# ON THE USE OF SEQUENTIAL MONTE CARLO METHODS FOR APPROXIMATING SMOOTHING FUNCTIONALS, WITH APPLICATION TO FIXED PARAMETER ESTIMATION *

JIMMY OLSSON[1], OLIVIER CAPPÉ[2], RANDAL DOUC[3] AND ÉRIC MOULINES[4]

**Abstract.** Sequential Monte Carlo (SMC) methods have demonstrated a strong potential for inference on the state variables in Bayesian dynamic models. In this context, it is also often needed to calibrate model parameters. To do so, we consider block maximum likelihood estimation based either on EM (Expectation-Maximization) or gradient methods. In this approach, the key ingredient is the computation of smoothed sum functionals of the hidden states, for a given value of the model parameters. It has been observed by several authors that using standard SMC methods for this smoothing task requires a substantial number of particles and may be unreliable for larger observation sample sizes. We introduce a simple variant of the basic sequential smoothing approach based on forgetting ideas. This modification, which is transparent in terms of computation time, reduces the variability of the approximation of the sum functional. Under suitable regularity assumptions, it is shown that this modification indeed allows a tighter control of the $L_p$ error of the approximation.

## INTRODUCTION

We consider a particular form of Bayesian dynamic models known as general *state-space* or *hidden Markov* models, which consists of a, partially observed, discrete-time bivariate process $\{(X_k, Y_k); k \geq 0\}$, where $\{X_k\}$ is a Markov chain and, conditional on $\{X_k\}$, $\{Y_k\}$ is a sequence of independent random observations such that the distribution of $Y_k$ is governed by $X_k$ only. This class of models is of great practical importance in fields as diverse as computational biology, computer vision or quantitative finance [3]. In the following, we consider a statistical model dominated by a measure $\mu(dx) \otimes \lambda(dy)$ and denote by $\nu_\theta(x)$ and $q_\theta(x, x')$ the initial and Markov transition densities of $\{X_k\}$, respectively, and by $g_\theta(x, y)$ the transition density of $Y_k$ given $X_k$. The task of interest consists in estimating the parameter $\theta$ from a fixed set of observations $\boldsymbol{Y}_{0:n} \triangleq (Y_0, \ldots, Y_n)$.

Due to the introduction of the latent state process $\{X_k\}$, the log-likelihood of the observations $\ell_n(\theta) \triangleq \log p_\theta(\boldsymbol{Y}_{0:n})$ cannot be easily maximized—even numerically—, except in some particular cases. A generic approach for latent variable models is provided by the Expectation-Maximization (EM) algorithm [6] which

requires the evaluation of the intermediate quantity

$$
\begin{aligned}
\mathcal{Q}_\theta(\theta') &\triangleq \mathbb{E}_\theta \left[ \log p_{\theta'}(\boldsymbol{X}_{0:n}, \boldsymbol{Y}_{0:n}) \middle| \boldsymbol{Y}_{0:n} \right] \\
&= \mathbb{E}_\theta \left[ \sum_{k=0}^{n-1} \log q_{\theta'}(X_k, X_{k+1}) \middle| \boldsymbol{Y}_{0:n} \right] + \mathbb{E}_\theta \left[ \sum_{k=0}^{n} \log g_{\theta'}(X_k, Y_k) \middle| \boldsymbol{Y}_{0:n} \right] + \mathbb{E}_\theta \left[ \log \nu_{\theta'}(X_0) \middle| \boldsymbol{Y}_{0:n} \right] ,
\end{aligned} \quad (1)
$$

where $\mathbb{E}_\theta$ denotes the expectation under the model distribution parameterized by $\theta$. Due to the underlying Markovian structure of the model, (1) is of the general form

$$
\gamma_n = \mathbb{E}_\theta \left[ \sum_{k=0}^{n-1} s_k(X_k, X_{k+1}) \middle| \boldsymbol{Y}_{0:n} \right] , \quad (2)
$$

where $\{s_k; k \geq 0\}$ is a sequence of functions, which depends on the observed values $\boldsymbol{Y}_{0:n}$ and on the parameter value $\theta'$. Although not of direct interest here, note that when $\nu_\theta$, $q_\theta$, and $g_\theta$ belong to exponential families, the EM algorithm may be rewritten in such a way that $s_k$ depends on $(Y_k, Y_{k+1})$ only; the function $\sum_{k=0}^{n-1} s_k(X_k, X_{k+1})$ is then referred to as the *complete-data sufficient statistic* [6]. The quantity featured in (2) is referred to as a *smoothing functional* and our primary goal with this paper is to discuss efficient schemes for approximating such quantities with sequential Monte Carlo (henceforth abbreviated to SMC) simulations.

Note that (2) also appears as a key ingredient when using gradient-based approaches to optimize the log-likelihood. Indeed, under appropriate differentiability assumptions, the *Fisher identity* states that

$$
\begin{aligned}
\nabla_\theta \ell_n(\theta) = \nabla_{\theta'} \mathcal{Q}_\theta(\theta')|_{\theta'=\theta} &= \mathbb{E}_\theta \left[ \nabla_\theta \log p_\theta(\boldsymbol{X}_{0:n}, \boldsymbol{Y}_{0:n}) \middle| \boldsymbol{Y}_{0:n} \right] \\
&= \mathbb{E}_\theta \left[ \sum_{k=0}^{n-1} \nabla_\theta \log q_\theta(X_k, X_{k+1}) \middle| \boldsymbol{Y}_{0:n} \right] \\
&\quad + \mathbb{E}_\theta \left[ \sum_{k=0}^{n} \nabla_\theta \log g_\theta(X_k, Y_k) \middle| \boldsymbol{Y}_{0:n} \right] + \mathbb{E}_\theta \left[ \nabla_\theta \log \nu_\theta(X_0) \middle| \boldsymbol{Y}_{0:n} \right] ,
\end{aligned} \quad (3)
$$

where $\nabla_\theta$ denotes the gradient with respect to $\theta$. Hence the score function (gradient of the log-likelihood) also has the form given in (2).[1] As a side comment, note that there is an important practical difference between the use of (1) and (3) for inference: As discussed above, when (1) is used in exponential families, it directly provides a parameter update equation which may be used to iteratively increase the log-likelihood. If (1) is approximated using stochastic simulations, one then obtains a so called Monte Carlo version of the EM algorithm. In contrast, approximating (3) with a simulation-based method only provides a noisy evaluation of the gradient of the log-likelihood in one point. To actually optimize the likelihood, one typically resorts to a stochastic approximation (or Robbins-Monro) scheme [3].

The rest of the paper is organized as follows: the principle of SMC methods is briefly recalled in Section 1; in Section 2, we discuss the limitations of SMC when applied to the approximation of smoothing functionals of the form given in (2) and propose a simple solution based on forgetting ideas; Section 3 provides some backup to our claim that the proposed approximation is significantly more reliable than the use of the standard SMC trajectory-based approximation.

---

[1]There exists a generic recursive implementation of (2), which can be used to compute $\gamma_n$ for increasing values of $n$ [10, 14, 18]. In our context however, this observation is not very useful as the recursive rewriting does not modify the fundamental nature of (2) and its implementation is only available in those specific cases where the filtering and smoothing relations may be implemented numerically.

# 1. Sequential Monte Carlo

Sequential Monte Carlo, also known as particle filtering, approximates the exact filtering and smoothing relations by propagating particle trajectories in the state space of the hidden chain. For more details, we refer to [3, 8, 9, 15, 17] and simply present below the systematic sequential importance sampling with resampling algorithm. In order to keep the notations simple, we fix the model parameters and omit $\theta$ in the following.

At time zero, a number $N$ of particles $\{\xi_0^{N,i}(0); 1 \leq i \leq N\}$ are drawn from a common probability density $\eta$. These *initial particles* are assigned the *importance weights* $\omega_0^{N,i} \triangleq W_0[\xi_0^{N,i}(0)]$, $i = 1, \ldots, N$, where $W_0(x) \triangleq g_0(x)\,\nu(x)/\eta(x)$, providing $\sum_{i=1}^{N} \omega_0^{N,i} f[\xi_0^{N,i}(0)] / \sum_{i=1}^{N} \omega_0^{N,i}$ as an importance sampling estimate of $\mathbb{E}\,[\,f(X_0)|\,Y_0]$. We define by $\boldsymbol{\xi}_m^{N,i} = (\xi_m^{N,i}(0), \ldots, \xi_m^{N,i}(m))$ the $i$th particle path, with in particular $\boldsymbol{\xi}_0^{N,i} = \xi_0^{N,i}(0)$.

At time $k$, we simulate $\xi_{k+1}^{N,i}(k+1) \sim r_k(\xi_k^{N,i}(k), \cdot)$ and set $\boldsymbol{\xi}_{k+1}^{N,i} = (\boldsymbol{\xi}_k^{N,i}, \xi_{k+1}^{N,i}(k+1))$, where $r_k$ is an instrumental transition density. In this *mutation step*, the new particles are simulated independently of each other. A popular choice is to set $r_k = q$, yielding the so-called *bootstrap filter*; more sophisticated techniques involve proposals depending on the new observation $y_{k+1}$. When the new observation is available, the importance weights are updated according to the formula $\omega_{k+1}^{N,i} = \omega_k^{N,i} \times W_{k+1}[\xi_k^{N,i}(k), \xi_{k+1}^{N,i}(k+1)]$, where $W_{k+1}(x, x') \triangleq q(x, x')g_{k+1}(x')/r_k(x, x')$. For a function $f$ of the state variables up to time $k+1$, the smoothed expectation $\mathbb{E}\,[\,f(\boldsymbol{X}_{0:k+1})|\,\boldsymbol{Y}_{0:k+1}]$ may be approximated by

$$\frac{1}{\sum_{i=1}^{N} \omega_{k+1}^{N,i}} \sum_{j=1}^{N} \omega_{k+1}^{N,j} f(\boldsymbol{\xi}_{k+1}^{N,j})\ .$$

As it is well established, the previous scheme fails because the distribution of the importance weights becomes more and more skewed as $k$ increases. To prevent degeneracy, a *selection mechanism* should be introduced. In its simpler form, this mechanism amounts to resample, when needed, the propagated particles by drawing, conditionally independently, indices $I_k^{N,1}, \ldots, I_k^{N,N}$ in the set $\{1, \ldots, N\}$ multinomially with respect to the normalized weights, that is, $\mathbb{P}(I_k^{N,i} = j) = \omega_k^{N,j} / \sum_{i=1}^{N} \omega_k^{N,i}$. Now, a new particle cloud $\{\tilde{\boldsymbol{\xi}}_k^{N,i}; 1 \leq i \leq N\}$ is formed by setting $\tilde{\boldsymbol{\xi}}_k^{N,j} = \boldsymbol{\xi}_k^{N,I_k^{N,j}}$. After the resampling procedure, the weights are all reset to $\tilde{\omega}_k^{N,i} = 1$, yielding the unweighted estimate

$$\frac{1}{N} \sum_{i=1}^{N} f(\tilde{\boldsymbol{\xi}}_k^{N,i})$$

of $\mathbb{E}\,[\,f(\boldsymbol{X}_{0:k+1})|\,\boldsymbol{Y}_{0:k+1}]$. In the following we assume that the resampling is done systematically at each step and we will use the notation $\boldsymbol{\xi}_k^{N,i}$ (without the tilde) to denote the $i$th particle path at time $k$ *after resampling*; accordingly the weight $\omega_k^{N,i}$ is then equal to 1. Note that the resampling mechanism may modify the whole trajectory of the particles, implying that in general, for $k \leq n$, $\xi_n^{N,i}(k) \neq \xi_{n+1}^{N,i}(k)$. Here again the multinomial resampling method is not the only conceivable way to carry out the selection step and we refer to [3, 8] for further reading on the topic.

# 2. The Fixed-Lag Approximation

For a function of the form given in (2), the natural SMC estimator is given by

$$\hat{\gamma}_n^N = \frac{1}{N} \sum_{i=1}^{N} \sum_{k=0}^{n-1} s_k \left[ \xi_n^{N,i}(k), \xi_n^{N,i}(k+1) \right]\ . \tag{4}$$

It is obvious that due to the additive nature of the functional, storing the whole particle trajectories is indeed not required: only the current particle positions $\xi_n^{N,i}(n)$ and the value of the functional along the trajectory $\boldsymbol{\xi}_n^{N,i}$ are really needed. Thus, the method necessitates only minor adaptations once the particle filter has been

implemented. Unfortunately, it has been noted by several authors [1, 3] that this method is not as stable as expected when $n$ increases due to the high variability of the part of the approximation $\hat{\gamma}_n^N$ that pertains to time indices $k \ll n$. The origin of this observation is illustrated in Figure 1. It is seen that the successive resamplings imply that, when $k \ll n$, the set $\{\xi_n^{N,i}(k); i = 1, \ldots, N\}$ is composed of a few elements only, hence yielding poorly reliable Monte Carlo estimates.
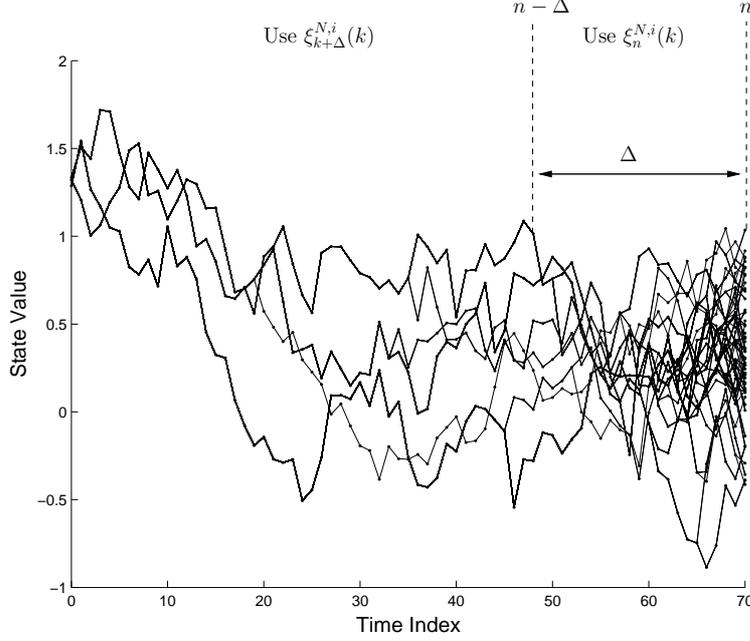


FIGURE 1. Principle of the fixed-lag approximation.

Several variants of the basic SMC approach presented in Section 1 have been proposed in the literature to improve the reliability of smoothing estimates. However these approaches compromise the sequential nature of the algorithm by requiring a backward pass through the data [2, 11–13]. In addition, these approaches are intended for estimating precisely the distribution of each state variable $X_k$ given the observations $\boldsymbol{Y}_{0:n}$ whereas we only want to estimate the smoothed expectation of a specific additive state functional. To improve the reliability of the smoothing estimate $\hat{\gamma}_n^N$, the proposed solution simply consists in substituting, in (4), $\xi_{k+\Delta}^{N,i}(k)$ for $\xi_n^{N,i}(k)$, whenever $k \leq n - \Delta$, where $\Delta$ is a fixed positive lag. The new *fixed-lag* estimator is thus given by

$$\hat{\gamma}_n^{N,\Delta_n} \triangleq \frac{1}{N} \sum_{i=1}^{N} \sum_{k=0}^{n-1} s_k \left[ \xi_{(k+\Delta)\wedge n}^{N,i}(k), \xi_{(k+\Delta)\wedge n}^{N,i}(k+1) \right] . \tag{5}$$

It is easily verified that this estimator has exactly the same numerical complexity as $\hat{\gamma}_n^N$ and can be implemented by storing the recent history of the particles $\{\xi_n^{N,i}(n-m); 0 \leq m \leq \Delta, 1 \leq i \leq N\}$.

The idea behind (5) is that for sufficiently large values of $\Delta$, $\mathbb{E}[s_k(X_k, X_{k+1})|\boldsymbol{Y}_{0:k+\Delta}]$ and $\mathbb{E}[s_k(X_k, X_{k+1})|\boldsymbol{Y}_{0:n}]$ (for $n \geq k+\Delta$) should be close. Such *forgetting* properties of the smoothing relations are indeed instrumental in studying both the statistical properties of likelihood-based parameter estimates [7] and the long-term stability of the SMC approach [5]. It is thus not unrealistic to assume that those forgetting properties do hold (see also [1] for a related use of the same idea).

In practice, it has been observed on several examples that the lag $\Delta$ controls a *bias/variance tradeoff*: when $\Delta$ is too small, $\hat{\gamma}_n^{N,\Delta_n}$ is a biased estimate of $\gamma_n$; as $\Delta$ is augmented, the bias disappears but the variance of

$\hat{\gamma}_n^{N,\Delta_n}$ raises. For a large range of values of $\Delta$ (usually between 10 and 100), the Monte Carlo variance of the error between $\hat{\gamma}_n^{N,\Delta_n}$ and $\gamma_n$ is and order of magnitude lower than when using the basic SMC estimator $\hat{\gamma}_n^N$ [3, 16].

## 3. ERROR BOUNDS

In this section, we state (without proof) the main result of [16] which bounds the Monte Carlo error $\hat{\gamma}_n^{N,\Delta_n}$ and $\gamma_n$ under strong mixing assumptions that guarantee that forgetting does hold for the exact smoothing distributions in the model under consideration. We begin by stating the required assumptions.

**(A1)**   *(i) $\sigma_- \triangleq \inf_{\theta\in\Theta} \inf_{(x,x')\in\mathsf{X}^2} q_\theta(x,x') > 0$, $\sigma_+ \triangleq \sup_{\theta\in\Theta} \sup_{(x,x')\in\mathsf{X}^2} q_\theta(x,x') < \infty$.*
*(ii) For all $y \in \mathsf{Y}$, $\sup_{\theta\in\Theta} \|g_\theta(\cdot,y)\|_{\mathsf{X},\infty} < \infty$ and $\inf_{\theta\in\Theta} \int_{\mathsf{X}} g_\theta(x,y)\,\mu(\mathrm{d}x) > 0$.*
*(iii) For all $k \geq 1$, $\|W_k\|_{\mathsf{X}^2,\infty} < \infty$; in addition, $\|W_0\|_{\mathsf{X},\infty} < \infty$.*

Assumptions *(i)* and *(ii)* guarantee that the posterior chain (state variables conditioned upon the observation sequence) forgets its initial condition at a *uniform* geometric rate $\rho \triangleq 1 - \sigma_-/\sigma_+$ (see, e.g., [7]). The additional assumption *(iii)* ensures that, for the filtering estimate, the approximation error of the SMC algorithm described in Section 1 may be bounded uniformly in time by a quantity of order $1/\sqrt{N}$ (which does not depend on $n$) [3, 5]. In many cases, the third assumption is indeed implied by *(i)* and *(ii)*. For instance, for the bootstrap filter, $W_k(x) = g_\theta(x,Y_k)$ and *(ii)* is implied by *(iii)*. We are now ready to state the main theorem of [16].

**Theorem 3.1.** *Under assumption **(A1)**, for $n \geq 0$, the following hold true for all $\Delta_n \geq 0$ and $N \geq 1$.*
*(i) For all $p \geq 2$,*

$$
\mathbb{E}^{1/p}\left[\left|\hat{\gamma}_n^{N,\Delta} - \gamma_n\right|^p \middle| \boldsymbol{Y}_{0:n}\right] \leq 2\rho^\Delta \sum_{k=0}^{n-\Delta} \|s_k\|_{\mathsf{X}^2,\infty} +
$$

$$
\frac{B_p}{\sqrt{N}(1-\rho)} \sum_{k=0}^{n-1} \|s_k\|_{\mathsf{X}^2,\infty} \left[\frac{1}{\sigma_-} \sum_{m=1}^{(k+\Delta)\wedge n} \frac{\|W_m\|_{\mathsf{X}^2,\infty}\rho^{0\vee(k-m)}}{G_m} + \frac{\|W_0\|_{\mathsf{X},\infty}\rho^k}{G_0} + 1\right],
$$

*(ii)*

$$
\left|\mathbb{E}\left[\hat{\gamma}_n^{N,\Delta}\middle|\boldsymbol{Y}_{0:n}\right] - \gamma_n\right| \leq 2\rho^\Delta \sum_{k=0}^{n-\Delta} \|s_k\|_{\mathsf{X}^2,\infty} +
$$

$$
\frac{B}{N(1-\rho)^2} \sum_{k=0}^{n-1} \|s_k\|_{\mathsf{X}^2,\infty} \left[\frac{1}{\sigma_-^2} \sum_{m=1}^{(k+\Delta)\wedge n} \frac{\|W_m\|_{\mathsf{X}^2,\infty}^2\rho^{0\vee(k-m)}}{G_m^2} + \frac{\|W_0\|_{\mathsf{X},\infty}^2\rho^k}{G_0^2}\right].
$$

*Here $B_p$ and $B$ are universal constants such that $B_p$ depends on $p$ only and $G_0 \triangleq \int g_\theta(x,Y_0)\,\nu(dx)$ and $G_m \triangleq \int g_\theta(x,Y_m)\,\mu(dx)$ are finite observation-dependent constants.*

Under some additional assumptions, the dependence with respect to the observation sequence $\boldsymbol{Y}_{0:n}$ may be integrated out as to obtain unconditional error bounds which account for the variability of both the observation sequence and the SMC simulations [16]. For the purpose of illustrating the bounds of Theorem 3.1 using simple arguments, assume that all $\|s_k\|_{\mathsf{X}^2,\infty}$ and all fractions $\|W_k\|_{\mathsf{X}^2,\infty}/G_k$ are uniformly bounded in $k$. We then draw the conclusion that if the lag $\Delta$ is increased with $n$ at a rate of $\log n$, then the error is dominated by the variability due to the particle filter—the second term of 3.1(i)—which is of order $O(N^{-1/2}n\log n)$. In contrast, setting $\Delta = n$, that is, using the direct trajectory-based approximation, would result in a stochastic error of order $O(N^{-1/2}n^2)$. Hence, with only moderate requirements on the lag $\Delta$ (due to the exponential nature of the forgetting), the proposed approximation is significantly more accurate than the standard trajectory-based approximation, for a comparable computational cost.

## References

[1] C. Andrieu, A. Doucet, and V. B. Tadic. Online simulation-based methods for parameter estimation in non linear non gaussian state-space models. In *Proc. IEEE Conf. Decis. Control*, 2005.

[2] M. Briers, A. Doucet, and S. Maskell. Smoothing algorithms for state-space models. Technical Report TR-CUED-F-INFENG 498, University of Cambridge, Department of Engineering, 2004.

[3] O. Cappé, E. Moulines, and T. Rydén. *Inference in Hidden Markov Models*. Springer, 2005.

[4] P. Del Moral. *Feynman-Kac Formulae. Genealogical and Interacting Particle Systems with Applications*. Springer, 2004.

[5] P. Del Moral and A. Guionnet. On the stability of interacting processes with applications to filtering and genetic algorithms. *Annales de l'Institut Henri Poincaré*, 37:155–194, 2001.

[6] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. Ser. B*, 39(1):1–38 (with discussion), 1977.

[7] R. Douc, E. Moulines, and T. Rydén. Asymptotic properties of the maximum likelihood estimator in autoregressive models with Markov regime. *Ann. Statist.*, 32(5):2254–2304, 2004.

[8] A. Doucet, N. De Freitas, and N. Gordon, editors. *Sequential Monte Carlo Methods in Practice*. Springer, New York, 2001.

[9] A. Doucet, S. Godsill, and C. Andrieu. On sequential Monte-Carlo sampling methods for Bayesian filtering. *Stat. Comput.*, 10:197–208, 2000.

[10] R. J. Elliott, L. Aggoun, and J. B. Moore. *Hidden Markov Models: Estimation and Control*. Springer, New York, 1995.

[11] S. J. Godsill, A. Doucet, and M. West. Monte carlo smoothing for non-linear time series. *J. Am. Statist. Assoc.*, 50:438–449, 2004.

[12] G. Kitagawa. Monte-Carlo filter and smoother for non-Gaussian nonlinear state space models. *J. Comput. Graph. Statist.*, 1:1–25, 1996.

[13] M. Klaas, M. Briers, N. De Freitas, A. Doucet, S. Maskell, and D. Lang. Fast particle smoothing: If i had a million particles. In *23rd Int. Conf. Machine Learning (ICML)*, Pittsburgh, Pennsylvania, June 25-29 2006.

[14] F. Le Gland and L. Mevel. Recursive estimation in HMMs. In *Proc. IEEE Conf. Decis. Control*, pages 3468–3473, 1997.

[15] J. Liu and R. Chen. Sequential Monte-Carlo methods for dynamic systems. *J. Roy. Statist. Soc. Ser. B*, 93:1032–1044, 1998.

[16] J. Olsson, O. Cappé, R. Douc, and E. Moulines. Sequential Monte Carlo smoothing with application to parameter estimation in non-linear state space models. Technical report, Lund University, 2006. arXiv:math.ST/0609514.

[17] B. Ristic, M. Arulampalam, and A. Gordon. *Beyond Kalman Filters: Particle Filters for Target Tracking*. Artech House, 2004.

[18] O. Zeitouni and A. Dembo. Exact filters for the estimation of the number of transitions of finite-state continuous-time Markov processes. *IEEE Trans. Inform. Theory*, 34(4), July 1988.