simplify the tuning. It is also possible to choose different step lengths for the drift term and for the covariance matrix of the proposal.

The final PMH2 algorithm is presented in Algorithm 1. It makes use of Algorithm 2, described in Sect. 3, to estimate the quantities needed for computing the proposal and the acceptance probability. Clearly, PMH0 and PMH1 are special cases obtained by using the corresponding proposal from (6) in the algorithm. Note that, while the algorithm make explicit reference to the auxiliary variable  $u$ , it only depends on this variable through the estimates  $\widehat{p}_{\theta'}(\mathbf{y})$ ,  $\widehat{\mathbf{S}}_T(\theta')$  and  $\widehat{\mathbf{l}}_T(\theta')$ .

### 2.3 Properties of the PMH1 and PMH2 proposals

In the sequel, we use a single step size  $\Gamma = \gamma^2 I_d$  for all the parameters in the (standard) proposal. This is done to illustrate the advantage of adding the Hessian information, which rescales the step lengths according to the local curvature. Hence, it allows for taking larger steps when the curvature is small and vice versa.

This property of PMH2 makes the algorithm scale-free in the same manner as a Newton algorithm in optimisation

**Algorithm 1** Second order particle Metropolis–Hastings

INPUTS: Algorithm 2.  $M > 0$  (no. MCMC steps),  $\theta_0$  (initial parameters),  $\gamma$  (step length).  
 OUTPUT:  $\theta = \{\theta_1, \dots, \theta_M\}$  (samples from the posterior).

- 1: Run Algorithm 2 to obtain  $\widehat{p}_{\theta_0}(\mathbf{y})$ ,  $\widehat{\mathbf{S}}_T(\theta_0)$  and  $\widehat{\Gamma}_T(\theta_0)$ .
- 2: **for**  $k = 1$  to  $M$  **do**
- 3:   Sample  $\theta' \sim q(\theta' | \theta_{k-1}, u_{k-1})$  by (6),  $\widehat{\mathbf{S}}_T(\theta_{k-1})$  and  $\widehat{\Gamma}_T(\theta_{k-1})$ .
- 4:   Run Algorithm 2 to obtain  $\widehat{p}_{\theta'}(\mathbf{y})$ ,  $\widehat{\mathbf{S}}_T(\theta')$  and  $\widehat{\Gamma}_T(\theta')$ .
- 5:   Sample  $\omega_k$  uniformly over  $[0, 1]$ .
- 6:   **if**  $\omega_k < \alpha(\theta', u', \theta_{k-1}, u_{k-1})$  given by (5) **then**
- 7:      $\theta_k \leftarrow \theta'$ . {Accept the parameter}
- 8:      $\{\widehat{p}_{\theta_k}(\mathbf{y}), \widehat{\mathbf{S}}_T(\theta_k), \widehat{\Gamma}_T(\theta_k)\} \leftarrow \{\widehat{p}_{\theta'}(\mathbf{y}), \widehat{\mathbf{S}}_T(\theta'), \widehat{\Gamma}_T(\theta')\}$ .
- 9:   **else**
- 10:      $\theta_k \leftarrow \theta_{k-1}$ . {Reject the parameter}
- 11:      $\{\widehat{p}_{\theta_k}(\mathbf{y}), \widehat{\mathbf{S}}_T(\theta_k), \widehat{\Gamma}_T(\theta_k)\} \leftarrow \{\widehat{p}_{\theta_{k-1}}(\mathbf{y}), \widehat{\mathbf{S}}_T(\theta_{k-1}), \widehat{\Gamma}_T(\theta_{k-1})\}$ .
- 12:   **end if**
- 13: **end for**

(Nocedal and Wright 2006, Chapter 3). That is, the proposal is invariant to affine transformations of the parameters. Note that, since the local information is used, this is different from scaling the proposal in PMH0 with the posterior covariance matrix estimated from a pilot run, as this only takes the geometry at the mode of the posterior into account.

Some analyses of the statistical properties are available for PMH0 (Sherlock et al. 2013), MH using a random walk (Roberts et al. 1997) and MALA (Roberts and Rosenthal 1998). It is known from these analyses that adding the gradient into the proposal can increase the mixing of the Markov chain. Note that these results are obtained under somewhat strict assumptions. Also, we know from numerical experiments (Girolami and Calderhead 2011) that there are further benefits of also taking the local curvature into account.

**3 Estimation of the likelihood, gradient, and Hessian**

In this section, we show how to estimate the likelihood together with the gradient and Hessian using SMC methods.

**3.1 Auxiliary particle filter**

An auxiliary particle filter (APF) (Pitt and Shephard 1999) can be used to approximate the sequence of joint smoothing distributions (JSDs)  $p_{\theta}(x_{1:t} | y_{1:t})$  for  $t = 1$  to  $T$ . The APF makes use of a particle system consisting of  $N$  weighted particles  $\{x_{1:t}^{(i)}, w_t^{(i)}\}_{i=1}^N$  to approximate the JSD at time  $t$  by

$$\widehat{p}_{\theta}(dx_{1:t} | y_{1:t}) \triangleq \sum_{i=1}^N \frac{w_t^{(i)}}{\sum_{k=1}^N w_t^{(k)}} \delta_{x_{1:t}^{(i)}}(dx_{1:t}). \tag{7}$$

Here,  $\delta_z(dx_{1:t})$  denotes the Dirac measure placed at  $z$ . The particle system is propagated from  $t - 1$  to  $t$  by first sampling an ancestor index  $a_t^{(i)}$ , with

$$\mathbb{P}(a_t^{(i)} = j) = v_{t-1}^{(j)} \left[ \sum_{k=1}^N v_{t-1}^{(k)} \right]^{-1}, \quad i, j = 1, \dots, N, \tag{8}$$

where  $v_{t-1}^{(i)}$  denotes the resampling weights. Given the ancestor index, a new particle is sampled according to

$$x_t^{(i)} \sim R_{\theta}(x_t | x_{1:t-1}^{a_t^{(i)}}, y_t), \quad i = 1, \dots, N. \tag{9}$$

Finally, we append the obtained sample to the trajectory by  $x_{1:t}^{(i)} = \{x_{1:t-1}^{a_t^{(i)}}, x_t^{(i)}\}$  and compute a new importance weight by

$$w_t^{(i)} \triangleq \frac{w_{t-1}^{a_t^{(i)}} g_{\theta}(y_t | x_t^{(i)}) f_{\theta}(x_t^{(i)} | x_{t-1}^{a_t^{(i)}})}{v_{t-1}^{a_t^{(i)}} R_{\theta}(x_t^{(i)} | x_{1:t-1}^{a_t^{(i)}}, y_t)}, \quad i = 1, \dots, N. \tag{10}$$

Hence, the empirical approximations of the smoothing distributions (7) can be computed sequentially for  $t = 1$  to  $T$  by repeating (8)–(10).

Note that the random variables  $u$  appearing in the extended target of the PMH algorithm correspond to all the random variables generated by the APF, i.e. all the particles and ancestor indices,

$$u = \left( \left\{ x_t^{(i)}, a_t^{(i)} \right\}_{i=1}^N, t = 1, \dots, T \right).$$

In this article, we make use of two important special cases of the APF: the bootstrap particle filter (bPF) (Gordon et al. 1993) and the fully adapted particle filter (faPF) (Pitt and Shephard 1999). For the bPF, we select the proposal kernel  $R_{\theta}(x_t | x_{1:t-1}, y_t) = f_{\theta}(x_t | x_{t-1})$  and the auxiliary weights  $v_t = w_t = g_{\theta}(y_t | x_t)$ . The faPF is obtained by  $R_{\theta}(x_t | x_{1:t-1}, y_t) = p_{\theta}(x_t | y_t, x_{t-1})$  and  $v_t = p_{\theta}(y_{t+1} | x_t)$ , resulting in the weights  $w_t \equiv 1$ . Note, that the faPF can only be used in models for which these quantities are available in closed-form.

**3.2 Estimation of the likelihood**

The likelihood for the SSM in (1) can be estimated using (3) by inserting estimated one-step predictors  $p_{\theta}(y_t | y_{1:t-1})$  obtained from the APF. The resulting likelihood estimator is given by

$$\widehat{p}_{\theta}(\mathbf{y} | u) = \frac{1}{N^T} \sum_{i=1}^N w_T^{(i)} \left\{ \prod_{t=1}^{T-1} \sum_{i=1}^N v_t^{(i)} \right\}. \tag{11}$$

It is known that this likelihood estimator is unbiased for any number of particles, see e.g. Pitt et al. (2012) and Proposition



7.4.1 in Del Moral (2004). As discussed in Sect. 2.1, this is exactly the property that is needed in order to obtain  $p(\theta|\mathbf{y})$  as the unique stationary distribution for the Markov chain generated by the PMH algorithm.

Consequently, PMH will target the correct distribution for any number of particles  $N \geq 1$ . However, the variance in the likelihood estimate is connected with the acceptance rate and the mixing of the Markov chain. Therefore it is important to determine the number of particles that balances a reasonable acceptance rate with a reasonable computational cost. This problem is studied for PMH0 in Pitt et al. (2012), Doucet et al. (2012).

### 3.3 Estimation of the gradient

As we shall see below, the gradient of the log-posterior can be estimated by solving a smoothing problem. The APF can be used directly to address this problem, since the particles  $\{x_{1:T}^{(i)}, w_T^{(i)}\}_{i=1}^N$  provide an approximation of the JSD at time  $T$  according to (7) (see also Poyiadjis et al. 2011). However, this method can give estimates with high variance due to *particle degeneracy*.

Instead, we make use of the FL smoother (Kitagawa and Sato 2001) which has the same linear computational cost, but smaller problems with *particle degeneracy* than the APF. Alternative algorithms for estimating this information are also available (Del Moral et al. 2010; Poyiadjis et al. 2011).

The gradient of the parameter log-posterior is given by

$$\mathbf{S}_T(\theta) = \nabla \log p(\theta) + \nabla \log p_\theta(\mathbf{y}), \tag{12}$$

where it is assumed that the gradient of the log-prior  $\nabla \log p(\theta)$  can be calculated explicitly. The gradient of the log-likelihood  $\nabla \log p_\theta(\mathbf{y})$  can, using *Fisher’s identity* (Cappé et al. 2005), be expressed as

$$\nabla \log p_\theta(\mathbf{y}) = \mathbb{E}_\theta \left[ \nabla \log p_\theta(\mathbf{x}, \mathbf{y}) \mid \mathbf{y} \right], \tag{13}$$

where for an SSM (1) we can write the gradient of the complete data log-likelihood as

$$\nabla \log p_\theta(\mathbf{x}, \mathbf{y}) = \sum_{t=1}^T \xi_\theta(x_t, x_{t-1}), \quad \text{where} \tag{14}$$

$$\xi_\theta(x_t, x_{t-1}) = \nabla \log f_\theta(x_t | x_{t-1}) + \nabla \log g_\theta(y_t | x_t).$$

Combining (14) with Fisher’s identity (13) yields

$$\nabla \log p_\theta(\mathbf{y}) = \sum_{t=1}^T \int \xi_\theta(x_t, x_{t-1}) p_\theta(x_{t-1:t} | \mathbf{y}) dx_{t-1:t},$$

which depends on the (intractable) two-step smoothing distribution  $p_\theta(x_{t-1:t} | \mathbf{y})$ . To approximate this quantity we use the FL smoother which relies on the assumption that there

is a decaying influence of future observations  $y_{t+\Delta:T}$  on the state  $x_t$ . This means that

$$p_\theta(x_{t-1:t} | \mathbf{y}) \approx p_\theta(x_{t-1:t} | y_{1:\kappa_t}),$$

holds for some large enough  $\kappa_t = \min\{t + \Delta, T\}$ . Here,  $\Delta$  denotes a pre-determined lag decided by the user, which depends on the forgetting properties of the model. By marginalisation of the empirical smoothing distribution  $\widehat{p}_\theta(x_{1:\kappa_t} | y_{1:\kappa_t})$  over  $x_{1:t-2}$  and  $x_{t+1:\kappa_t}$ , we obtain the approximation

$$\widehat{p}_\theta^\Delta(dx_{t-1:t} | \mathbf{y}) \triangleq \sum_{i=1}^N w_{\kappa_t}^{(i)} \delta_{\tilde{x}_{\kappa_t, t-1:t}^{(i)}}(dx_{t-1:t}). \tag{15}$$

Here, we use the notation  $\tilde{x}_{\kappa_t, t}^{(i)}$  to denote the ancestor at time  $t$  of particle  $x_{\kappa_t}^{(i)}$  and  $\tilde{x}_{\kappa_t, t-1:t}^{(i)} = \{\tilde{x}_{\kappa_t, t-1}^{(i)}, \tilde{x}_{\kappa_t, t}^{(i)}\}$ . Inserting (14)–(15) into (13) provides an estimator of (12),

$$\widehat{\mathbf{S}}_T(\theta | u) = \nabla \log p(\theta) + \sum_{t=1}^T \sum_{i=1}^N w_{\kappa_t}^{(i)} \xi_\theta(\tilde{x}_{\kappa_t, t}^{(i)}, \tilde{x}_{\kappa_t, t-1}^{(i)}), \tag{16}$$

which is used in the proposal distributions in (6).

### 3.4 Estimation of the Hessian

The negative Hessian of the parameter log-posterior can be written as

$$l_T(\theta) = -\nabla^2 \log p(\theta) - \nabla^2 \log p_\theta(\mathbf{y}), \tag{17}$$

where it is assumed that the Hessian of the log-prior  $\nabla^2 \log p(\theta)$  can be calculated analytically. The negative Hessian of the log-likelihood, also known as the *observed information matrix*, can using *Louis’ identity* (Cappé et al. 2005) be expressed as

$$-\nabla^2 \log p_\theta(\mathbf{y}) = \nabla \log p_\theta(\mathbf{y})^2 - \mathbb{E}_\theta \left[ \nabla^2 \log p_\theta(\mathbf{x}, \mathbf{y}) \mid \mathbf{y} \right] - \mathbb{E}_\theta \left[ \nabla \log p_\theta(\mathbf{x}, \mathbf{y})^2 \mid \mathbf{y} \right]. \tag{18}$$

Here, we have introduced the notation  $v^2 = vv^\top$  for a vector  $v$ . From this, we can construct an estimator of (17) using the estimator of the gradient in (16), of the form

$$\widehat{l}_T(\theta | u) = -\nabla^2 \log p(\theta) + \widehat{\mathbf{S}}_T(\theta | u)^2 - \widehat{l}_T^{(1)}(\theta | u) - \widehat{l}_T^{(2)}(\theta | u), \tag{19}$$

where we introduce  $l_T^{(1)}(\theta) = \mathbb{E}_\theta \left[ \nabla^2 \log p_\theta(\mathbf{x}, \mathbf{y}) \mid \mathbf{y} \right]$  and  $l_T^{(2)}(\theta) = \mathbb{E}_\theta \left[ \nabla \log p_\theta(\mathbf{x}, \mathbf{y})^2 \mid \mathbf{y} \right]$ . We obtain the estimator of

the first term analogously to (16) as

$$\widehat{I}_T^{(1)}(\theta|u) = \sum_{t=1}^T \sum_{i=1}^N w_{\kappa_t}^{(i)} \zeta_\theta(\tilde{x}_{\kappa_t,t}^{(i)}, \tilde{x}_{\kappa_t,t-1}^{(i)}), \quad \text{where} \quad (20)$$

$$\zeta_\theta(x_t, x_{t-1}) = \nabla^2 \log f_\theta(x_t|x_{t-1}) + \nabla^2 \log g_\theta(y_t|x_t).$$

The estimator of the second term needs a bit more work and we start by rewriting the last term in (18) as

$$\begin{aligned} & \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}_\theta \left[ \xi_\theta(x_t, x_{t-1}) \xi_\theta(x_s, x_{s-1})^\top | \mathbf{y} \right] \\ &= \sum_{t=1}^T \left\{ \mathbb{E}_\theta \left[ \xi_\theta(x_t, x_{t-1})^2 | \mathbf{y} \right] \right. \\ & \quad \left. + \sum_{s=1}^{t-1} \mathbb{E}_\theta \left[ (\xi_\theta(x_t, x_{t-1}), \xi_\theta(x_s, x_{s-1}))^\dagger | \mathbf{y} \right] \right\}, \quad (21) \end{aligned}$$

where we have introduced the operator  $(a, b)^\dagger = ab^\top + ba^\top$  for brevity. Consider the last term appearing in this expression, we can rewrite it as

$$\begin{aligned} & \sum_{s=1}^{t-1} \mathbb{E}_\theta \left[ \xi_\theta(x_t, x_{t-1}) \xi_\theta(x_s, x_{s-1})^\top | \mathbf{y} \right] \\ &= \mathbb{E}_\theta \left[ \xi_\theta(x_t, x_{t-1}) \underbrace{\left\{ \sum_{s=1}^{t-1} \mathbb{E}_\theta \left[ \xi_\theta(x_s, x_{s-1}) | x_{t-1}, y_{1:t-1} \right] \right\}^\top}_{\triangleq \alpha_\theta(x_{t-1})^\top} | \mathbf{y} \right]. \end{aligned}$$

From this, we see that (21) can be written as an additive functional of the form

$$\sum_{t=1}^T \mathbb{E}_\theta \left[ (\xi_\theta(x_t, x_{t-1}))^2 + ((\xi_\theta(x_t, x_{t-1}), \alpha_\theta(x_{t-1}))^\dagger) | \mathbf{y} \right],$$

which can be estimated using the FL smoother as before. However, for this we need to compute the quantities  $\alpha_\theta(x_{t-1})$ . One option is to make use of a type of fixed-lag approximation for  $\alpha_\theta(x_{t-1})$ , by assuming that  $x_s$  and  $x_t$  are conditionally independent given  $y_{1:\kappa_t}$ , whenever  $|s - t| > \Delta$ . This approach has previously been used by Doucet et al. (2013). Alternatively, we can use a filter approximation according to

$$\widehat{\alpha}_\theta(x_t^{(i)}) = \widehat{\alpha}_\theta(x_{t-1}^{a_t(i)}) + \xi_\theta(x_t^{(i)}, x_{t-1}^{a_t(i)}), \quad (22)$$

for  $i = 1, \dots, N$ . Note that this approach suffers from the same particle degeneracy as the APF. However, this only affects a small number of terms and in our experience this approximation works sufficiently well to give estimates with

reasonable variance. The resulting estimator using (21) is

$$\widehat{I}_T^{(2)}(\theta|u) = \sum_{t=1}^T \sum_{i=1}^N w_{\kappa_t}^{(i)} \eta_\theta(\tilde{x}_{\kappa_t,t}^{(i)}, \tilde{x}_{\kappa_t,t-1}^{(i)}), \quad \text{where} \quad (23)$$

$$\eta_\theta(x_t, x_{t-1}) = \xi_\theta(x_t, x_{t-1})^2 + (\xi_\theta(x_t, x_{t-1}), \widehat{\alpha}_\theta(x_{t-1}))^\dagger.$$

Hence, the Hessian can be estimated using (19) by inserting the estimators from (20), (22) and (23).

### 3.5 Regularisation of the estimate of the Hessian

The PMH2 proposal (6) relies on the assumption that the observed information matrix is positive definite (PD). The estimator given in (19) does not always satisfy this, especially when the Markov chain is located far from the posterior mode. Typically, the amount of information is limited in such regions and this results in that the curvature is difficult to estimate. To cope with this issue, one alternative is to regularize the Hessian by adding a diagonal matrix to shift the eigenvalues to be positive. The diagonal matrix can e.g. be selected such that

$$\Delta \widehat{I}_T = \max \left\{ 0, -2\lambda_{\min}(\widehat{I}_T) \right\} I_d, \quad (24)$$

where  $\lambda_{\min}(\widehat{I}_T)$  denotes the smallest eigenvalue of  $\widehat{I}_T(\theta|u)$ . In this article, we make use of this method for handling non-PD estimates of the negative Hessian for the PMH2 algorithm. This heuristic is common for Newton-type optimisation algorithms (Nocedal and Wright 2006, Chapter 3.4).

Note, that there are other solutions available for ensuring positive definiteness that only shifts the negative eigenvalues, see (Nocedal and Wright 2006, Chapter 3). We emphasise that this type of regularization keeps the Markov chain invariant, i.e. still targets the correct posterior distribution (recall Sect. 2.1).

Another alternative is to replace the estimate of the negative Hessian with the inverse sample covariance matrix calculated using the trace of Markov chain when the estimate is not PD. This can be seen as a hybrid between the PMH2 algorithm and a *pre-conditioned* PMH1 algorithm. This resembles some other adaptive MH algorithms (Andrieu and Thoms 2008) in which the same procedure is used to adapt the covariance matrix of a random walk proposal. For this, we can make use of the last  $L$  iterations of the MH algorithm after that the Markov chain has reached stationarity. During the burn-in phase, non-PD estimates can be handled using a regularization approach or by rejecting the proposed parameter. In this article, we refer to this method for handling non-PD estimates of the negative Hessian as the *hybrid* PMH2 algorithm, where we use the latter alternative during the burn-in phase. Note that this pre-conditioning can also

**Algorithm 2** Estimation of the likelihood, the gradient and the Hessian of the log-posterior

INPUTS:  $\mathbf{y}$  (data),  $R(\cdot)$  (propagation kernel),  $\nu(\cdot)$  (weight function),  $N > 0$  (no. particles),  $0 < \Delta \leq T$  (lag).  
 OUTPUTS:  $\widehat{p}_\theta(\mathbf{y})$  (est. of the likelihood),  $\widehat{\mathbf{S}}_T(\theta)$  (est. of the gradient),  $\widehat{\mathbf{l}}_T(\theta)$  (est. of the negative Hessian).

- 1: Initialise each particle  $x_0^{(i)}$ .
- 2: **for**  $t = 1$  to  $T$  **do**
- 3:   Resample and propagate each particle using (9).
- 4:   Calculate the weights for each particle using (10).
- 5: **end for**
- 6: Compute  $\widehat{p}_\theta(\mathbf{y})$  by (11).
- 7: Compute  $\widehat{\mathbf{S}}_T(\theta)$  and  $\widehat{\mathbf{l}}_T(\theta)$  by (16), and (19).
- 8: **if**  $\widehat{\mathbf{l}}_T(\theta) \leq 0$  **then**
- 9:   [standard] Regularize  $\widehat{\mathbf{l}}_T(\theta)$  by adding  $\widehat{\Delta\mathbf{l}}_T$  computed by (24)
- 10:   [hybrid] Replace  $\widehat{\mathbf{l}}_T(\theta)$  by the inverse covariance matrix computed using the  $L$  final samples of the Markov chain during the burn-in.
- 11: **end if**

be applied to the PMH0 and PMH1 algorithm, we return to this in Sect. 4.4.

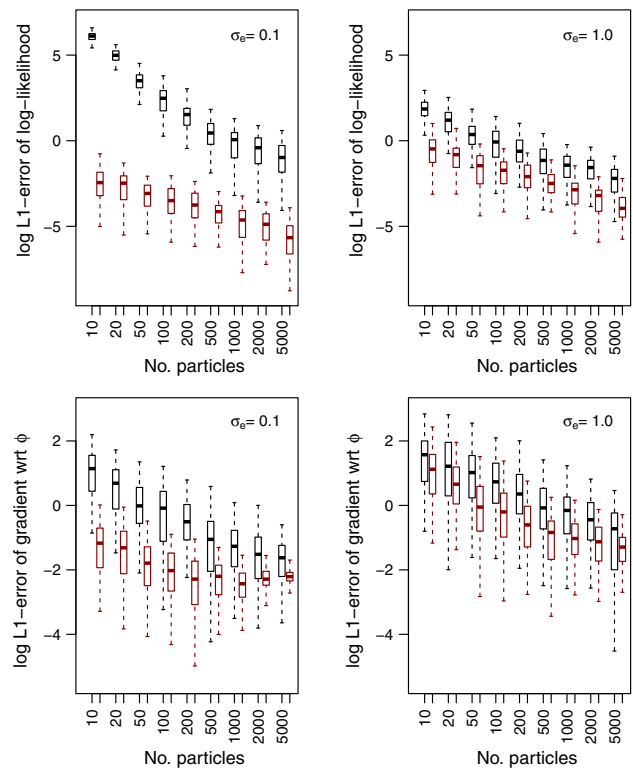
3.6 Resulting SMC algorithm

In Algorithm 2, we present the complete procedure that combines the APF with the FL smoother to compute the estimates needed for the PMH2 proposal (6). Note that the two different methods to handle non-PD estimates of the negative Hessian matrix results in the *standard* and *hybrid* PMH2 algorithm, respectively.

We end this section by briefly discussing the statistical properties of the estimates of the gradient and Hessian obtained from the FL smoother. From Olsson et al. (2008), we know that the FL smoother gives biased estimates of the gradient and Hessian for any number of particles. Remember that this does not effect the invariance of the Markov chain (recall Sect. 2.1). The main advantage of the FL smoother over the APF (which gives a consistent estimate) is that the former enjoys a smaller variance than the APF, i.e. we obtain a favourable bias-variance trade-off for a certain choice of lag  $\Delta$ . Note that a too small lag gives a large bias in the estimate and a too large lag gives a large variance in the estimate; we return to this choice in Sect. 4.

4 Numerical illustrations

In this section, we provide illustrations of the properties of the FL smoother and the different proposed algorithms. The source code in Python and the data used for some of the numerical illustrations are available for download at: <http://liu.johandahlin.com/>.



**Fig. 1** The log  $L_1$ -error in the log-likelihood estimates and the estimates of the gradient with respect to  $\phi$  in the LGSS model with  $\sigma_e = 0.1$  (left) and  $\sigma_e = 1$  (right). The bPF (black) and faPF (red) are evaluated by 1,000 MC iterations using a fixed data set with  $T = 100$ . (Color figure online)

4.1 Estimation of the log-likelihood and the gradient

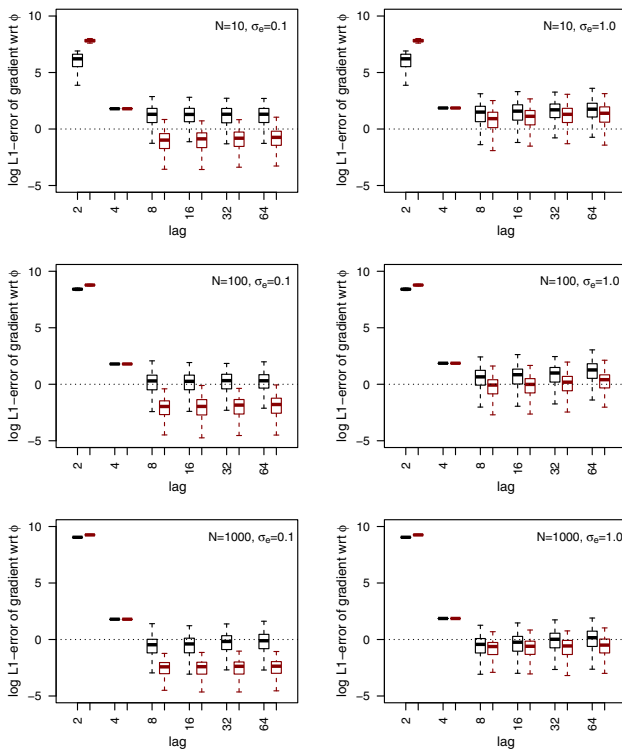
We begin by illustrating the use of the FL smoother for estimating the log-likelihood and the gradient. Here, we consider a linear Gaussian state space (LGSS) model given by

$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \phi x_t, \sigma_v^2), \tag{25a}$$

$$y_t|x_t \sim \mathcal{N}(y_t; x_t, \sigma_e^2). \tag{25b}$$

We generate two data realisations of length  $T = 100$  using parameters  $\theta^{(1)} = \{\phi, \sigma_v^2, \sigma_e^2\} = \{0.5, 1.0, 0.1^2\}$  and  $\theta^{(2)} = \{0.5, 1.0, 1.0\}$  with a known initial zero state. We use the lag  $\Delta = 5$  and run the PFs with systematic resampling (Carpenter et al. 1999).

For this model, we can compute the true values of the log-likelihood and the gradient by running an RTS smoother (Rauch et al. 1965). In Fig. 1, we present boxplots of the  $L_1$ -errors in the estimated log-likelihood and the gradient of the log-posterior with respect to  $\phi$ , evaluated at the true parameters. When  $\sigma_e = 0.1$ , we observe that the faPF has a large advantage over the bPF for all choices of  $N$ . When  $\sigma_e = 1.0$ , we get smaller difference in the error of the gradient estimates, but the log-likelihood estimates are still better for



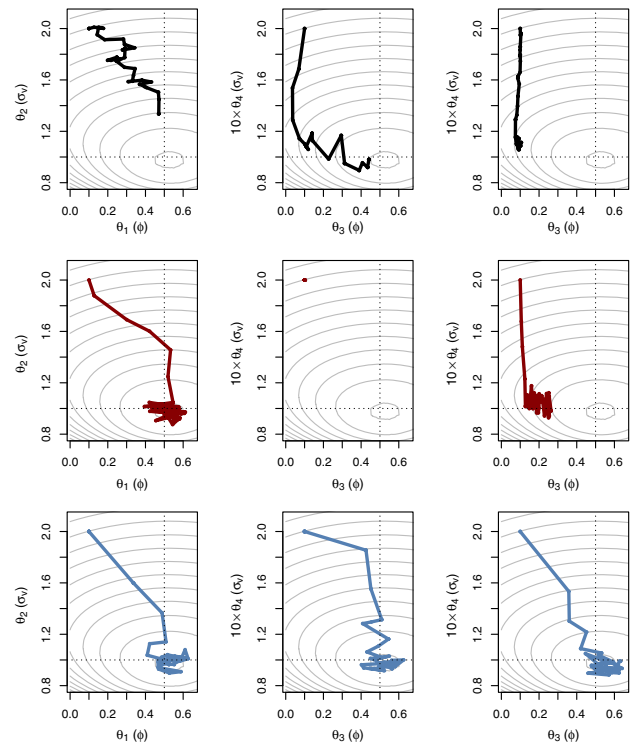
**Fig. 2** The log  $L_1$ -error in the estimates of the gradient with respect to  $\phi$  in the LGSS model with  $\sigma_e = 0.1$  (left) and  $\sigma_e = 1$  (right). The bPF (black) and faPF (red) are evaluated by 1,000 Monte Carlo iterations using a fixed data set with  $T = 100$ . (Color figure online)

the faPF. Similar results are also obtained for the gradient with respect to  $\sigma_v$ .

In Fig. 2, we present the error in the gradient estimates with respect to  $\phi$  using a varying lag  $\Delta$  and a varying number of particles  $N$ . The results are obtained by 1,000 Monte Carlo runs on a single data set generated from the previously discussed LGSS model with  $T = 100$ . We conclude again that faPF is preferable when available. The results are largely robust to the lag, as long as this is chosen large enough when using the faPF. A lag of about 12 seems to be a good choice for this model when  $T = 100$  and when using the faPF with systematic resampling.

#### 4.2 Burn-in and scale-invariance

Consider the problem of inferring  $\{\theta_1, \theta_2\} = \{\phi, \sigma_v\}$  in the LGSS model (25a, b). We simulate a single data set with parameters  $\theta^{(1)}$  (as defined in the previous section) of length  $T = 250$ . We use an uniform parameter prior over  $|\phi| < 1, \sigma_v > 0$  and initialise in  $\theta_0 = \{0.1, 2\}$ . We use faPF with systematic resampling,  $N = 100$  and  $\Delta = 12$ . Here, we use the standard version in Algorithm 2 to adjust the estimate of the Hessian in the cases when it is not PD, resulting in the PMH2 algorithm.



**Fig. 3** The trace plots of the first 50 steps using PMH0 (upper), PMH1 (middle) and PMH2 (lower). The dotted lines show the true parameters of the LGSS model. The gray contours show the log-posterior.

We adjust the step lengths  $\gamma$  to give an acceptance rate during a pilot run of between 0.7 and 0.8 in the stationary phase. We obtain  $\gamma = \{0.04, 0.065, 1.0\}$  for PMH $\{0, 1, 2\}$ , respectively. Note that a single step length is used for each proposal to simplify the tuning. Of course, different step lengths can be used for each parameter, and we could also use different step lengths during the burn-in and the stationary phase of the algorithm using the approach discussed in Sect. 2.2. As previously mentioned, the PMH2 algorithm avoids this (potentially difficult and time-consuming) procedure, by taking the local geometric information into account.

In the left column of Fig. 3, we present the first 50 iterations of the Markov chain from the three different algorithms. We note that the added information in the proposals of PMH1 and PMH2 aids the Markov chain in the burn-in phase. This results in that the Markov chains for the proposed algorithms reach the mode of the posterior quicker than the random walk used in PMH0.

To illustrate the scale invariance of the PMH2 algorithm, we reparametrise the LGSS model by  $\{\theta_3, \theta_4\} = \{\phi, \sigma_v/10\}$ . We keep the same settings as for the previous parametrisation and rerun the algorithms. From this run we obtain the middle column in Fig. 3. We see clearly that the PHM1-algorithm does not perform well and gets stuck at the initial parameter value. The reason is that the second component of the gradient is increased by a factor 10 for the rescaled model. Since



we still use the same step length, this will cause the PMH1 algorithm to overshoot the region of high posterior probability when proposing new values, and these will therefore never be accepted.

Finally, to improve the performance we recalibrate the three algorithms on the new parametrisation using the same procedure as before. We then obtain the new step lengths {0.005, 0.0075, 1.0}. The resulting Markov chains are presented in the right column of Fig. 3. Despite the new step lengths, PMH0 and PMH1 continue to struggle. The reason is that the step lengths are limited by the small posterior variance in the  $\theta_4$ -parameter, resulting in a very slow progression in the  $\theta_3$ -direction. Again, for PMH2, the added Hessian information is used to rescale the proposal in each dimension resulting in a more efficient exploration of the posterior than for PMH0 and PMH1.

### 4.3 The mixing of the Markov chains at stationarity

We continue by investigating the mixing of the Markov chains at stationarity using an estimate of the integrated autocorrelation time (IACT) given by

$$\widehat{\text{IACT}}(\theta_{1:M}) = 1 + 2 \sum_{k=1}^K \widehat{\rho}_k(\theta_{1:M}), \tag{26}$$

where  $\widehat{\rho}_k(\theta_{1:M})$  denotes the empirical autocorrelation at lag  $k$  of  $\theta_{1:M}$  (after the burn-in has been discarded). A low value of the IACT indicates that we obtain many uncorrelated samples from the target distribution, implying that the chain is mixing well. Here,  $K$  is determined as the first index for which the empirical autocorrelation satisfies  $|\widehat{\rho}_K(\theta_{1:M})| < 2/\sqrt{M}$ , i.e. when the coefficient is statistical insignificant.

We return to the LGSS model in (25a, b) with the original parameterisation  $\{\theta_1, \theta_2\} = \{\phi, \sigma_v\}$  using the same settings as before. A total of 25 data sets are generated using the parameters  $\theta^{(1)}$  and the algorithms are initialised at the true parameter values to avoid a long burn-in phase. The step sizes are determined using a series of pilot runs on the first generated dataset to minimise the total IACT for each algorithm. This is done to make a fair comparison between the different algorithms at their near *optimal* performance. The resulting step sizes are obtained as {0.08, 0.075, 1.50}.

Finally, we estimate the mixing in each of the 25 simulated data sets during  $M=30,000$  MCMC iterations (discarding the first 10,000 iterations as burn-in). The results are presented in Table 1, where the median and interquartile range (IQR; the distance between the 25 and 75 % quartiles) are presented for each PMH algorithm. Here, we present the results for the standard version of Algorithm 2.

We see that the added information decreases the IACT about 2 times for PMH1 and PMH2 compared with PMH0.

**Table 1** Median and IQR for the acceptance rate and IACT using different SMC algorithms

	Acc. rate	IACT ( $\phi$ )		IACT ( $\sigma_v$ )	
		Median	Median	IQR	Median
<b>PMH0</b>					
bPF(500)	0.02	257	146	265	371
bPF(1000)	0.06	83	129	79	118
bPF(2000)	0.15	29	23	15	24
faPF(50)	0.37	9	8	8	5
faPF(100)	0.38	9	6	7	4
faPF(200)	0.38	7	6	7	4
<b>PMH1</b>					
bPF(500)	0.02	187	271	203	347
bPF(1000)	0.10	64	85	49	72
bPF(2000)	0.22	23	16	12	24
faPF(50)	0.58	<b>3</b>	2	<b>3</b>	1
faPF(100)	0.59	4	2	<b>3</b>	1
faPF(200)	0.58	<b>3</b>	1	<b>3</b>	1
<b>PMH2</b>					
bPF(500)	0.03	170	211	164	190
bPF(1000)	0.10	59	73	65	80
bPF(2000)	0.24	13	10	19	17
faPF(50)	0.66	<b>3</b>	1	4	2
faPF(100)	0.66	<b>3</b>	1	5	2
faPF(200)	0.66	<b>3</b>	1	4	2

The values are computed using 25 different data sets from the LGSS model  
 Bold values indicate the best values in each column (i.e. the lowest values of IACT)

We conclude that the extra information brought by the gradient and the Hessian improves the mixing of the Markov chains in this model, which results in a more efficient exploration of the posterior. Note that, for this parametrisation of the LGSS model the posterior is quite isotropic (which can also be seen in the left column of Fig. 3). Hence, the conditions are in fact rather favourable for PMH0 and PMH1.

### 4.4 Parameter inference in a Poisson count model

In this section, we analyse the annual number of major earthquakes<sup>1</sup> (over 7 on the Richter scale) during the period from year 1900 to 2014. Following Langrock (2011), we model

<sup>1</sup> The data is obtained from the Earthquake Data Base System of the U.S. Geological Survey, which can be accessed at <http://earthquake.usgs.gov/earthquakes/eqarchives/>.

**Table 2** Median and IQR for the acceptance rate and IACT using different number of particles

	Version	SMC alg.	Acc. rate	IACT ( $\phi$ )		IACT ( $\sigma$ )		IACT ( $\beta$ )	
			Median	Median	IQR	Median	IQR	Median	IQR
PMH0	Standard	bPF(500)	0.26	497	712	16	3	2,639	1,163
	Standard	bPF(1000)	0.30	89	150	<b>15</b>	3	2,680	438
	Pre-cond.	bPF(500)	0.73	1,063	912	982	274	2,815	706
	Pre-cond.	bPF(1000)	0.81	831	290	1,106	671	2,374	1,446
PMH1	Standard	bPF(500)	0.76	665	442	277	162	2,651	364
	Standard	bPF(1000)	0.82	490	134	205	30	2,875	1,007
	Pre-cond.	bPF(500)	0.62	266	187	<b>9</b>	3	1,728	1638
	Pre-cond.	bPF(1000)	0.70	98	209	<b>9</b>	3	1,480	1732
PMH2	Standard	bPF(500)	0.24	91	17	53	14	222	37
	Standard	bPF(1000)	0.28	60	14	47	17	139	59
	Hybrid	bPF(500)	0.45	20	3	17	4	30	15
	Hybrid	bPF(1000)	0.49	<b>17</b>	4	18	3	<b>23</b>	5

The values are computed using 10 runs on the Earthquake count data model  
 Bold values indicate the best values in each column (i.e. the lowest values of IACT)

the data using

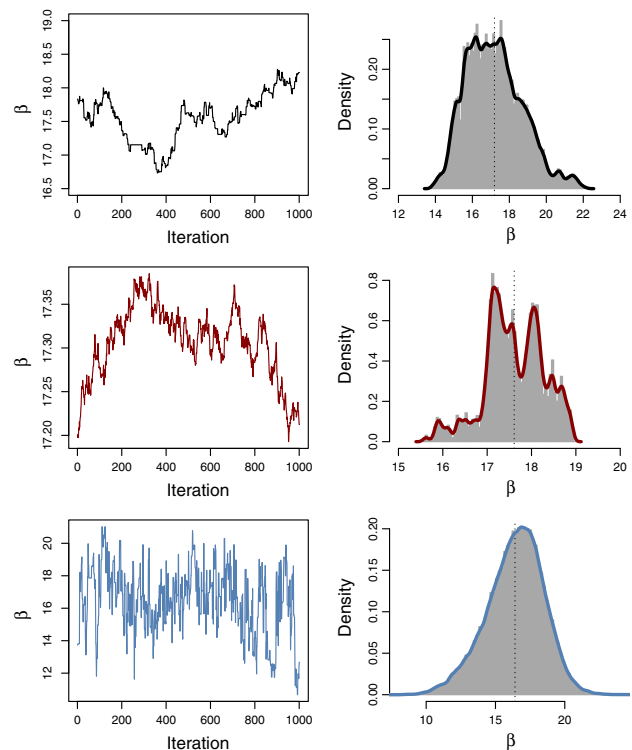
$$x_{t+1}|x_t \sim N(x_{t+1}; \phi x_t, \sigma^2), \tag{27a}$$

$$y_t|x_t \sim P(y_t; \beta \exp(x_t)), \tag{27b}$$

with parameters  $\theta = \{\phi, \sigma, \beta\}$  and uniform priors over  $|\phi| < 1$ ,  $\sigma > 0$  and  $\beta > 0$ . Here,  $P(\lambda)$  denotes a Poisson distribution with parameter  $\lambda$ .

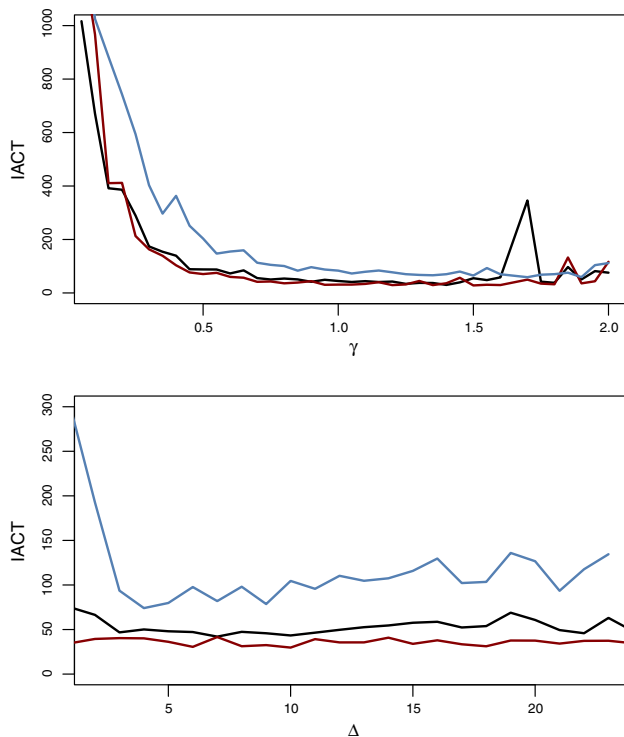
We repeat the procedure from the previous subsection and obtain the step lengths  $\{0.06, 0.006, 0.85\}$ . Here, we use  $M = 30,000$  MCMC iterations (discarding the first 10,000 iterations as burn-in), the bPF with systematic resampling,  $\Delta = 12$ ,  $\theta_0 = \{0.5, 0.5, 18\}$  and  $L = 2,500$ . In this model, the estimate of the negative Hessian is often non-PD (during about half of the iterations) and the choice of regularisation is therefore important. To explore the properties of the regularisation, we apply both the standard and hybrid version of the PMH2 algorithm discussed in Sect. 3.5. We compare these methods to standard and pre-conditioned versions of the the PMH0 and PMH1 algorithms, using the sample posterior covariance matrix calculated in the same manner as for the hybrid PMH2 algorithm.

In Table 2, we present the resulting acceptance rates and IACT values for each parameter and algorithm. We note the large decrease in IACT for  $\beta$  when using the Hessian information, where the hybrid PMH2 seems to perform better than standard version for this model. The improved mixing by using PMH2 is due to the scale invariance property, as the parameter  $\beta$  is at least an order of magnitude larger than  $\phi$  and  $\sigma$  (c.f. Fig. 3). Note that a reparameterisation or using separate step lengths for the parameters could possibly have helped in improving the mixing in  $\beta$  for the standard versions of PMH0 and PMH1.



**Fig. 4** Part of the trace (left) and posterior estimates (right) for the  $\beta$  parameter in the earthquake count model using standard versions of PMH0 (upper), PMH1 (middle) and hybrid version of PMH2 (lower). Dotted lines indicate the posterior means.

Using the standard and hybrid version of PMH2, decreases the overall computational cost by a factor of about 100 for a specific number of effective samples. The poor performance of the pre-conditioned algorithms is probably due to that the



**Fig. 5** The IACT for  $\phi$  (black),  $\sigma$  (red) and  $\beta$  (blue) for varying step sizes  $\gamma$  (upper) and lag  $\Delta$  (lower). The values are computed as the median of 10 runs using standard PMH2 with the same data. (Color figure online)

sample posterior covariance matrix does not fully capture the geometry of the posterior distribution.

In Fig. 4, we present the trace and posterior estimates for  $\beta$  using the standard versions of PMH0 and PMH1 as well as hybrid PMH2. The posterior estimates are obtained by pooling the 10 parallel Markov chains after the burn-ins have been discarded. We see that the traces behave rather differently with hybrid PMH2 exploring the space well compared with the other methods.

Using the parameter posterior estimate, we can compute point estimates for the parameters of the model. The posterior mean for hybrid PMH2 is obtained as  $\{0.88, 0.15, 16.58\}$  with standard deviations  $\{0.07, 0.03, 2\}$ . The parameter estimate is comparable to the estimate  $\{0.88, 0.15, 17.65\}$  obtained by a maximum likelihood-based method using the same data and model in Dahlin (2014, Example 4.9).

#### 4.5 Robustness in the lag and step size

The PMH2 algorithm requires a number of parameters to be select by the user for each parameter inference problem. It is therefore interesting to discuss the robustness of the method with respect to these parameters. In the previous illustrations, we have seen that the number of particles  $N$  is an important factor in determining the mixing.

Two other important parameters are the step length  $\gamma$  and the lag in the FL-smoother  $\Delta$ . To illustrate the impact of these quantities on the IACT, we return to the Earthquake model in (27a, b) using the standard PMH2 algorithm with the same settings but with  $M = 15,000$  (discarding the first 5,000 iterations as burn-in) and  $N = 1,500$ . In Fig. 5, we present the IACT for the three parameters in the model when varying  $\gamma$  and  $\Delta$ , keeping everything else fixed. The standard PMH2 algorithm seems to be rather robust to both the choice of  $\Delta$  and  $\gamma$  after a certain threshold. Recall the discussion in Sect. 4.1 for the FL smoother. We conclude that a suitable standard choice for the step length could be  $\gamma = 1$ , which can be fine tuned if the performance is not good enough. This recommendation is also common in the literature concerning Newton-type algorithms.

## 5 Discussion and future work

Adding the gradient and Hessian information to the PMH proposal can have beneficial results including: (i) a shorter burn-in phase, (ii) a better mixing of the Markov chain, and (iii) scale-invariance of the proposal which simplifies tuning. The latter point is true in particular for PMH2, since this method takes the local curvature of the posterior into account, effectively making the method invariant to affine transformations.

It is common to distinguish between two phases of MCMC algorithms: the burn-in and stationary phases. We have seen empirically that the proposed methods can improve upon the original PMH0 during both of these phases but the *best* choices for the step lengths can differ between these two phases. Typically, a smaller step length is preferred during burn-in and a larger during stationarity (the opposite holds for PMH0). The reason for this is that during burn-in, the (natural) gradient information will heavily skew the proposal in a direction of increasing posterior probability. That is, the methods tend to be *aggressive* and propose large steps to make rapid progression toward regions of high posterior probability. While this is intuitively appealing, the problem is that we require the Markov chains to be reversible at all times. The reverse of these large steps can have very low probability which prevents them from being accepted.

One interesting direction for future work is therefore to pursue adaptive algorithms (see e.g. Andrieu and Thoms 2008; Peters et al. 2010; Pitt et al. 2012), to automatically tune the step lengths during the different phases of the algorithms. Another interesting possibility is to relax the reversibility requirement during burn-in; see Diaconis et al (2000) for a related reference. This would cause the methods to behave like optimisation procedures during the initial phase, but transition into samplers during the second phase.

Finally, another very interesting direction for future work is to extend the proposed methods to develop a particle-version of the manifold Hamiltonian Monte Carlo (mHMC) algorithm (Duane et al. 1987; Neal 2010; Girolami and Calderhead 2011). The reason for this is motivated by the large improvement in mixing seen by e.g. Neal (2010), Girolami and Calderhead (2011) for high dimensional problems in *vanilla* MH sampling.

**Acknowledgments** This work was supported by: Learning of complex dynamical systems (Contract number: 637-2014-466) and Probabilistic modeling of dynamical systems (Contract number: 621-2013-5524) and CADICS, a Linnaeus Center, all funded by the Swedish Research Council.

## References

- Andrieu, C., Roberts, G.O.: The pseudo-marginal approach for efficient Monte Carlo computations. *Ann. Stat.* **37**(2), 697–725 (2009)
- Andrieu, C., Thoms, J.: A tutorial on adaptive MCMC. *Stat. Comput.* **18**(4), 343–373 (2008)
- Andrieu, C., Vihola, M.: Convergence properties of pseudo-marginal Markov chain Monte Carlo algorithms. Pre-print [arXiv:1012.1484v1](https://arxiv.org/abs/1012.1484v1) (2011)
- Andrieu, C., Doucet, A., Holenstein, R.: Particle Markov chain Monte Carlo methods. *J. R. Stat. Soc.* **72**(3), 269–342 (2010)
- Beaumont, M.A.: Estimation of population growth or decline in genetically monitored populations. *Genetics* **164**(3), 1139–1160 (2003)
- Cappé, O., Moulines, E., Rydén, T.: *Inference in Hidden Markov Models*. Springer, Berlin (2005)
- Carpenter, J., Clifford, P., Fearnhead, P.: Improved particle filter for nonlinear problems. *IEE Proc. Radar Sonar Navig.* **146**(1), 2–7 (1999)
- Dahlin, J.: Sequential Monte Carlo for inference in nonlinear state space models. Licentiate’s thesis no. 1652, Linköping University (2014)
- Dahlin, J., Lindsten, F., Schön, T.B.: Particle Metropolis Hastings using Langevin dynamics. In: Proceedings of the 38th International Conference on Acoustics, Speech, and Signal Processing (ICASSP), Vancouver (2013)
- Dahlin, J., Lindsten, F., Schön, T.B.: Second-order particle MCMC for Bayesian parameter inference. In: Proceedings of the 19th IFAC World Congress, Cape Town (2014)
- Del Moral, P.: Feynman-Kac Formulae—Genealogical and Interacting Particle Systems with Applications. Probability and its applications. Springer, Berlin (2004)
- Del Moral, P., Doucet, A., Singh, S.: Forward smoothing using sequential Monte Carlo. Pre-print [arXiv:1012.5390v1](https://arxiv.org/abs/1012.5390v1) (2010)
- Diaconis, P., Holmes, S., Neal, R.: Analysis of a nonreversible Markov chain sampler. *Ann. Appl. Probab.* **10**(3), 685–1064 (2000)
- Doucet, A., Johansen, A.: A tutorial on particle filtering and smoothing: fifteen years later. In: Crisan, D., Rozovsky, B. (eds.) *The Oxford Handbook of Nonlinear Filtering*. Oxford University Press, Oxford (2011)
- Doucet, A., Jacob, P., Johansen, A.M.: Discussion on Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *J. R. Stat. Soc.* **73**(2), 162 (2011)
- Doucet, A., Pitt, M.K., Deligiannidis, G., Kohn, R.: Efficient implementation of Markov chain Monte Carlo when using an unbiased likelihood estimator. [arXiv.org](https://arxiv.org/abs/1210.1871v3), Pre-print [arXiv:1210.1871v3](https://arxiv.org/abs/1210.1871v3) (2012)
- Doucet, A., Jacob, P.E., Rubenthaler, S.: Derivative-free estimation of the score vector and observed information matrix with application to state-space models. Pre-print [arXiv:1304.5768v2](https://arxiv.org/abs/1304.5768v2) (2013)
- Duane, S., Kennedy, A.D., Pendleton, B.J., Roweth, D.: Hybrid Monte Carlo. *Phys. Lett. B* **195**(2), 216–222 (1987)
- Everitt, R.G.: Bayesian parameter estimation for latent Markov random fields and social networks. *J. Comput. Gr. Stat.* **21**(4), 940–960 (2012)
- Flury, T., Shephard, N.: Bayesian inference based only on simulated likelihood: particle filter analysis of dynamic economic models. *Econom. Theory* **27**(5), 933–956 (2011)
- Girolami, M., Calderhead, B.: Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *J. R. Stat. Soc.* **73**(2), 1–37 (2011)
- Golightly, A., Wilkinson, D.J.: Bayesian parameter inference for stochastic biochemical network models using particle Markov chain Monte Carlo. *Interface Focus* **1**(6), 807–820 (2011)
- Gordon, N.J., Salmond, D.J., Smith, A.F.M.: Novel approach to nonlinear/non-Gaussian Bayesian state estimation. *IEEE Proc. Radar Signal Process.* **140**(2), 107–113 (1993)
- Kitagawa, G., Sato, S.: Monte Carlo smoothing and self-organising state-space model. In: Doucet, A., de Fretias, N., Gordon, N. (eds.) *Sequential Monte Carlo methods in practice*, pp. 177–195. Springer, Berlin (2001)
- Langrock, R.: Some applications of nonlinear and non-Gaussian state-space modelling by means of hidden Markov models. *J. Appl. Stat.* **38**(12), 2955–2970 (2011)
- Neal, R.M.: MCMC using Hamiltonian dynamics. In: Brooks, S., Gelman, A., Jones, G., Meng, X.L. (eds.) *Handbook of Markov Chain Monte Carlo*. Chapman & Hall, London (2010)
- Nemeth, C., Fearnhead, P.: Particle Metropolis adjusted Langevin algorithms for state-space models. Pre-print [arXiv:1402.0694v1](https://arxiv.org/abs/1402.0694v1) (2014)
- Nocedal, J., Wright, S.: *Numerical Optimization*, 2nd edn. Springer, Berlin (2006)
- Olsson, J., Cappé, O., Douc, R., Moulines, E.: Sequential Monte Carlo smoothing with application to parameter estimation in nonlinear state space models. *Bernoulli* **14**(1), 155–179 (2008)
- Peters, G.W., Hosack, G.R., Hayes, K.R.: Ecological non-linear state space model selection via adaptive particle Markov chain Monte Carlo. Pre-print [arXiv:1005.2238v1](https://arxiv.org/abs/1005.2238v1) (2010)
- Pitt, M.K., Shephard, N.: Filtering via simulation: auxiliary particle filters. *J. Am. Stat. Assoc.* **94**(446), 590–599 (1999)
- Pitt, M.K., Silva, R.S., Giordani, P., Kohn, R.: On some properties of Markov chain Monte Carlo simulation methods based on the particle filter. *J. Econom.* **171**(2), 134–151 (2012)
- Poyiadjis, G., Doucet, A., Singh, S.S.: Particle approximations of the score and observed information matrix in state space models with application to parameter estimation. *Biometrika* **98**(1), 65–80 (2011)
- Rauch, H.E., Tung, F., Striebel, C.T.: Maximum likelihood estimates of linear dynamic systems. *AIAA J.* **3**(8), 1445–1450 (1965)
- Robert, C.P., Casella, G.: *Monte Carlo Statistical Methods*, 2nd edn. Springer, Berlin (2004)
- Roberts, G.O., Rosenthal, J.S.: Optimal scaling of discrete approximations to Langevin diffusions. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **60**(1), 255–268 (1998)
- Roberts, G.O., Stramer, O.: Langevin diffusions and Metropolis–Hastings algorithms. *Methodol. Comput. Appl. Probab.* **4**(4), 337–357 (2003)
- Roberts, G.O., Gelman, A., Gilks, W.R.: Weak convergence and optimal scaling of random walk Metropolis algorithms. *Ann. Appl. Probab.* **7**(1), 110–120 (1997)
- Sherlock, C., Thiery, A.H., Roberts, G.O., Rosenthal, J.S. On the efficiency of pseudo-marginal random walk Metropolis algorithms. Pre-print [arXiv:1309.7209v1](https://arxiv.org/abs/1309.7209v1) (2013)