

MULTI-INDEX SEQUENTIAL MONTE CARLO METHODS FOR PARTIALLY OBSERVED STOCHASTIC PARTIAL DIFFERENTIAL EQUATIONS

Ajay Jasra,¹ Kody J. H. Law,^{2,*} & Yaxian Xu¹

¹Computer, Electrical and Mathematical Science and Engineering Division, King Abdullah University of Science and Technology, Thuwal, 23955-6900, KSA.

²Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN, USA, 37831

³Department of Statistics & Applied Probability National University of Singapore Singapore, Singapore

*Address all correspondence to: Kody J. H. Law, E-mail: kody.law@manchester.ac.uk, URL: <https://sites.google.com/view/kodylaw/home>

Original Manuscript Submitted: ; Final Draft Received:

In this paper we consider sequential joint state and static parameter estimation given discrete time observations associated to a partially observed stochastic partial differential equation. It is assumed that one can only estimate the hidden state using a discretization of the model. In this context, it is known that the multi-index Monte Carlo (MIMC) method of [1] can be used to improve over direct Monte Carlo from the most precise discretization. However, in the context of interest, it cannot be directly applied, but rather must be used within another method such as sequential Monte Carlo (SMC). We show how one can use the MIMC method by renormalizing the standard identity and approximating the resulting identity using the SMC² method of [2], which is an exact method that can be used in this context. We prove that our approach can reduce the cost to obtain a given mean square error, relative to just using SMC² on the most precise discretization. We demonstrate this with some numerical examples.

KEY WORDS: Stochastic Partial Differential Equations; Multi-Index Monte Carlo, Sequential Monte Carlo

1. INTRODUCTION

We consider joint state and static parameter estimation for discrete time observations, associated to a partially observed stochastic partial differential equation (SPDE). Such models can be considered a form of hidden Markov model (HMM), and these have a significant number of practical applications; see e.g. [3] for instance. See Figure 1 for a schematic. The objective is to estimate the states $(X_0, X_1, \dots, X_n) \in \mathcal{X}^{n+1}$ and parameters $\theta \in \Theta$ given the data (y_0, \dots, y_n) , i.e. we want to find

$$\mathbb{P}(X_0 \in A_0, X_1 \in A_1, \dots, X_n \in A_n, \theta \in A_T | y_0, \dots, y_n)$$

recursively in time n .

In this article we focus upon the scenario where one will have to discretize the time and space element of the SPDE. In this scenario, one is faced with the problem of joint state and static parameter estimation for a HMM with a high-dimensional state; a problem which is notoriously challenging. The main issue is that for any fixed static parameter, one can seldom calculate the joint density (the smoother), given the data, of the hidden states over the observation times. Joint inference on the parameter is even more challenging, due to the dependence of all the hidden states on the parameter (see Figure 1). The state of the art method for consistently solving this problem for a fixed time n is the particle Markov chain Monte Carlo (PMCMC) method [4], which involves using sequential Monte Carlo (SMC)

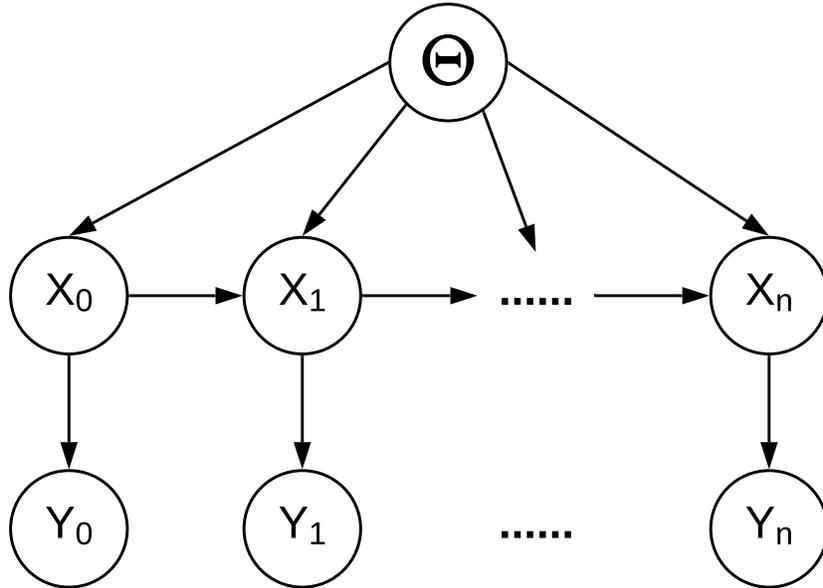


FIG. 1: A graphical model of the HMM studied in this paper.

within Markov chain Monte Carlo (MCMC). SMC methods can approximate the sequence of joint distributions on $X_{0:n}$ for θ fixed (smoothers). They consist of sampling N samples (also called particles) in parallel, sequentially in time. SMC methods involve the recursion of a mutation step, an importance sampling step, and a resampling step. They provide consistent (as $N \rightarrow \infty$) approximations of expectations w.r.t. the smoother. In many contexts, SMC is referred to as a particle filter, hence the name PMCMC. We remark, however, one also seeks to perform statistical inference sequentially in time, which adds yet another complication. There exists a methodology which can extend PMCMC to a dynamic context, called SMC² [2]. This method is an SMC algorithm which sequentially samples from a sequence of auxiliary target distributions (the targets of a certain PMCMC algorithm), which admit the distribution of interest as a marginal. The main reason why one would consider such a level of complexity, is that conventional SMC methodology which is designed for joint state and parameter inference associated to HMMs suffers from the so-called path degeneracy issue (see e.g. [5]), which renders it useless for long time intervals.

In the problem of interest, we are dealing with an expectation w.r.t. a probability measure which is defined on a high-dimensional continuous state-space as a result of the space and time discretization. It is assumed that the computational cost associated to performing any simulation-based numerical method will increase as the precision of the discretization is enhanced. In such a context, the multilevel Monte Carlo (MLMC) method is very effective for reducing the cost to achieve a given level of accuracy [6–8]. Following the success of the MLMC method, the work [1] revisited the MLMC identity through the lense of sparse grids for multi-dimensional discretizations. This more general method, which reduces to MLMC for one dimensional (discretizations) problems, is called multi-index Monte Carlo (MIMC). This method can be much more efficient than MLMC for higher dimensional problems with suitable regularity, for example providing canonical complexity $\mathcal{O}(\varepsilon^{-2})$ to achieve mean square error (MSE) $\mathcal{O}(\varepsilon^2)$. In this approach one rewrites the expectation of interest as a sum of difference of differences (DOD) w.r.t. independent refinements of the discretization levels of different dimensions, for example different spatial dimensions and time. If the number of dimensions which are discretized is d , then this DOD involves expectations w.r.t. at most 2^d different discretization levels. Then, under appropriate assumptions and given an efficient *coupling* of these 2^d distributions, the cost to achieve a prespecified MSE is reduced with respect to considering the Monte Carlo method using a single discretization level [1]. Indeed, under appropriate assumptions, relying strongly on the mixed regularity of the solution of the SPDE, and with an appropriate index set, the MIMC approach can achieve a substantial improvement in cost with respect to MLMC [1]. The coupling is essential to ensuring that the *variance* of the estimates of the DODs decay

appropriately w.r.t. the discretization indices. This allows one to use fewer samples in approximating higher index residuals, hence balancing the cost in an optimal way.

The sampling of a “good coupling” of the 2^d distributions is the main key to cost reduction in MIMC, as in MLMC. The problem of approximating expectations w.r.t. the joint distribution of interest for a single discretization level/index is already challenging, and coupling 2^d such distributions is naturally much more complex. A method for using MCMC within the MIMC framework was developed in [9], based upon an idea in [10] developed originally for PMCMC (see also [11]). The method involves constructing an approximate coupling of the 2^d targets and using this to approximate a re-normalized multi-index (MI) identity. See [12] for a pedagogical introduction to this general approximate coupling strategy. In the present context, the aim is to perform inference sequentially as data arrives, where the unobserved process is an SPDE. This is done using the SMC² method described above. It involves first extending the MLMC PMCMC method of [10] to the MIMC context, in order to accommodate an SPDE model. Second, this new MIMC PMCMC is deployed within an SMC² algorithm, which yields the MISMC² algorithm. The latter generalization of course yields MLSMC² as a byproduct, since ML is a particular case of MI. Under appropriate assumptions, we prove that our approach can reduce the cost required to obtain a given MSE, relative to just using SMC² on the most precise discretization, or even using the MLMC version. This is demonstrated with a numerical example.

Before proceeding, and to avoid confusion, we briefly digress on the distinction from existing applications of MLMC to SMC. For a more comprehensive discussion, the reader is referred to the recent review [12]. In particular, in [13] the authors employ an SMC algorithm over *levels* (i.e. level there is analogous to time in the present), and subsequently constructed MLMC estimators using importance sampling estimators of the level increments. In contrast, in [14] the authors couple pairs of SMC algorithms in the particle filtering context using a coupled resampling mechanism.

This article is structured as follows. In Section 2, we provide a high-level introduction into the underlying idea of the article. In Section 3 we describe the problem and how it may be solved, if numerical approximation were not required. In Section 4 we show how our approach can be numerically approximated. In Section 5 the theoretical result is given, with the proofs in the appendix. Numerical results are presented in Section 6 .

2. HIGH-LEVEL DISCUSSION OF THE APPROACH

The method presented in this article is quite complicated. To assist non-experts, we provide a high-level description of the basic idea with minimal notations, in which one dimension is discretized. First we explain the approximate coupling strategy used to enable MIMC estimation. Then we explain the PMCMC and SMC² methods.

2.1 Approximate Coupling

Consider a probability density on state-space X

$$p(x) := \frac{J(x)F(x)}{Z}$$

where J, F are two positive, real-valued functions, $\int_{\mathsf{X}} F(x)dx = 1$, and $Z = \int_{\mathsf{X}} J(x)F(x)dx$ (assumed to be finite). It is of interest to compute expectation of real-valued functions $\varphi : \mathsf{X} \rightarrow \mathbb{R}$ that are p -integrable:

$$\mathbb{E}_p[\varphi(X)] = \int_{\mathsf{X}} \varphi(x)p(x)dx.$$

Suppose that one only has access to a sequence $J_l(x)F_l(x)$, $l \in \{0, 1, \dots\}$ which are positive, real-valued functions such that

$$\lim_{l \rightarrow \infty} \mathbb{E}_{p_l}[\varphi(X)] = \mathbb{E}_p[\varphi(X)]$$

where $\mathbb{E}_{p_l}[\varphi(X)] = \int_{\mathsf{X}} \varphi(x)p_l(x)dx$, $p_l(x) = [J_l(x)F_l(x)]/Z_l$, $Z_l = \int_{\mathsf{X}} J_l(x)F_l(x)dx$.

Now, the MLMC identity

$$\mathbb{E}_{p_L}[\varphi(X)] = \sum_{l=1}^L \{\mathbb{E}_{p_l}[\varphi(X)] - \mathbb{E}_{p_{l-1}}[\varphi(X)]\} + \mathbb{E}_{p_0}[\varphi(X)]$$

can be very useful to reduce the computational effort in the Monte Carlo approximation of $\mathbb{E}_{p_L}[\varphi(X)]$, to achieve a given error (versus considering only $\mathbb{E}_{p_L}[\varphi(X)]$). The key to this method is the ability to construct a coupling $\check{p}_{l,l-1}$ of (p_l, p_{l-1}) for each $l \in \{1, 2, \dots\}$, i.e. a probability density function on $\mathbb{X} \times \mathbb{X}$ such that for every $(x, x') \in \mathbb{X} \times \mathbb{X}$

$$\begin{aligned} p_l(x) &= \int_{\mathbb{X}} \check{p}_{l,l-1}(x, x') dx' \\ p_{l-1}(x') &= \int_{\mathbb{X}} \check{p}_{l,l-1}(x, x') dx. \end{aligned}$$

Then, one has

$$\mathbb{E}_{p_l}[\varphi(X)] - \mathbb{E}_{p_{l-1}}[\varphi(X)] = \int_{\mathbb{X} \times \mathbb{X}} \varphi(x) \check{p}_{l,l-1}(x, x') d(x, x') - \int_{\mathbb{X} \times \mathbb{X}} \varphi(x') \check{p}_{l,l-1}(x, x') d(x, x'). \quad (1)$$

If the coupling is sufficiently good, so that for instance

$$\int_{\mathbb{X} \times \mathbb{X}} (\varphi(x) - \varphi(x'))^2 \check{p}_{l,l-1}(x, x') d(x, x') \leq h(l),$$

where $\lim_{l \rightarrow \infty} h(l) = 0$, h is a positive, real-valued, monotonically decreasing function on $\{0, 1, \dots\}$, then the aforementioned benefits are possible; see e.g. [6–8]. The Monte Carlo method would rely on exact sampling from the distribution associated to the coupling $\check{p}_{l,l-1}$.

Let $l \geq 1$ be fixed. In many practical problems of interest, such as the one considered in this article, deriving a suitable coupling $\check{p}_{l,l-1}$ which is amenable to known simulation methodology can be very challenging. The basic idea used in this paper and as adopted in [10] is as follows. Suppose one can find a coupling $\check{F}_{l,l-1}$ of (F_l, F_{l-1}) such that

$$\begin{aligned} F_l(x) &= \int_{\mathbb{X}} \check{F}_{l,l-1}(x, x') dx' \\ F_{l-1}(x') &= \int_{\mathbb{X}} \check{F}_{l,l-1}(x, x') dx, \end{aligned} \quad (2)$$

and

$$\int_{\mathbb{X} \times \mathbb{X}} (\varphi(x) - \varphi(x'))^2 \check{F}_{l,l-1}(x, x') d(x, x') \leq h(l).$$

Now set

$$\check{p}_{l,l-1}(x, x') = \frac{\max\{J_l(x), J_{l-1}(x')\} \check{F}_{l,l-1}(x, x')}{\check{Z}_{l,l-1}},$$

with $\check{Z}_{l,l-1} = \int_{\mathbb{X} \times \mathbb{X}} \max\{J_l(x), J_{l-1}(x')\} \check{F}_{l,l-1}(x, x') d(x, x')$ (assumed to be finite). Note that

$$\begin{aligned} \mathbb{E}_{p_l}[\varphi(X)] &= \int_{\mathbb{X}} \varphi(x) p_l(x) dx \\ &= \frac{1}{\check{Z}_l} \int_{\mathbb{X} \times \mathbb{X}} \varphi(x) J_l(x) \check{F}_{l,l-1}(x, x') d(x, x') \\ &= \frac{\check{Z}_{l,l-1}}{\check{Z}_l} \int_{\mathbb{X} \times \mathbb{X}} \varphi(x) \frac{J_l(x)}{\max\{J_l(x), J_{l-1}(x')\}} \check{p}_{l,l-1}(x, x') d(x, x') \\ &= \mathbb{E}_{\check{p}_{l,l-1}} \left[\varphi(X) \frac{J_l(X)}{\max\{J_l(X), J_{l-1}(X')\}} \right] / \mathbb{E}_{\check{p}_{l,l-1}} \left[\frac{J_l(X)}{\max\{J_l(X), J_{l-1}(X')\}} \right] \end{aligned} \quad (3)$$

where the last line uses that $\mathbb{E}_{\tilde{p}_{l,t-1}} \left[\frac{J_l(X)}{\max\{J_l(X), J_{l-1}(X')\}} \right] = Z_l / \tilde{Z}_{l,t-1}$. Thus

$$\begin{aligned} \mathbb{E}_{p_l}[\varphi(X)] - \mathbb{E}_{p_{l-1}}[\varphi(X)] &= \mathbb{E}_{\tilde{p}_{l,t-1}} \left[\varphi(X) \frac{J_l(X)}{\max\{J_l(X), J_{l-1}(X')\}} \right] / \mathbb{E}_{\tilde{p}_{l,t-1}} \left[\frac{J_l(X)}{\max\{J_l(X), J_{l-1}(X')\}} \right] - \\ &\quad \mathbb{E}_{\tilde{p}_{l,t-1}} \left[\varphi(X') \frac{J_{l-1}(X')}{\max\{J_l(X), J_{l-1}(X')\}} \right] / \mathbb{E}_{\tilde{p}_{l,t-1}} \left[\frac{J_{l-1}(X')}{\max\{J_l(X), J_{l-1}(X')\}} \right]. \end{aligned}$$

The main interest of this identity is the fact that one may be able to construct a coupling like $\tilde{F}_{l,t-1}$ and very efficient sampling methods for $\tilde{p}_{l,t-1}$, whereas this may not be the case for $\tilde{p}_{l,t-1}$. It is then possible (e.g. [10]) that the benefits of the MLMC method can be achieved, even though one does not know how to sample from a good coupling $\tilde{p}_{l,t-1}$.

2.2 Monte Carlo Methods

First a simplified description of PMCMC is given. Suppose one aims to estimate expectations with respect to the joint state and parameter smoothing distribution associated to an HMM, as described above, i.e. the aim is to approximate

$$p(x_{1:n}, \theta | y_{1:n}) \propto p(x_{1:n} | y_{1:n}, \theta) p(y_{1:n} | \theta) p(\theta).$$

The first factor on the right-hand side is the smoothing distribution for a fixed parameters, which can be well-approximated consistently with a particle filter. The second term is the marginal likelihood, for which an unbiased and non-negative estimator is available via the particle filter. It seems reasonable then to consider approximating the target distribution using a particle filter. The PMCMC method leverages this intuition by using a non-negative unbiased estimate of the un-normalized joint distribution above derived from the particle filter within an MCMC method. Let $\mathbb{P}_n(v_{1:n} | \theta)$ denote the distribution of *all auxiliary variables* $v_{1:n} = (x_{1:n}^{1:N}, a_{1:n}^{1:N})$ of an N -particle filter targeting the smoothing distribution $p(x_{1:n} | y_{1:n}, \theta)$ (the variables $a_{1:n}^{1:N} \in [1, \dots, N]^{n \times N}$ denote the N resampled indices at times $i = 1, \dots, n$). Let $\hat{p}(y_{1:n} | \theta)$ denote the particle filter estimate of the marginal likelihood. PMCMC is an MCMC method with the target distribution $\mathbb{P}_n(v_{1:n} | \theta) \hat{p}(y_{1:n} | \theta) p(\theta)$.

Now, we describe a simplified version of the SMC² algorithm. Let $\mathbb{Q}_n(v_{n+1} | v_{1:n}, \theta)$ denote the 1-step transition kernel of the particle filter, so that $\mathbb{P}_n(v_{1:n+1} | \theta) = \mathbb{Q}_n(v_{n+1} | v_{1:n}, \theta) \mathbb{P}_n(v_{1:n} | \theta)$. Suppose K_n is an MCMC kernel targeting $\Pi_n(v_{1:n}, \theta) \propto \mathbb{Q}_n(v_{1:n}, \theta)$, where one can construct estimates of $p(x_{1:n}, \theta | y_{1:n})$ from an appropriate marginal of Π_n . Define

$$M_n(v_{1:n}, dv'_{1:n+1}) = K_n(v_{1:n}, dv'_{1:n}) \otimes \mathbb{Q}_n(v'_{1:n}, v'_{n+1}) dv'_{n+1}. \quad (4)$$

It is now possible to run an SMC sampler to sequentially target Π_n . For $i = 1, \dots, N$, one draws $\hat{v}_1^i := v^i \sim \Pi_0$, and then iterates for $n \geq 1$

- Simulate $v_{n+1}^i \sim M_n(\hat{v}_{1:n}^i, \cdot)$;
- Resample $\hat{v}_{1:n+1}^i = v_{1:n+1}^j$, with probability proportional to

$$\frac{\mathbb{Q}_{n+1}(v_{1:n+1}^j)}{\mathbb{Q}_n(v_{1:n}^j) \mathbb{Q}_n(v_{1:n}^j, v_{n+1}^j)}.$$

Note that one does not need to be able evaluate either \mathbb{Q}_n or \mathbb{Q}_n , but only simulate from them and evaluate the ratio above. A very similar method, details of which will be given in Sec. 4, is referred to as SMC²: the first (inner) SMC appears in the PMCMC kernel K_n , for each n , and the second (outer) SMC appears in the SMC sampler on the extended state-space.

3. RIGOROUS PROBLEM FORMULATION

3.1 Model

Let (Y, \mathcal{Y}) and (X, \mathcal{X}) be measurable spaces. We consider a pair of stochastic processes indexed by a time parameter. We are given a sequence of observations y_0, y_1, \dots which are realizations of a discrete-time process $\{Y_n\}_{n \in \mathbb{N}_0}$, $Y_n \in Y$, where the time between observations is one unit. These observations are associated with a continuous-time (Markov) stochastic process $\{X_t\}_{t \geq 0}$, with $X_t \in X$. The process would typically arise from the finite-time evolution of an SPDE, although we do not make this constraint at this time.

We now present the stochastic model which describes the probabilistic structure of the processes $\{Y_n\}_{n \in \mathbb{N}_0}$ and $\{X_t\}_{t \geq 0}$. In our model, $\theta \in \Theta \subseteq \mathbb{R}^k$ is a static parameter associated to the model. We will define the aforementioned structure, conditional upon θ and then define a prior probability distribution on this static parameter. Let $X_{0:n} = (X_0, \dots, X_n)$ correspond to a discrete-time skeleton of $\{X_t\}_{t \geq 0}$ on the grid $0 : n$. We are interested in the posterior probability distribution of $(X_{0:n}, \theta)$ conditional on observed data y_0, \dots, y_n , sequentially over discrete unit times (n). It is supposed that for any $n \geq 0$, $A \in \mathcal{Y}$

$$\mathbb{P}(Y_n \in A | y_{0:n-1}, \{x_t\}_{t \in [0, \dots, n]}, \theta) = \int_A g_\theta(x_n, y) dy$$

where dy is a σ -finite measure on (Y, \mathcal{Y}) and for each $(\theta, x) \in \Theta \times X$, $g_\theta(x, \cdot) : Y \rightarrow \mathbb{R}_+$ is a probability density on Y . For each $\theta \in \Theta$, $f_\theta : X^2 \rightarrow \mathbb{R}_+$, (resp. $\mu_\theta : X \rightarrow \mathbb{R}_+$) are the transition density over unit time (resp. initial density) of $\{X_t\}_{t \geq 0}$ (resp. X_0) w.r.t. a dominating σ -finite measure on (X, \mathcal{X}) . Note that for every $(\theta, x) \in \Theta \times X$, $f_\theta(x, \cdot) : X \rightarrow \mathbb{R}_+$ is a probability density, and for any $n \geq 1$

$$\mathbb{P}(X_n \in A | \{x_r\}_{0 \leq r \leq n-1}, \theta) = \int_A f_\theta(x_{n-1}, x) dx.$$

Let ν be a probability density w.r.t. Lebesgue measure (written $d\theta$) on $(\Theta, \mathcal{B}(\Theta))$ with $\mathcal{B}(\Theta)$ the Borel sets.

For $n \geq 0$, the posterior probability density on $X^{n+1} \times \Theta$ that is induced by this construction is given by

$$\pi_n(x_{0:n}, \theta) \propto \nu(\theta) \mu_\theta(x_0) g_\theta(x_0, y_0) \prod_{p=1}^n f_\theta(x_{p-1}, x_p) g_\theta(x_p, y_p). \quad (5)$$

In other words for $A \in \mathcal{V}^{n+1} \mathcal{X} \vee \mathcal{B}(\Theta)$

$$\mathbb{P}((X_{0:n}, \theta) \in A | y_{0:n}) = \int_A \pi_n(x_{0:n}, \theta) d(x_{0:n}, \theta).$$

Henceforth, we will suppress the dependence on $y_{0:n}$ throughout the article. Let $\varphi : X^{n+1} \times \Theta \rightarrow \mathbb{R}$ be integrable w.r.t. π_n . Our objective is to compute, recursively in n

$$\mathbb{E}_{\pi_n}[\varphi(X_{0:n}, \theta)] := \int_{X^{n+1} \times \Theta} \varphi(x_{0:n}, \theta) \pi_n(x_{0:n}, \theta) d(x_{0:n}, \theta)$$

where we use the notation \mathbb{E}_π to denote expectations w.r.t. a probability density/measure π . The role of the function φ is as a summary or quantity of interest, relating to the random variables $X_{0:n}, \theta$. That is, our objective is to compute expectations, such as moments, w.r.t. the posterior, with density defined in (5). We note that one often must use Monte Carlo methods to approximate the sequence of expectations.

Before concluding this section we mention a canonical statistical model, the conditionally Gaussian model.

Example 1. Let $N(m, C)$ denote a (possibly infinite-dimensional) Gaussian random variable with mean m and covariance operator C , and let $\phi(\cdot; m, C)$ denote its density (with respect to some dominating measure, which may be

taken as Lebesgue in finite dimensions). Assume $X_0 \sim N(m_0, \Sigma_0)$. For each $\theta \in \Theta$, let $\Psi_\theta : X \rightarrow X$ and $h_\theta : X \rightarrow Y$ be continuous and let $\Sigma_\theta, \Gamma_\theta$ be symmetric positive definite operators. An example of a model is, for $n \geq 0$

$$\begin{aligned} X_{n+1}|X_n &\sim N(\Psi_\theta(X_n), \Sigma_\theta), \\ Y_n|X_n &\sim N(h_\theta(X_n), \Gamma_\theta). \end{aligned} \quad (6)$$

This model is ubiquitous in the data assimilation literature [15]. Once ν is specified, the model (5) is given by $\mu_\theta(x_0) = \phi(x_0; m_0, \Sigma_0)$, $f_\theta(x, x') = \phi(x'; \Psi_\theta(x), \Sigma_\theta)$, and $g_\theta(x, y) = \phi(y; h_\theta(x), \Gamma_\theta)$.

This model fits into the context of this paper when X is infinite dimensional, e.g. a Hilbert space, and Ψ_θ is the solution of a PDE parametrized by θ .

3.2 Discretized Model

The exposition here closely follows that developed in [9]. Here we explicitly assume one must work with a discretized version of the model, that is, there does not (currently) exist an unbiased and non-negative approximation of $\pi_n(x_{0:n}, \theta)$. We remark that if the latter approximations are available, then the strategy to be outlined is not required.

Set $\alpha \in \mathbb{N}_0^d$, which will refer to a collection of indices which will denote the level of discretization of our model in each of d dimensions. That is, as the components of α increase, so does the accuracy of the approximation. Precise examples are given in Section 6. More explicitly, for any fixed $\alpha \in \mathbb{N}_0^d$, let $(X_\alpha, \mathcal{X}_\alpha)$ and $(Y_\alpha, \mathcal{Y}_\alpha)$ be measurable spaces such that for every $n \geq 0$ one can obtain a biased approximation $X_\alpha \in X_\alpha \subseteq X$ of X_n , and $Y_\alpha \in Y_\alpha \subseteq Y$ of Y_n . That is, one can define the probability density for $n \geq 0$, on $X_\alpha^{n+1} \times \Theta$:

$$\pi_{n,\alpha}(x_{0:n}, \theta) \propto \nu(\theta) \mu_{\theta,\alpha}(x_0) g_{\theta,\alpha}(x_0, y_0) \prod_{p=1}^n f_{\theta,\alpha}(x_{p-1}, x_p) g_{\theta,\alpha}(x_p, y_p) \quad (7)$$

where $y_n \in Y_\alpha$ for each $n \geq 0$. Here for every $(\alpha, \theta) \in \mathbb{N}_0^d \times \Theta$

- For all $x \in X_\alpha$, $g_{\theta,\alpha}(x, \cdot)$ is a probability density on Y_α ;
- For all $x \in X_\alpha$, $f_{\theta,\alpha}(x, \cdot)$ is a probability density on X_α ;
- $\mu_{\theta,\alpha}$ is a probability density on X_α .

Consider $\varphi : \mathbb{N}_0^d \times X^{n+1} \times \Theta \rightarrow \mathbb{R}$, where for any $(x_{0:n}, \theta) \in X^{n+1} \times \Theta$

$$\lim_{\min_{1 \leq i \leq d} \alpha_i \rightarrow +\infty} \varphi_\alpha(x_{0:n}, \theta) = \varphi(x_{0:n}, \theta).$$

It is assumed that we have

$$\begin{aligned} \mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)] &\neq \mathbb{E}_{\pi_n}[\varphi(X_{0:n}, \theta)]. \\ \lim_{\min_{1 \leq i \leq d} \alpha_i \rightarrow +\infty} |\mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)] - \mathbb{E}_{\pi_n}[\varphi(X_{0:n}, \theta)]| &= 0. \end{aligned} \quad (8)$$

As remarked in the introduction, the computational cost associated with X_α, Y_α (sampling, or evaluating the densities $g_{\theta,\alpha}, f_{\theta,\alpha}, \mu_{\theta,\alpha}$) increases as any index of α increases. Our objective is now to compute $\mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)]$ recursively for each $n \geq 0$.

3.3 Multi-Index Methods

The approach to be described, provides an approach to approximating $\mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)]$ for any fixed $\alpha \in \mathbb{N}_0^d$. Define the difference operator $\Delta_i, i \in \{1, \dots, d\}$ as

$$\Delta_i \mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)] := \begin{cases} \mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)] - \mathbb{E}_{\pi_{n,\alpha-e_i}}[\varphi_{\alpha-e_i}(X_{0:n}, \theta)] & \text{if } \alpha_i > 0 \\ \mathbb{E}_{\pi_{n,\alpha}}[\varphi_\alpha(X_{0:n}, \theta)] & \text{otherwise} \end{cases}$$

where e_i are the canonical vectors on \mathbb{R}^d . Set

$$\Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)] := (\Delta_1 \circ \Delta_2 \circ \cdots \circ \Delta_d) \left(\mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)] \right) \quad (9)$$

where $(\Delta_1 \circ \Delta_2 \circ \cdots \circ \Delta_d)$ denotes the composition of $\Delta_1, \dots, \Delta_d$. Note that the order of applying the operators Δ_i in Δ does not matter.

We now consider the identity

$$\mathbb{E}_{\pi_n} [\varphi(X_{0:n}, \theta)] = \sum_{\alpha \in \mathbb{N}_0^d} \Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)].$$

The work [1] proposes to leverage this identity by constructing a biased estimator of $\mathbb{E}_{\pi_n} [\varphi(X_{0:n}, \theta)]$ for some finite multi-index set $\mathcal{I} \subset \mathbb{N}_0^d$ as follows

$$\mathbb{E}_{\pi_n, \mathcal{I}} [\varphi(X_{0:n}, \theta)] := \sum_{\alpha \in \mathcal{I}} \Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)]. \quad (10)$$

Such an approximation strategy is inspired by work in the sparse grids literature [16]. This estimator can (in principle) be approximated by a Monte Carlo method. This can be achieved by (if possible) sampling from a coupling of the (at most) 2^d different probability measures for a given $\alpha \in \mathcal{I}$ (see [1] for details). It is counterintuitive at first to construct a single estimator from a sum of so many other estimators, but in fact if the coupling is strong, and under appropriate assumptions, then this can be substantially more efficient than a single term estimator, in the sense that smaller MSE can be achieved for the same cost. It is remarked that sampling from such a coupling is very challenging, especially in the context of the model considered in Section 3.2. The residual error is given by

$$\mathbb{E}_{\pi_n} [\varphi(X_{0:n}, \theta)] - \mathbb{E}_{\pi_n, \mathcal{I}} [\varphi(X_{0:n}, \theta)] = \sum_{\alpha \notin \mathcal{I}} \Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)].$$

It is shown in [1,10] that under appropriate assumptions on the convergence of estimates of the individual terms in (9) one can gain significant ‘speed-up’ relative to single term or even single index MLMC methods. The key point we emphasize here is that *all results of [1], pertaining to all different index sets \mathcal{I}* , rely solely on the convergence properties of estimates of the individual terms (9). Since there are significant other difficulties to deal with in the present work, the finer properties of estimators with various different index sets \mathcal{I} will not be considered here, although we note this is a crucial consideration in practice. Our objective here will rather be to establish a general proof of principle method which provides convergence of estimates of the individual terms in (9) under suitable assumptions. The results will be illustrated in Section 6.

3.4 Renormalized Multi-Index Identity

The following idea builds upon the approaches in [9] and [10]. Consider (10) and in particular consider a single given summand (9) for $\alpha \in \mathcal{I}$, with n fixed. This summand is itself a linear combination of expectations with respect to $1 < k_\alpha \leq 2^d$ probability measures. These k_α probability measures induce k'_α differences in (10); if $k_\alpha = 2^d$, then $k'_\alpha = 2^{d-1}$. We remark that in the case that $k_\alpha = 1$, one does not need to consider how to construct a coupling for an MLMC method as the summand (9) is only an expectation w.r.t. a single probability measure.

For simplicity of notation we will denote the k_α multi-indices by $\alpha(1), \dots, \alpha(k_\alpha)$, where for $i \in \{1, \dots, k_\alpha\}$, $\alpha(i) \in \mathcal{I}$. Let $\alpha(i)_j$ denote the j^{th} -element of $\alpha(i)$. The convention of the (non-unique) labelling is such that, $\sum_{j=1}^d [\alpha(2i) - \alpha(2i-1)]_j = 1$ for each $i \in \{1, \dots, k'_\alpha\}$, $\alpha(k_\alpha) = \alpha$ and $\alpha(1) = (\max\{\alpha_1 - 1, 0\}, \dots, \max\{\alpha_d - 1, 0\})$. This labelling will provide a convenient way to write $\Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)]$ below.

Example 2. Suppose $d = 3$ and $\alpha = (2, 2, 2)$, so $k_\alpha = 8$, $k'_\alpha = 4$. Then a labelling which satisfies the above constraints is

$$\begin{aligned} \alpha(1) &= (1, 1, 1), \alpha(2) = (2, 1, 1), \alpha(3) = (1, 1, 2), \alpha(4) = (2, 1, 2), \\ \alpha(5) &= (1, 2, 1), \alpha(6) = (2, 2, 1), \alpha(7) = (1, 2, 2), \alpha(8) = (2, 2, 2). \end{aligned}$$

Recall the form of the target (7) and the multi-increment summand to be estimated (9). We suppose that there exists a coupling of the discretized dynamics. That is, there exists a Markov density $\check{f}_{\theta, \alpha(1:k_\alpha)}(x(1:k_\alpha), x'(1:k_\alpha))$ such that for any $x(1:k_\alpha) = (x(1), \dots, x(k_\alpha)) \in \bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}$ and any $i \in \{1, \dots, k_\alpha\}$, $A_i \in \mathcal{X}_{\alpha(i)}$, we have:

$$\int_{\bigotimes_{j=1}^{i-1} \mathcal{X}_{\alpha(j)} \times A_i \times \bigotimes_{j=i+1}^{k_\alpha} \mathcal{X}_{\alpha(j)}} \check{f}_{\theta, \alpha(1:k_\alpha)}(x(1:k_\alpha), x'(1:k_\alpha)) dx'(1:k_\alpha) = \int_{A_i} f_{\theta, \alpha(i)}(x(i), x'(i)) dx'(i). \quad (11)$$

In other words, for a given α , coupled Markov dynamics are performed on the hierarchy of k_α meshes in such a way that the marginal of the coupled dynamics with respect to any of the k_α meshes $\alpha(i)$ corresponds to the exact Markov dynamics on mesh $\alpha(i)$. Importantly, we do not assume that we can evaluate this Markov density. We will only require being able to simulate from it.

Similarly, we suppose that there exists a probability density $\check{\mu}_\theta$ on $\bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}$ such that for any $i \in \{1, \dots, k_\alpha\}$, $A_i \in \mathcal{X}_{\alpha(i)}$, we have:

$$\int_{\bigotimes_{j=1}^{i-1} \mathcal{X}_{\alpha(j)} \times A_i \times \bigotimes_{j=i+1}^{k_\alpha} \mathcal{X}_{\alpha(j)}} \check{\mu}_{\theta, \alpha(1:k_\alpha)}(x(1:k_\alpha)) dx(1:k_\alpha) = \int_{A_i} \mu_{\theta, \alpha(i)}(x(i)) dx(i). \quad (12)$$

We remark that one can find scenarios for which this is true. We give an example below and another later in Section 6.

Example 3. As a concrete example, consider the setting of (6) in Example 1, where Ψ_θ is the forward solution of a PDE over a unit time interval, for example Navier-Stokes equation as in [17]. Suppose $\mu_\theta = N(0, S^a)$, where $a > 0$ and S is the (possibly weak) solution operator $S : f \mapsto u$ of an elliptic PDE on a cubic domain $\Omega = [0, 1]^d$, with convex boundary $\partial\Omega$:

$$\left(- \sum_{i=1}^d \partial^2 u / \partial x_i^2 \right) (x) = f(x), \quad x \in \Omega, \\ u(x) = 0, \quad x \in \partial\Omega.$$

Suppose we have eigenfunctions $\{v_k\}_{k \in \mathbb{Z}_+^d}$ such that $Sv_k = \lambda_k v_k$ for $k \in \mathbb{Z}_+^d$, and then we have a spectral (Karhunen-Loève) expansion of $X_0 \sim \mu_\theta$ as $X_0 = \sum_{k \in \mathbb{Z}_+^d} \lambda_k^{a/2} \xi_k v_k$, where $\xi_k \sim N(0, 1)$ i.i.d. (see e.g. [18]). Note the eigenfunctions can be decomposed as products of eigenfunctions $\{\psi_l\}_{l \in \mathbb{Z}_+}$ of the $d = 1$ problem $(\partial^2 \psi_l / \partial x_1^2)(x_1) = \lambda_l \psi_l(x_1)$, i.e. $v_k(x_1, \dots, x_d) = \prod_{i=1}^d \psi_{k_i}(x_i)$. Furthermore, suppose that $X_{0, \alpha} = \sum_{k \in \mathcal{T}_\alpha} \lambda_k^{a/2} \xi_k v_k$, where $\mathcal{T}_\alpha = \{k \in \mathbb{Z}_+^d; 1 \leq k_1 \leq 2^{\alpha_1}, \dots, 1 \leq k_d \leq 2^{\alpha_d}\}$. One can simulate

$$X_{0, \alpha(k_\alpha)} = \sum_{k \in \mathcal{T}_{\alpha(k_\alpha)}} \lambda_k^{a/2} \xi_k v_k \sim \mu_{\theta, \alpha(k_\alpha)},$$

and then coarsen this realization appropriately such that $X_{0, \alpha(i)} = \sum_{k \in \mathcal{T}_{\alpha(i)}} \lambda_k^{a/2} \xi_k v_k$ is a realization of $\mu_{\theta, \alpha(i)}$, for each of the other targets $i < k_\alpha$. Note that $\mathcal{T}_{\alpha(i)} \subset \mathcal{T}_{\alpha(k_\alpha)}$ so this just consists in using the same $\{\xi_k\}_{k \in \mathcal{T}_{\alpha(k_\alpha)}}$ and setting $\xi_k = 0$ for all $k \in \mathcal{T}_{\alpha(k_\alpha)} \setminus \mathcal{T}_{\alpha(i)}$. Hence we have a strong coupling which satisfies (12).

For simulating the forward kernel $\Psi_{\theta, \alpha}$, assume that we have a Galerkin spectral solver [19] for any given spectral truncation level α and a given time-step size (if time-discretization is considered in the approximation, its discretization level will be given by index α_{d+1}). Assume $\Sigma_\theta = S^b$, for $b > 0$, with S defined as above, and

discretizations defined similarly. One can then simulate a single realization $\chi_{\alpha(k_\alpha)} \sim N(0, \Sigma_{\theta, \alpha(k_\alpha)})$ just as for the initial condition and coarsen this to $\chi_{\alpha(i)}$ for driving each of the other dynamics with $i < k_\alpha$. Hence we have a strong coupling which satisfies (11).

Based on (11) and (12) we have an exact k_α -fold coupling of

$$\nu(\theta) \check{\mu}_{\theta, \alpha(1:k_\alpha)}(x_0(1:k_\alpha)) \prod_{p=1}^n \check{f}_{\theta, \alpha(1:k_\alpha)}(x_{p-1}(1:k_\alpha), x_p(1:k_\alpha)),$$

which takes the role of the 2-fold coupling $\check{F}_{l,l-1}$ given in (2). Let $\check{g} : \mathbb{N}_0^{d \times k_\alpha} \times \bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)} \times \Theta \times \mathcal{Y} \rightarrow (0, \infty)$ be arbitrary for the moment. We consider the following probability density on the space $(\bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)})^{n+1} \times \Theta$

$$\begin{aligned} \xi_{n, \alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta) &\propto \nu(\theta) \check{\mu}_{\theta, \alpha(1:k_\alpha)}(x_0(1:k_\alpha)) \prod_{p=1}^n \check{f}_{\theta, \alpha(1:k_\alpha)}(x_{p-1}(1:k_\alpha), x_p(1:k_\alpha)) \times \\ &\prod_{p=0}^n \check{g}_{\theta, \alpha(1:k_\alpha)}(x_p(1:k_\alpha), y_p). \end{aligned} \quad (13)$$

In the works [9,10] the following choice is made, for $p = 0, \dots, n$,

$$\check{g}_{\theta, \alpha(1:k_\alpha)}(x_p(1:k_\alpha), y_p) = \max\{g_{\theta, \alpha(1)}(x_p(1), y_p), \dots, g_{\theta, \alpha(k_\alpha)}(x_p(k_\alpha), y_p)\}, \quad (14)$$

and this is the choice used in this article. The motivation will become apparent shortly. It is noted that other choices are possible (see e.g. [11] for further discussion and some alternatives).

The expression (9) will be approximated using samples distributed according to (13). We start by considering a single expectation in (9). Note that (11) and (12), and the form of $\xi_{n, \alpha(1:k_\alpha)}$, immediately imply that for any $\alpha(i)$, $i \in \{1, \dots, k_\alpha\}$

$$\mathbb{E}_{\pi_{n, \alpha(i)}}[\varphi_{\alpha(i)}(X_{0:n}, \theta)] = \frac{\mathbb{E}_{\xi_{n, \alpha(1:k_\alpha)}}[\varphi_{\alpha(i)}(X_{0:n}(i), \theta) \prod_{p=0}^n \frac{g_{\theta, \alpha(i)}(X_p(i), y_p)}{\check{g}_{\theta, \alpha(1:k_\alpha)}(X_p(1:k_\alpha), y_p)}]}{\mathbb{E}_{\xi_{n, \alpha(1:k_\alpha)}}[\prod_{p=0}^n \frac{g_{\theta, \alpha(i)}(X_p(i), y_p)}{\check{g}_{\theta, \alpha(1:k_\alpha)}(X_p(1:k_\alpha), y_p)}]}. \quad (15)$$

This is recognizable as the analogue of (3). The following notation will be used for the weights, for any $\alpha(i)$, $i \in \{1, \dots, k_\alpha\}$

$$H_{i, n, \alpha, \theta}(x_{0:n}(1:k_\alpha)) := \prod_{p=0}^n \frac{g_{\theta, \alpha(i)}(x_p(i), y_p)}{\check{g}_{\theta, \alpha(1:k_\alpha)}(x_p(1:k_\alpha), y_p)}.$$

Note that the form of (14) ensures the weights are bounded (by 1), which is desirable for stability of the algorithm. By combining (15) with $H_{i, n, \alpha, \theta}(x_{0:n}(1:k_\alpha))$ and recalling the definitions of Δ and Δ_i given in and above (9), one can then deduce that

$$\begin{aligned} \Delta \mathbb{E}_{\pi_{n, \alpha}}[\varphi_\alpha(X_{0:n}, \theta)] &= \sum_{i=1}^{k'_\alpha} \tau_{i, \alpha} \left(\frac{\mathbb{E}_{\xi_{n, \alpha(1:k_\alpha)}}[\varphi_{\alpha(2i)}(X_{0:n}(2i), \theta) H_{2i, n, \alpha, \theta}(X_{0:n}(1:k_\alpha))]}{\mathbb{E}_{\xi_{n, \alpha(1:k_\alpha)}}[H_{2i, n, \alpha, \theta}(X_{0:n}(1:k_\alpha))]} \right. \\ &\quad \left. \frac{\mathbb{E}_{\xi_{n, \alpha(1:k_\alpha)}}[\varphi_{\alpha(2i-1)}(X_{0:n}(2i-1), \theta) H_{2i-1, n, \alpha, \theta}(X_{0:n}(1:k_\alpha))]}{\mathbb{E}_{\xi_{n, \alpha(1:k_\alpha)}}[H_{2i-1, n, \alpha, \theta}(X_{0:n}(1:k_\alpha))]} \right) \end{aligned} \quad (16)$$

where $|\alpha| = \sum_{j=1}^d \alpha_j$ and $\tau_{i, \alpha} = (-1)^{|\alpha(k_\alpha) - \alpha(2i)|}$. This is the multi-index analogue of (1), and has been used before in [9]. This approach is called approximate coupling; see [12] for further discussion.

Our strategy for approximating $\mathbb{E}_{\pi_{n, \mathcal{I}}}[\varphi(X_{0:n}, \theta)]$ is then the following. Noting (10) and (16), we will approximate each summand in (10), by approximating the r.h.s. of (16). This will be achieved as follows. Independently for each $\alpha \in \mathcal{I}$ (with $k_\alpha > 1$), and serially for each n , we will sample (approximately) from $\xi_{n, \alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)$, and compute a Monte Carlo estimate of $\Delta \mathbb{E}_{\pi_{n, \alpha}}[\varphi_\alpha(X_{0:n}, \theta)]$. As noted above, in the case $k_\alpha = 1$, one can simply sample from $\pi_{n, \alpha}$ as no coupling is required.

4. SIMULATION STRATEGY

The purpose of this Section is to describe a method to approximate $\mathbb{E}_{\pi_{n,\mathcal{I}}}[\varphi(X_{0:n}, \theta)]$. This will be achieved by using the SMC² method to approximate expectations w.r.t. $\xi_{n,\alpha(1:k_\alpha)}$ for each $\alpha \in \mathcal{I}$ (with $k_\alpha > 1$), recursively in n . If $k_\alpha = 1$, then one can simply use the standard SMC² method considering $\pi_{n,\alpha}$. The SMC² approach uses the particle MCMC method [4], which in turn relies upon particle filters. Therefore, to develop our approach, we first provide a review of particle filters, followed by particle MCMC in the present context.

4.1 Particle Filter

In this section, we focus upon the approximation of the density of the state conditional on fixed parameter θ , given by $\xi_{n,\alpha(1:k_\alpha),\theta}(x_{0:n}(1:k_\alpha)) \propto \xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)$. It is natural to adopt an SMC approach, in which we sequentially perform importance sampling and then resampling. This procedure is the standard particle filter which is given in Algorithm 1. The number of particles N will be fixed once and for all, but it is noted that this is an important consideration for the efficiency of the proposed method.

The joint density of all the variables sampled in Algorithm 1, up to time n , is written

$$\psi_{\alpha,\theta}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}) \quad (17)$$

where, for $0 \leq k \leq n$, $x_k^{1:N}(1:k_\alpha) = (x_k^1(1:k_\alpha), \dots, x_k^N(1:k_\alpha)) \in (\otimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)})^N$, $a_k^{1:N} = (a_k^1, \dots, a_k^N) \in \{1, \dots, N\}^N$, $x_{0:n}^{1:N}(1:k_\alpha) = (x_0^{1:N}(1:k_\alpha), \dots, x_n^{1:N}(1:k_\alpha)) \in (\otimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)})^{(n+1)N}$, $a_{0:n-1}^{1:N} = (a_0^{1:N}, \dots, a_{n-1}^{1:N}) \in \{1, \dots, N\}^{nN}$. In particular, a_{n-1}^i is the index at time $n-1$ of the resampled particle which has the index i at time n . For each $i \in \{1, \dots, N\}$ define the ancestral lineage indices as

$$b_n^i = i \text{ and } b_k^i = a_k^{b_{k+1}^i}, \quad k \in \{0, \dots, n-1\}. \quad (18)$$

For each $i = 1, \dots, N$, define

$$\bar{x}_{0:n}^i(1:k_\alpha) := x_{0:n}^{b_{0:n}^i}(1:k_\alpha). \quad (19)$$

Algorithm 1: The Particle Filter

- **Initialize:** Set $p = 0$, for $i \in \{1, \dots, N\}$ sample $x_0^i(1:k_\alpha)$ from $\check{\mu}_{\theta,\alpha(1:k_\alpha)}$ and evaluate the weight

$$w_{p,\alpha,\theta}^i = \left(\check{g}_{\theta,\alpha(1:k_\alpha)}(x_0^i(1:k_\alpha), y_0) \right) \left(\sum_{j=1}^N \check{g}_{\theta,\alpha(1:k_\alpha)}(x_0^j(1:k_\alpha), y_0) \right)^{-1}$$

- **Iterate:** Set $p = p + 1$,

- Sample $(a_{p-1}^1, \dots, a_{p-1}^N) \in \{1, \dots, N\}^N$, where, independently for each $i \in \{1, \dots, N\}$, $\mathbb{P}(a_{p-1}^i = j) = w_{p-1,\alpha,\theta}^j$.

- Sample $x_p^i(1:k_\alpha) | x_{p-1}^{a_{p-1}^i}(1:k_\alpha)$ from $\check{f}_{\theta,\alpha(1:k_\alpha)}(x_{p-1}^{a_{p-1}^i}(1:k_\alpha), \cdot)$, for $i \in \{1, \dots, N\}$, and evaluate the weight

$$w_{p,\alpha,\theta}^i = \left(\check{g}_{\theta,\alpha(1:k_\alpha)}(x_p^i(1:k_\alpha), y_p) \right) \left(\sum_{j=1}^N \check{g}_{\theta,\alpha(1:k_\alpha)}(x_p^j(1:k_\alpha), y_p) \right)^{-1}.$$

The following empirical measure then provides an approximation of $\xi_{n,\alpha(1:k_\alpha),\theta}(x_{0:n}(1:k_\alpha))$

$$\sum_{i=1}^N w_{n,\alpha,\theta}^i \delta_{\bar{x}_{0:n}^i(1:k_\alpha)}(dx_{0:n}). \quad (20)$$

This will prove useful in the next section.

The normalization constant $Z_{n,\alpha,\theta} = \int_{\otimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1}} \xi_{n,\alpha(1:k_\alpha),\theta}(x_{0:n}(1:k_\alpha)) dx_{0:n}(1:k_\alpha)$ can be unbiasedly estimated [20] by

$$Z_{n,\alpha,\theta}^N = \prod_{p=0}^n \left(\frac{1}{N} \sum_{i=1}^N \check{g}_{\theta,\alpha(1:k_\alpha)}(x_p^i(1:k_\alpha), y_p) \right). \quad (21)$$

It is noted that particle filters often do not work well in high-dimensions (e.g. [21]). However, in some cases, where the target probability is a high and finite dimensional discretization of an infinite dimensional distribution (as will be the case in the context of this article), the algorithm can work quite well; see e.g. [22].

4.2 Particle MCMC

In this section, we focus on approximating expectations w.r.t. $\xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)$, for a fixed n . The notation $\bar{x}_{0:n}^{(i)}(1:k_\alpha), \theta^i$ will refer to the i^{th} sample of a Markov chain designed to approximate expectations w.r.t. $\xi_{n,\alpha(1:k_\alpha)}$. In the algorithms to be presented, $r(\theta^{i-1}, \cdot)$ is a proposal density on Θ which we will assume is a postive probability density w.r.t. $d\theta$ for any θ^{i-1} . In Algorithm 2 we present an approach to approximate expectations w.r.t. $\xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)$. The method in Algorithm 2 is called particle MCMC (PMCMC).

Algorithm 2: PMCMC Algorithm

- **Initialize:**

- (0) Set $i = 0$ and sample θ^0 from the prior. Given θ^0 run the particle filter in Algorithm 1 and record the estimate of Z_{n,α,θ^0}^N from eq. (21).
- (I) Select a trajectory $x_{0:n}^i(1:k_\alpha)$ from the particle filter just run using (20), denote the stored state $\bar{x}_{0:n}^{(0)}(1:k_\alpha)$.

- **Iterate:**

- (II) Set $i = i + 1$ and propose θ' given θ^{i-1} from a proposal $r(\theta^{i-1}, \cdot)$ (described in the main text).
- (III) Given θ' run the particle filter in Algorithm 1 and record the estimate $Z_{n,\alpha,\theta'}^N$.
- (IV) Select a trajectory $x_{0:n}^{s'}(1:k_\alpha)$ from the particle filter just run using (20).
- (V) Set $\theta^i = \theta'$ with probability

$$\min \left\{ 1, \frac{Z_{n,\alpha,\theta'}^N \nu(\theta') r(\theta', \theta^{i-1})}{Z_{n,\alpha,\theta^{i-1}}^N \nu(\theta^{i-1}) r(\theta^{i-1}, \theta')} \right\}$$

otherwise $\theta^i = \theta^{i-1}$.

- (VI) Let $\bar{x}_{0:n}^{(i)}(1:k_\alpha) = x_{0:n}^{s'}(1:k_\alpha)$ if $\theta^i = \theta'$, otherwise let $\bar{x}_{0:n}^{(i)}(1:k_\alpha) = \bar{x}_{0:n}^{(i-1)}(1:k_\alpha)$ if $\theta^i = \theta^{i-1}$.
-

The target density associated to the PMCMC kernel described in Algorithm 2 on the state-space $\Theta \times \left(\otimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1} \right)^N \times$

$\{1, \dots, N\}^{Nn+1}$ is given by

$$\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta, s) := \frac{1}{Z_{n,\alpha}} w_{n,\alpha,\theta}^s Z_{n,\alpha,\theta}^N \Psi_{\alpha,\theta}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}),$$

where we recall that $\Psi_{\alpha,\theta}$ is the density of the particle filter (17), and $Z_{n,\alpha} = \int_{\Theta} Z_{n,\alpha,\theta}$ is the normalizing constant, due to the unbiased property of (21). This is shown in [4, Theorem 4], as well as the fact that

$$\begin{aligned} & \tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta, s) \\ &= \frac{\xi_{n,\alpha(1:k_\alpha)}(\bar{x}_{0:n}^s(1:k_\alpha), \theta)}{N^{n+1}} \frac{\Psi_{\alpha,\theta}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N})}{\check{\mu}_{\theta,\alpha(1:k_\alpha)}(x_0^{b_p^s}(1:k_\alpha)) \prod_{p=1}^n w_{p-1,\alpha,\theta}^{b_p^s} \check{f}_{\theta,\alpha(1:k_\alpha)}(x_{p-1}^{b_p^s}(1:k_\alpha), x_p^{b_p^s}(1:k_\alpha))}, \end{aligned} \quad (22)$$

where b_p^s is defined in (18). In other words $(\bar{x}_{0:n}^s(1:k_\alpha), \theta)$ has marginal density $\xi_{n,\alpha(1:k_\alpha)}$. As a result, consider estimating the integral

$$\int_{\otimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1} \times \Theta} \Phi(x_{0:n}(1:k_\alpha), \theta) \xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta) dx_{0:n}(1:k_\alpha) d\theta$$

for an integrable and real-valued function $\Phi : \otimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1} \times \Theta \rightarrow \mathbb{R}$. [4] show that, for Algorithm 2, this above quantity is consistently estimated (that is, the estimate converges almost surely as $N \rightarrow +\infty$) by:

$$\frac{1}{N} \sum_{i=1}^N \Phi(\bar{x}_{0:n}^{(i)}(1:k_\alpha), \theta^i). \quad (23)$$

As a result, for n fixed one can approximate the r.h.s. of (16) as,

$$\sum_{l=1}^{k'_\alpha} \tau_{l,\alpha} \left(\frac{\sum_{i=1}^N \varphi_{\alpha(2l)}(\bar{x}_{0:n}^{(i)}(l), \theta^i) H_{2l,n,\alpha,\theta}(\bar{x}_{0:n}^{(i)}(1:k_\alpha), \theta^i)}{\sum_{i=1}^N H_{2l,n,\alpha,\theta}(\bar{x}_{0:n}^{(i)}(1:k_\alpha), \theta^i)} - \frac{\sum_{i=1}^N \varphi_{\alpha(2l-1)}(\bar{x}_{0:n}^{(i)}(l), \theta^i) H_{2l-1,n,\alpha,\theta}(\bar{x}_{0:n}^{(i)}(1:k_\alpha), \theta^i)}{\sum_{i=1}^N H_{2l-1,n,\alpha,\theta}(\bar{x}_{0:n}^{(i)}(1:k_\alpha), \theta^i)} \right),$$

where we remind the reader that $\tau_{i,\alpha} = (-1)^{|\alpha(k_\alpha) - \alpha(2i)|}$. We hence refer to using Algorithm 2 in the context of (16) within (10) as above as MIPMCMC. Note that steps (I), (IV), (VI) can be ignored if one is only interested in estimation related to θ . The resulting chain has the marginal of (22) as its target: $\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta)$ ([2,4]). In the context of M—PMCMC one must compute the weights in the expression above regardless, and so the joint distribution is required.

4.3 SMC²

To consider the method to be discussed, we start with some definitions. We define the spaces:

$$\begin{aligned} \mathbb{E}_0 &:= \left(\bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)} \right)^N \times \Theta \\ \mathbb{E}_n &:= \left(\bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1} \right)^N \times \{1, \dots, N\}^{nN} \times \Theta \quad n \geq 1 \end{aligned}$$

and states

$$\begin{aligned} U_{0,\alpha} &:= (X_0^{1:N}(1:k_\alpha), \theta) \\ U_{n,\alpha} &:= (X_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta) \quad n \geq 1 \end{aligned}$$

and note $U_{0,\alpha} \in \mathbb{E}_0$, $U_{n,\alpha} \in \mathbb{E}_n$, $n \geq 1$.

We now introduce the SMC² method in Algorithm 3. This method is a type of particle filter which targets the sequence of probability distributions $\{\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta)\}_{n \geq 0}$ (see [2, Proposition 1] for a justification) which admit $\{\xi_{n,\alpha(1:k_\alpha)}(\theta)\}_{n \geq 0}$ as a particular marginal; we explain how one can estimate expectations w.r.t. $\xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)$ below. The convergence (as $N_\alpha \rightarrow \infty$) of such an algorithm, then follows the theory of particle approximations of Feynman-Kac formulae, as described in [20].

4.3.1 Intuition of the Algorithm

To understand the intuition of the algorithm of [2], consider the first step, where the objective is to approximate expectations w.r.t.

$$\tilde{\xi}_{0,\alpha(1:k_\alpha)}(x_0^{1:N}(1:k_\alpha), \theta) \propto \left(\frac{1}{N} \sum_{j=1}^N \check{g}_{\theta,\alpha(1:k_\alpha)}(x_0^j(1:k_\alpha), y_0) \right) \left(\prod_{j=1}^N \check{\mu}_{\theta^i,\alpha(1:k_\alpha)}(x_0^j(1:k_\alpha)) \right) \nu(\theta).$$

This can be achieved by (self-normalized) importance sampling, just as in the initialization step of Algorithm 3. This is because the term $G_{0,\alpha}(u_{0,\alpha})$ is an importance weight, which allows one to correct for the discrepancy between the distribution sampled (in the initialization step of Algorithm 3) and the one of interest.

Algorithm 3: An SMC² Algorithm targeting $\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta)$

- **Initialize.** Set $n = 0$, for $i \in \{1, \dots, N_\alpha\}$ sample θ^i from the prior ν , and $X_0^{i,j}(1:k_\alpha)$ from $\check{\mu}_{\theta^i,\alpha(1:k_\alpha)}(\cdot)$, $j \in \{1, \dots, N\}$. Compute the weight:

$$G_{0,\alpha}(u_{0,\alpha}^i) = \frac{1}{N} \sum_{j=1}^N \check{g}_{\theta^i,\alpha(1:k_\alpha)}(x_0^{i,j}(1:k_\alpha), y_0).$$

- **Iterate:**

- (I) **Select:** Set $n = n + 1$, resample $u_{n-1,\alpha}^{1:N_\alpha}$ using the normalized $\{G_{n-1,\alpha}(u_{n-1,\alpha}^i)\}_{i=1}^{N_\alpha}$, denoting the resulting samples $\hat{u}_{n-1,\alpha}^{1:N_\alpha}$.
- (II) **Mutate:** For $i \in \{1, \dots, N_\alpha\}$ generate $\tilde{u}_{n-1,\alpha}^i | \hat{u}_{n-1,\alpha}^i$ using one iteration of Algorithm 2 (here we only require samples from the marginal $\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta)$ and so steps (I), (IV), (VI) of Algorithm 2 can be ignored).
- (III) **Extend:** For $i \in \{1, \dots, N_\alpha\}$ sample $X_n^{i,j}(1:k_\alpha), a_{n-1}^{i,j}$, $j \in \{1, \dots, N\}$ from

$$\prod_{j=1}^N \frac{\check{g}_{\hat{\theta}^i,\alpha(1:k_\alpha)}(\tilde{x}_{n-1}^{i,a_{n-1}^{i,j}}(1:k_\alpha), y_{n-1})}{\sum_{l=1}^N \check{g}_{\hat{\theta}^i,\alpha(1:k_\alpha)}(\tilde{x}_{n-1}^{i,l}(1:k_\alpha), y_{n-1})} \check{f}_{\hat{\theta}^i,\alpha(1:k_\alpha)}(\tilde{x}_{n-1}^{i,a_{n-1}^{i,j}}(1:k_\alpha), x_n^{i,j}(1:k_\alpha)).$$

Set $u_{n,\alpha}^i = (\tilde{u}_{n-1,\alpha}^i, x_n^{i,1:N}(1:k_\alpha), a_{n-1}^{i,1:N})$.

- (IV) **Compute the weight:** For $i \in \{1, \dots, N_\alpha\}$

$$G_{p,\alpha}(u_{n,\alpha}^i) = \frac{1}{N} \sum_{j=1}^N \check{g}_{\hat{\theta}^i,\alpha(1:k_\alpha)}(x_n^{i,j}(1:k_\alpha), y_n).$$

We now want to move our samples, in such a way as to approximate expectations w.r.t. $\tilde{\xi}_{1,\alpha(1:k_\alpha)}(x_{0:1}^{1:N}(1:k_\alpha), a_0^{1:N}, \theta)$. This can be achieved in the iterate step of Algorithm 3, as we now explain. In the select step, this is a resampling of the samples, just as in the particle filter. The resulting samples are approximately sampled from $\tilde{\xi}_{0,\alpha(1:k_\alpha)}(x_0^{1:N}(1:k_\alpha), \theta)$. The mutate step will now produce new samples which are still approximately sampled from $\tilde{\xi}_{0,\alpha(1:k_\alpha)}(x_0^{1:N}(1:k_\alpha), \theta)$, as the transition kernel leaves this probability invariant. The extend step now produces the additional random variables needed to approximate expectations w.r.t. $\tilde{\xi}_{1,\alpha(1:k_\alpha)}(x_{0:1}^{1:N}(1:k_\alpha), a_0^{1:N}, \theta)$. The strategy employed leads to the convenient weight function $G_{1,\alpha}(u_{1,\alpha})$ in the next step. This corresponds to the ratio, up-to a normalizing constant, of $\tilde{\xi}_{1,\alpha(1:k_\alpha)}(x_{0:1}^{1:N}(1:k_\alpha), a_0^{1:N}, \theta)$ to the product of $\tilde{\xi}_{0,\alpha(1:k_\alpha)}(x_0^{1:N}(1:k_\alpha), \theta)$ the proposal used in the extend step. Expectations w.r.t. $\tilde{\xi}_{1,\alpha(1:k_\alpha)}(x_{0:1}^{1:N}(1:k_\alpha), a_0^{1:N}, \theta)$ can now be approximated again by self-normalized importance sampling. The algorithm then just continues for the rest of the sequence $\{\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta)\}_{n \geq 2}$.

The reason why one considers the sequence $\{\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta)\}_{n \geq 0}$, instead of the original $\{\xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)\}_{n \geq 0}$, is because the algorithm associated with the former is expected to be more efficient than a related SMC algorithm for the latter. To explain further, one can consider the naive algorithm which samples the initial θ from the prior (i.e. N samples) and then run a particle filter with N associated trajectories $x_{0:n}(1:k_\alpha)$. The main issue here is of course that one never updates the θ , so that estimates of expectations associated to θ would be very poor. This is further exacerbated by the path degeneracy problem for particle filters; one does not update the trajectory in the past, and due to the resampling operation the distinctness of the trajectories in the past will be essentially lost. These latter issues can be circumvented by applying an MCMC kernel of invariant measure $\xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta)$ at each time step to each sample. However, as we have remarked, in general PMCMC is considered to be more efficient than this. The choice of N is an important tuning parameter, which is considered in detail in the work [23]. The recommendation there is $N \propto n$. See [2,5] for further insights. We do not consider biased methods such as [24] here.

4.3.2 Estimating Expectations with Respect to the Joint Target

We now concisely describe how to use the samples $\tilde{u}_{p,\alpha}^{1:N_\alpha}$ in order to estimate expectations with respect to the joint (coupled) state and parameter, using the SMC analogue of the PMCMC estimator (23), hence enabling estimation of (16). As described in Section 4.2 and Algorithm 2, we require samples from $\tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta, s)$ to achieve this. As suggested in Algorithm 3 step (II), assume for the moment that we have ignored steps (I), (IV), (VI) of Algorithm 2, so for each SMC sampler particle we now need to sample S and construct $\bar{x}_{0:n}^s(1:k_\alpha)$. Of course the full PMCMC Algorithm 2 could be used, in which case no further work is required. However, since s and $\bar{x}_{0:n}^s(1:k_\alpha)$ do not appear in the recursion, and are not required until one estimates a joint expectation, this is considered as a separate step here.

At any time n , sample $S_n^i \in \{1, \dots, N\}$, $i \in \{1, \dots, N_\alpha\}$ with probability

$$\mathbb{P}(S_n^i = j | \tilde{u}_{n,\alpha}^i) = \frac{\check{g}_{\tilde{\theta}^i, \alpha(1:k_\alpha)}(\tilde{x}_n^{i,j}(1:k_\alpha), y_n)}{\sum_{l=1}^N \check{g}_{\tilde{\theta}^i, \alpha(1:k_\alpha)}(\tilde{x}_p^{i,l}(1:k_\alpha), y_n)}. \quad (24)$$

For $p \leq n$ recall the definition (18), and define

$$\bar{x}_p^{(i)}(1:k_\alpha) := \tilde{x}_p^{i, S_n^i}(1:k_\alpha). \quad (25)$$

Note this construction corresponds to sampling from (20), as is done in Algorithm 2. The notation $\bar{x}_p^{(i)}(1:k_\alpha)$ (which is redundant with the PMCMC notation) has been used on purpose. Now an SMC consistent estimator can be constructed analogous to the PMCMC estimator given in (23), using the following empirical measure

$$\eta_{n,\alpha}^{N_\alpha} := \frac{1}{N} \sum_{i=1}^{N_\alpha} \delta_{\bar{x}_{0:n}^{(i)}(1:k_\alpha)}.$$

Let $l \in \{1, \dots, k_\alpha\}$, $\varphi : \mathbb{N}_0^d \times \mathcal{X}^{n+1} \times \Theta \rightarrow \mathbb{R}$, $H : \mathbb{N}_0^d \times \bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1} \times \Theta \rightarrow \mathbb{R}$ (measurable and integrable w.r.t. $\xi_{n,\alpha(1:k_\alpha)}$). Recall that following from (22) one has

$$\begin{aligned} & \int_{E_p \times \{1, \dots, N\}} \varphi_{\alpha(l)}(\bar{x}_{0:n}^s(l), \theta) H_{\alpha}(\bar{x}_{0:n}^s(1:k_\alpha), \theta) \tilde{\xi}_{n,\alpha(1:k_\alpha)}(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta, s) d(x_{0:n}^{1:N}(1:k_\alpha), a_{0:n-1}^{1:N}, \theta, s) \\ &= \int_{\bigotimes_{l=1}^{k_\alpha} \mathcal{X}_{\alpha(l)}^{n+1} \times \Theta} \varphi_{\alpha(l)}(x_{0:n}(l), \theta) H_{\alpha}(x_{0:n}(1:k_\alpha), \theta) \xi_{n,\alpha(1:k_\alpha)}(x_{0:n}(1:k_\alpha), \theta) d(x_{0:n}(1:k_\alpha), \theta). \end{aligned}$$

One can now consistently estimate the above expectation analogous to (23) using (25) as follows

$$\eta_{n,\alpha}^{N_\alpha}(\varphi_{\alpha(l)} H_{\alpha}) = \frac{1}{N_\alpha} \sum_{i=1}^{N_\alpha} \varphi_{\alpha(l)}(\bar{x}_{0:n}^{(i)}(l), \tilde{\theta}^i) H_{\alpha}(\bar{x}_{0:n}^{(i)}(1:k_\alpha), \tilde{\theta}^i).$$

Hence we have the following estimate of (16)

$$\Delta \mathbb{E}_{\pi_{n,\alpha}^{N_\alpha}}[\varphi_{\alpha}(X_{0:n}, \theta)] := \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \left(\frac{\eta_{n,\alpha}^{N_\alpha}(\varphi_{\alpha(2i)} H_{2i,n,\alpha,\theta})}{\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})} - \frac{\eta_{n,\alpha}^{N_\alpha}(\varphi_{\alpha(2i-1)} H_{2i-1,n,\alpha,\theta})}{\eta_{n,\alpha}^{N_\alpha}(H_{2i-1,n,\alpha,\theta})} \right),$$

where we remind the reader that $\tau_{i,\alpha} = (-1)^{|\alpha(k_\alpha) - \alpha(2i)|}$.

5. THEORETICAL RESULTS

We now consider the MISMC² procedure in the previous section, however the results naturally extend to the static MIPMCMC method as well, which appears as a component of this method. We will analyze the variance of our MI method(s), under the following assumptions.

(A1) There exist $0 < \underline{C} < \bar{C} < +\infty$ such that for every $\alpha \in \mathcal{I}$, $\theta \in \Theta$, $(x, y) \in \mathcal{X}_\alpha \times \mathcal{Y}_\alpha$

$$\underline{C} \leq g_{\theta,\alpha}(x, y) \leq \bar{C}.$$

(A2) For every $n \geq 0$, $\varphi : \mathbb{N}_0^d \times \mathcal{X}^{n+1} \times \Theta \rightarrow \mathbb{R}$ bounded, every $\alpha \in \mathcal{I}$, there exist a $C(\alpha(1:k_\alpha))$, with $\lim_{\min_{1 \leq i \leq d} \alpha_i \rightarrow +\infty} C(\alpha(1:k_\alpha)) = 0$, such that for any collection of scalar, bounded random variables $\beta(\alpha(1:k_\alpha), 2i, 2i-1)$, $i \in \{1, \dots, k'_\alpha\}$ we have almost surely

$$\begin{aligned} & \sup_{(x_{0:n}(1:k_\alpha), \theta) \in (\bigotimes_{i=1}^{k_\alpha} \mathcal{X}_{\alpha(i)}^{n+1}) \times \Theta} \left| \left\{ \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \beta(\alpha(1:k_\alpha), 2i, 2i-1) \left[\varphi_{\alpha(2i)}(x_{0:n}(1:k_\alpha), \theta) - \right. \right. \right. \\ & \left. \left. \left. \varphi_{\alpha(2i-1)}(x_{0:n}(1:k_\alpha), \theta) \right] \right\} \right| \leq C(\alpha(1:k_\alpha)) \sum_{i=1}^{k'_\alpha} |\beta(\alpha(1:k_\alpha), 2i, 2i-1)|^2. \end{aligned}$$

We remind the reader again that $\tau_{i,\alpha} = (-1)^{|\alpha(k_\alpha) - \alpha(2i)|}$.

(A1) is a strong, but standard, assumption that has been used in the HMM literature, particularly in the SMC context; see for instance [20]. (A2) is certainly non-standard and in general one would like to deduce under simpler hypotheses. In our efforts to achieve this, we have not found a suitable technical approach and leave this more involved analysis to future work. The following result is the culmination of our efforts. The expectation below is w.r.t. the randomness in the SMC² algorithm.

Theorem 1. Assume (A1)-(A2). Then for every $n \geq 0$, $\varphi : \mathbb{N}_0^d \times \mathbb{X}^{n+1} \times \Theta \rightarrow \mathbb{R}$ bounded and every $\alpha \in \mathcal{I}$, there exist a $C(\alpha(1 : k_\alpha))$, with $\lim_{\min_{1 \leq i \leq d} \alpha_i \rightarrow +\infty} C(\alpha(1 : k_\alpha)) = 0$, such that:

$$\mathbb{E} \left[\left(\Delta \mathbb{E}_{\pi_n, \alpha}^{N_\alpha} [\varphi_\alpha(X_{0:n}, \theta)] - \Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)] \right)^2 \right] \leq \frac{C(\alpha(1 : k_\alpha))}{N_\alpha}$$

and

$$\left| \mathbb{E} \left[\Delta \mathbb{E}_{\pi_n, \alpha}^{N_\alpha} [\varphi_\alpha(X_{0:n}, \theta)] - \Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)] \right] \right| \leq \frac{C(\alpha(1 : k_\alpha))}{N_\alpha}.$$

Proof. Follows directly from Lemma 1 and Proposition 2 in the appendix. \square

It is noted that our bound depends upon the time parameter and d and we do not address these aspects in our subsequent discussion.

5.1 MIMC Considerations

Define a multi-index estimator as

$$\widehat{\varphi}_{\mathcal{I}}^{\text{MI}} := \sum_{\alpha \in \mathcal{I}} \Delta \mathbb{E}_{\pi_n, \alpha}^{N_\alpha} [\varphi_\alpha(X_{0:n}, \theta)].$$

Below the necessary assumptions are given, which are common for multi-index methods. $\text{Cost}(X_\alpha)$ denotes the cost of sampling the discretized random variable X_α . Recall $C(\alpha(1 : k_\alpha))$ appears in Theorem 1 and Assumption (A2).

Assumption 1.

For every $n \geq 0$, there exists $C < +\infty$, $w_i, \beta_i, \gamma_i > 0$ for $i = 1, \dots, d$, such that for every $\alpha \in \mathbb{N}_0^d$:

- (a) $|\Delta \mathbb{E}_{\pi_n, \alpha} [\varphi_\alpha(X_{0:n}, \theta)]| \leq C \prod_{i=1}^d 2^{-w_i \alpha_i}$;
- (b) $C(\alpha(1 : k_\alpha)) \leq C \prod_{i=1}^d 2^{-\beta_i \alpha_i}$;
- (c) $\text{Cost}(X_\alpha) \leq C \prod_{i=1}^d 2^{\gamma_i \alpha_i}$.

Before presenting the main MISMC² theorem, we need to introduce some index sets, which relate to Assumption 1. The tensor product index set is defined by

$$\mathcal{I}_{\alpha^*} := \{ \alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}_0^d : \alpha_1 \in \{0, \dots, \alpha_1^*\}, \dots, \alpha_d \in \{0, \dots, \alpha_d^*\} \}. \quad (26)$$

The total degree index set for $L \in \mathbb{R}_+$ and $\zeta \in \mathbb{R}_+^d$ is defined as

$$\mathcal{I}_{\zeta, L}^{\text{TD}} := \left\{ \alpha \in \mathbb{N}_0^d : \sum_{i=1}^d \alpha_i \zeta_i \leq L \right\}. \quad (27)$$

It is suggested in Sec. 2.2 of [1] (and verified in numerical examples) that the optimal index set is given by $\mathcal{I}_L^{\text{TD}} = \mathcal{I}_{\zeta^*, L}^{\text{TD}}$, where $z_i^* \propto \log(2)(w_i + (\gamma_i - \beta_i)/2)$ and $\sum_{i=1}^d z_i^* = 1$. See [1] for a detailed investigation of the various relationships between the rates of convergence and these index sets. The methodology developed is applicable to general index sets, but we will present the proposition for only a simplified set of circumstances in the interest of clarity and simplicity.

Proposition 1. Assume (A1)-(A2), Assumption 1 and that $\beta_i > \gamma_i$, for all $i = 1, \dots, d$, and one of the following cases holds

$$[\text{A}] \quad \mathcal{I} = \mathcal{I}_{\alpha^*} \text{ and } \sum_{i=1}^d \gamma_i / w_i \leq 2; \text{ or}$$

*To be precise, there would typically be different constants, but it obviously suffices to take the largest.

[B] $\mathcal{I} = \mathcal{I}_L^{\text{TD}}$.

Then there exist $C < +\infty$, and either $\alpha^* = (m_1, \dots, m_d) \in \mathbb{N}_0^d$ in case [A] or $L \in \mathbb{R}_+$ in case [B], and $\{N_\alpha\}_{\alpha \in \mathcal{I}}$, such that for any $\varepsilon > 0$:

$$\mathbb{E} \left[\left(\widehat{\varphi}_{\mathcal{I}}^{\text{MI}} - \mathbb{E}[\varphi(X_{0:n}, \theta)] \right)^2 \right] \leq C\varepsilon^2,$$

for a cost of $\mathcal{O}(\varepsilon^{-2})$.

Proof. Standard optimization of cost as a function of N_α for a fixed variance yields that

$$N_\alpha = \varepsilon^{-2} \left(\sum_{\alpha \in \mathcal{I}} \sqrt{C(\alpha(1:k_\alpha)) \text{Cost}(X_\alpha)} \right)^{-1} \sqrt{C(\alpha(1:k_\alpha)) / \text{Cost}(X_\alpha)},$$

where $C(\alpha(1:k_\alpha))$ and $\text{Cost}(X_\alpha)$ are defined in Assumption 1 (b-c). Under the assumptions above, and following from Theorem 1, the proof for case [A] is the same as that of Proposition 3.2 in [9]. Case [B] follows from Theorem 2.2 of [1] (see also Theorem 2 of [7]). \square

Note that ε^2 here represents the asymptotic MSE and the proposition above relates this to the cost. Simply put, the proposition states that the cost is proportional to the inverse of the MSE, which is called the *canonical case* because it is the best one can expect from any Monte Carlo method. This can be readily generalized to different relationship between the coefficients (w_i, β_i, γ_i) . There are many different cases in general, but the rules of thumb are that (i) the complexity has a logarithmic penalty if $\beta_j \leq \gamma_j$ for any j , and (ii) there is a smaller exponent on ε as well if $\beta_j < \gamma_j$ for any j . The various conditions can be derived in a similar manner as in [1] (see also [7] and [9] for some discussion). Note that if $\sum_{i=1}^d \gamma_i/w_i > 2$ instead in case [A] then the cost is $\varepsilon^{-\sum_{i=1}^d \gamma_i/w_i}$, corresponding to the cost of a single realization at the finest discretization. In this case, the cost of MLMC will still be at least as large, because a single realization at the finest discretization of the tensor product index set will always be required.

6. NUMERICAL RESULTS

6.1 Modelling

We illustrate the performance of the proposed methods on the Bayesian parameter inference problem of a partially observed stochastic system which is the solution to a SPDE. Comparisons are made with sampling from the most precise discretization of the underlying stochastic system using either PMCMC or SMC². The objective here is to illustrate the theory and test the applicability of the method under weaker assumptions than provided by the theory. Therefore, we will restrict attention to the total degree index set \mathcal{I}_α , despite its suboptimality in this example.

We consider the stochastic heat equation on a one-dimensional domain $[0, 1]$ over the time interval $[0, T]$, i.e.

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + au + \theta \dot{W}_t$$

with the Dirichlet boundary condition and initial value $u(x, 0) = u_0(x) = \sum_{k=1}^{\infty} u_{k,0} e_k(x)$ for $x \in (0, 1)$. The eigenfunction $e_k(x) = \sqrt{2} \sin(k\pi x)$ has the corresponding eigenvalue $\lambda_k = k^2\pi^2$ and the noise W_t is the space-time white noise, i.e. the cylindrical Brownian motion given by $W_t = \sum_{k=1}^{\infty} \sqrt{q_k} e_k \beta_t^k$, where β_t^k ($k \geq 1$) are i.i.d. scalar Brownian motions. The hidden process is assumed to be modelled by the solution to this SPDE with $q_k = 1$ and $u_{k,0} = 1$ for $1 \leq k \leq K_{\max}$ and $u_{k,0} = 0$ otherwise.

Pointwise observations of the process are obtained at times $t(n) = n\delta$ for $n = 1, 2, \dots, 100$ and $\delta = 0.001$, and at the locations $x_1 = 1/3$ and $x_2 = 2/3$ under an additive Gaussian noise with mean zero and variance $\tau^2 = 1$. If we denote the observation vector at time $t(n)$ by $y_n = (y_{n,1}, y_{n,2})^T$, the corresponding likelihood function is

$$g(x_n, y_n) \propto \prod_{i=1}^2 \exp \left(-\frac{1}{2\tau^2} (y_{n,i} - u(x_i, t(n)))^2 \right)$$

where $u(x_i, t(n))$ is the solution of the above SPDE at time $t(n)$ and location x_i and note that $u(x, t) = \sum_{k=1}^{\infty} u_{k,t} e_k(x)$. The model parameter θ is assumed to be unknown and is assigned a prior distribution $\text{Gamma}(1, \sqrt{0.1})$ where $\text{Gamma}(a, b)$ represents the Gamma distribution with shape parameter a and scale parameter b . A fixed sequence of observations $y_{1:100}$ is simulated with $a = 1/2$ and $\theta = \sqrt{0.1}$.

The problem of interest is the Bayesian static parameter estimation of θ from the above-mentioned model sequentially for each n as the data arrives. Our ultimate goal is to approximate $\mathbb{E}_{\pi_n}[\varphi(\theta)]$, where $\varphi(\theta) = \theta$ and π_n is the posterior density of θ , given $y_{1:n}$, induced by the HMM with no discretization bias. In this case, we are interested in the posterior mean of the model parameter θ .

Given the approximation multi-index $\alpha^* = (m_x, m_t)$, we will estimate $\mathbb{E}_{\pi_n, \mathcal{I}_{\alpha^*}}[\varphi(\theta)] = \mathbb{E}_{\pi_n, \alpha^*}[\varphi_{\alpha^*}(\theta)]$ to approximate $\mathbb{E}_{\pi_n}[\varphi(\theta)]$, where π_n, α^* is the posterior distribution associated with the multi-index α^* and π_n is the target posterior distribution.

We adopt the exponential Euler scheme developed in [25] for discretizing the underlying hidden process. To be precise, at a multi-index $\alpha = (\alpha_x, \alpha_t)$, the above SPDE is solved with the first $K_\alpha = K_0 \times 2^{\alpha_x}$ eigenfunctions and $M_\alpha = M_0 \times 2^{\alpha_t}$ time steps as follows

$$u_{\alpha,k,i+1} = e^{-\lambda_k h} u_{\alpha,k,i} + \frac{1 - e^{-\lambda_k h}}{\lambda_k} a u_{\alpha,k,i} + r_{k,i} \quad (28)$$

where $r_{k,i} \sim N\left(0, \frac{\theta^2(1 - e^{-2\lambda_k h})}{2\lambda_k}\right)$ for $k = 1, \dots, K_\alpha$ and $i = 0, 1, \dots, M_\alpha - 1$. The time step-size $h = \delta/M_\alpha$ and $u_{\alpha,k,i}$ is the solution for the coefficient associated with the k^{th} eigenfunction, i^{th} time step and the discretization index α .

The coupling of the k_α ($1 \leq k_\alpha \leq 4$) discretized probability laws is constructed as follows. We start with the simulation of the most expensive random variable that corresponds to the multi-index α . For simulations involving $\alpha_x - 1$, only the subset of the first $K_{\alpha - e_x}$ components are retained. For simulations involving $\alpha_t - 1$, $r_{k,i}$ in (28) is replaced by $\hat{r}_{k,i} = e^{-\lambda_k h} \tilde{r}_{k,2i} + \tilde{r}_{k,2i+1}$ [26] for $i = 0, 1, \dots, M_{\alpha - e_t} - 1$, where $\{\tilde{r}_{k,i}\}_{i=0}^{M_\alpha - 1}$ are simulated with respect to the multi-index α .

Assumption 1 (b) was verified directly by estimating the quantity in Theorem 1 using the empirical variance over 20 multi-increment estimators. The values $\beta_x = 1$ and $\beta_t = 2$ were fit, for $\gamma_x = \gamma_t = 1$, which is consistent with the results in [9]. We also assume $w_i = \beta_i/2$, as in [9], and this is verified numerically. It is noted that assumption (A2) is likely not satisfied in this example, and so the numerical results are testing the applicability of the method under weaker assumptions than provided by the theory.

Following from the linear Gaussian form of (28), $u_{\alpha,k,i}$ is Gaussian. Since the observations are also linear and Gaussian, the posterior on the state path is Gaussian can be computed exactly (for fixed parameters θ) using the classical Kalman smoother [15]. In fact, it can be computed without time discretization error, as shown in [9]. Following from standard identities for Gaussian random variables, its normalizing constant (the true marginal likelihood $p(y_{0:n}|\theta)$) can be computed exactly as well. As a result, the true likelihood calculated from the Kalman smoother with high spatial resolution and no time discretization error is used within standard MCMC to produce a ground truth, denoted by $\mathbb{E}(\varphi)$. The MSE (denoted ε^2 in Proposition 1) of the approximations is then computed as follows. For a given estimator $\hat{\varphi}$ the MSE is estimated using

$$\frac{1}{R} \sum_{r=1}^R (\hat{\varphi}^{(r)} - \mathbb{E}(\varphi))^2,$$

where $\hat{\varphi}^{(r)}$ is a realization of the estimator, i.e. using MCMC, MIPMCMC, or MISMC².

6.2 Results Using PMCMC

We consider the estimation of the posterior mean of the model parameter θ in this section with n fixed and $n = 100$. The proposed MIPMCMC method is implemented, as well as a standard PMCMC at the finest discretization level. The number of particles $N = 500$ is fixed as well.

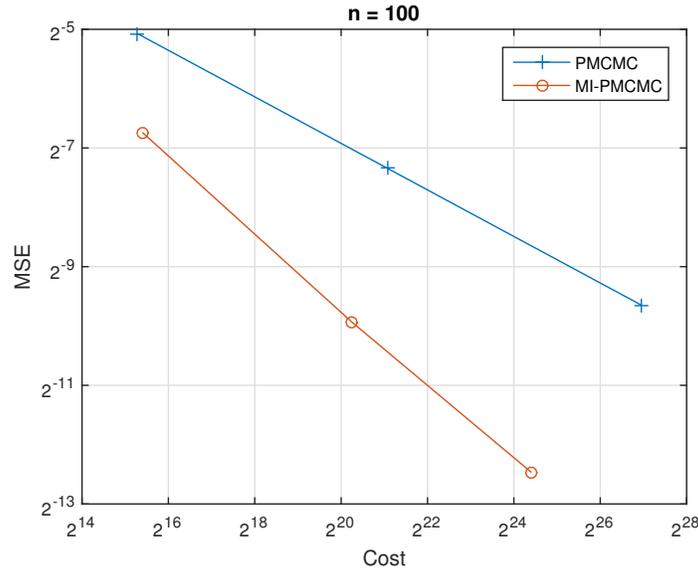


FIG. 2: MSE v.s. Cost at time index $n = 100$ for the PMCMC method

Following the optimal choice of discretization $K = M^2$ as discussed in [25], the cost for a single realization is proportional to M^3 . This results in the optimal cost for the ordinary PMCMC of $\mathcal{O}(\varepsilon^{-5})$. For the MIPMCMC method, we let $m_x = 2m_t \geq 2\log(\varepsilon/2)$ and use the optimal $N_\alpha \propto \varepsilon^{-2}m_x 2^{-\alpha_x - 3\alpha_t/2}$, as mentioned in the proof of Proposition 1 and discussed further in [9] and references therein. The proportionality sign arises because the constants are unknown in the terms appearing in Assumption 1 (b-c). In practice these are estimated along with the rates using simulations at lower levels. The cost is dominated by $\mathcal{O}(\varepsilon^{-3})$ (where $3 = \sum_{i=x,t} \gamma_i/w_i$ – see discussion following Proposition 1). Both algorithms are implemented for 20 runs and with the most precise discretization index $\alpha^* = (2, 1), (4, 2), (6, 3)$. The MSE is then estimated using these $R = 20$ realizations.

The MSE vs cost plot is illustrated in Figure 2. The cost rates are verified numerically as in Figure 2 and are consistent with [9]. The fitted rate is about -5.1 for the ordinary PMCMC method and -3.1 for MIPMCMC. It is noted that in the context of MLPMCMC for this example, i.e. refining once in both (x, t) at each level, one will find $\gamma = 3, \beta = 2, \alpha = \beta/2$, and $2 + (\gamma - \beta)/\alpha = 3 (= \gamma/\alpha)$. In other words, the rate is the same as we obtain here for MIPMCMC [7].

6.3 Results on SMC²

The proposed SMC² method as well as the ordinary SMC² method are implemented with the most precise discretization indices $\alpha^* = (2, 1), (4, 2), (6, 3), (8, 4)$ and as above the number of particles is fixed at $N = 500$. The proposed SMC² method is run with the optimal choice of $N_\alpha \propto \varepsilon^{-2}m_x 2^{-\alpha_x - 3\alpha_t/2}$ as discussed in [9] and subsection 6.2. The ground truth in this case is calculated by the weighted average of the θ particles from the iterated batch importance sampling algorithm [27] with true likelihood increments derived from the Kalman techniques, which is used for computing the MSE of the approximations. Both algorithms are implemented for 20 runs and the MSE is estimated using these realizations.

The same rate is expected as in subsection 6.2, under the same choices of (m_x, m_t) and N_α . This is verified numerically, as illustrated in Figure 3, which displays the MSE vs cost plot at different time index $n \in \{50, 65, 80, 100\}$. The fitted rate is about -5.2 for the ordinary SMC² and -3 for the multi-index SMC² method.

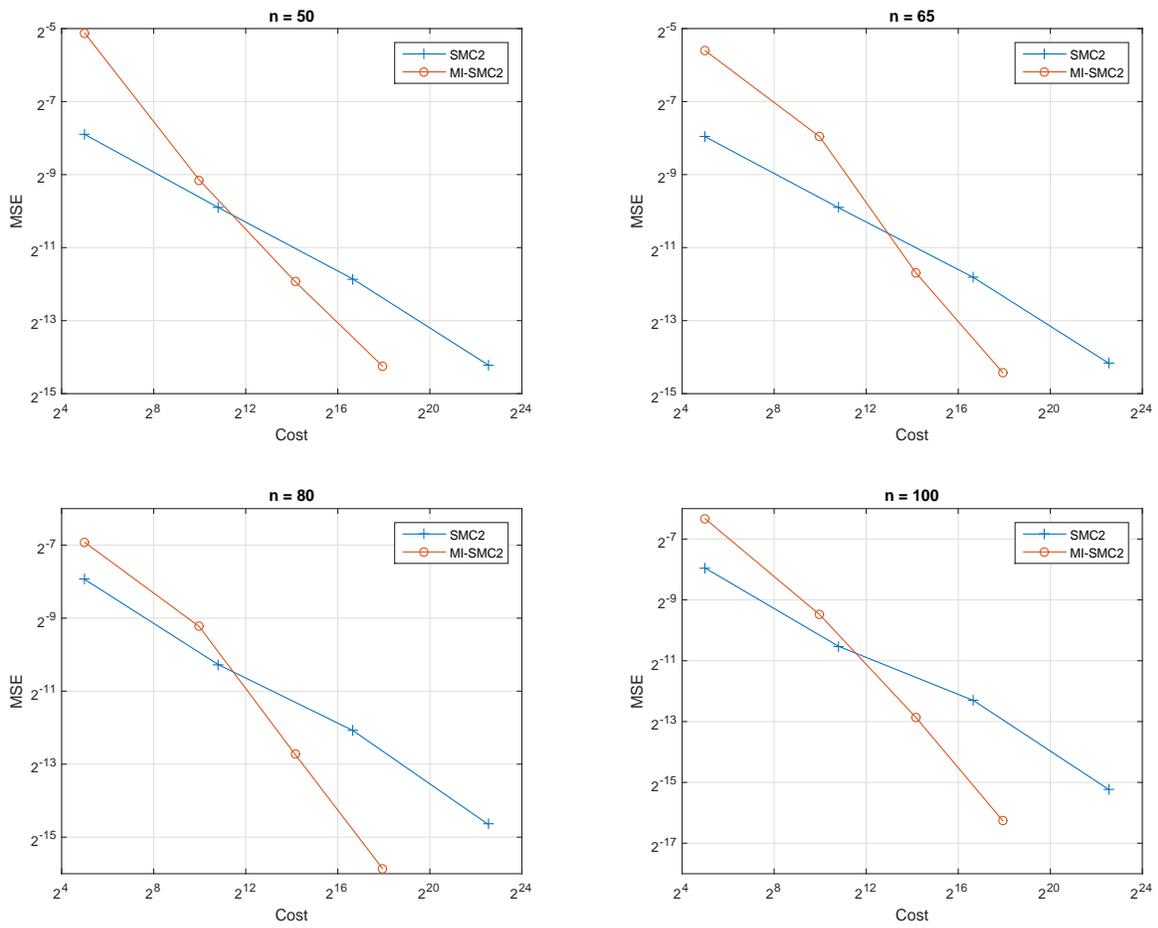


FIG. 3: MSE v.s. Cost at different time index $n \in \{50, 65, 80, 100\}$

7. CONCLUSION

MISMC² and MIPMCMC are introduced. It is proven under strong assumptions that these methods achieve a better rate of complexity with respect to MSE than their single level counterparts, and this can even be canonical 1/MSE. The algorithm is tested on a typical numerical example which may not satisfy the required assumptions, and the results are verified. This makes us optimistic that theoretical results for the present algorithm can be obtained under weaker assumptions. Another interesting direction for future research is further exploration of the method in practical scenarios and with optimal index sets.

Acknowledgements

We would like to thank Abdul-Lateef Haji-Ali for useful discussions relating to the material in this paper. AJ was supported by a KAUST CRG4 grant ref: 2584 and KAUST baseline funding. K.J.H.L. & A.J. were supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research (ASCR), under field work proposal number ERKJ333. KJHL was additionally supported by The Alan Turing Institute under the EPSRC grant EP/N510129/1. He was also funded in part by Oak Ridge National Laboratory Directed Research and Development Seed funding.

APPENDIX A. MAIN PROOFS

Let (E, \mathcal{E}) be a measurable space. The supremum norm is written as $\|f\| = \sup_{u \in E} |f(u)|$. We will consider a non-negative operator $K : E \times \mathcal{E} \rightarrow \mathbb{R}_+$, finite measure μ on (E, \mathcal{E}) and a real-valued, measurable $f : E \rightarrow \mathbb{R}$ and define the operations:

$$\mu K : A \mapsto \int K(u, A) \mu(du); \quad Kf : u \mapsto \int f(v) K(u, dv).$$

We also write $\mu(f) = \int f(u) \mu(du)$.

Recall the definition of E_p in Section 4.3 and denote by \mathcal{E}_p the associated σ -algebra. Let $p \geq 1$ and denote by $M_p : E_{p-1} \times \mathcal{E}_p \rightarrow [0, 1]$ the Markov kernel which is a composition of

1. A marginal PMCMC kernel $\bar{M}_p : E_{p-1} \times \mathcal{E}_{p-1} \rightarrow [0, 1]$, as in Algorithm 2 (ignoring s),
2. Followed by the sampling of, for $i \in \{1, \dots, N_\alpha\}$ $X_p^{i,j}(1 : k_\alpha), a_{p-1}^{i,j}, j \in \{1, \dots, N\}$ in Algorithm 3, the iterate step.

Denote by $\tilde{\eta}_0$ as the initial probability measure (on (E_0, \mathcal{E}_0)) of θ^i and $X_0^{i,j}(1 : k_\alpha), j \in \{1, \dots, N\}$ in Algorithm 3, the initialization step. Define the probability measure on $(E_p \times \{1, \dots, N\}), \mathcal{E}_p \vee 2^{\{1, \dots, N\}}$, $\eta_{p,\alpha}$:

$$\eta_{p,\alpha}(d(u_{p,\alpha}, s)) := \left(\int_{E_0 \times \dots \times E_{p-1}} \left[\prod_{q=0}^p G_p(u_{p,\alpha}) \right] \tilde{\eta}_0(du_0, \alpha) \left[\prod_{q=1}^{p-1} M_q(u_{q-1, \alpha}, du_{q,\alpha}) \right] \bar{M}_p(u_{p-1, \alpha}, du_{p,\alpha}) \times \right. \\ \left. \mathbb{P}(s|u_{p,\alpha}) ds \right) / \left(\int_{E_0 \times \dots \times E_{p-1}} \left[\prod_{q=0}^p G_p(u_{p,\alpha}) \right] \tilde{\eta}_0(du_0, \alpha) \left[\prod_{q=1}^{p-1} M_q(u_{q-1, \alpha}, du_{q,\alpha}) \right] \right)$$

where ds is counting measure and $\mathbb{P}(s|u_{p,\alpha})$ is as (24).

Note that one can easily show that (16) is equal to

$$\sum_{i=1}^{k'_\alpha} (-1)^{|\alpha(k_\alpha) - \alpha(2i)|} \left(\frac{\eta_{n,\alpha}(\varphi_{\alpha(2i)} H_{2i,n,\alpha,\theta})}{\eta_{n,\alpha}(H_{2i,n,\alpha,\theta})} - \frac{\eta_{n,\alpha}(\varphi_{\alpha(2i-1)} H_{2i-1,n,\alpha,\theta})}{\eta_{n,\alpha}(H_{2i-1,n,\alpha,\theta})} \right).$$

Recall $\tau_{i,\alpha} = (-1)^{|\alpha(k_\alpha) - \alpha(2i)|}$ and set, for each $\varphi, \alpha(i)$,

$$\zeta_{i,n,\varphi}(x_{0:n}(1 : k_\alpha), s, \theta) = \varphi_{\alpha(i)}(x_{0:n}^s(i), \theta) H_{i,n,\alpha,\theta}(x_{0:n}(1 : k_\alpha)).$$

Now set

$$\begin{aligned}\psi_{n,i,\alpha}^{N_\alpha} &:= \frac{\eta_{n,\alpha}^{N_\alpha}(\zeta_{2i-1,n,\varphi})}{\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})\eta_{n,\alpha}^{N_\alpha}(H_{2i-1,n,\alpha,\theta})} \\ \psi_{n,i,\alpha} &:= \frac{\eta_{n,\alpha}(\zeta_{2i-1,n,\varphi})}{\eta_{n,\alpha}(H_{2i,n,\alpha,\theta})\eta_{n,\alpha}(H_{2i-1,n,\alpha,\theta})} \\ \bar{\psi}_{n,i,\alpha}^{N_\alpha} &:= \psi_{n,i,\alpha}^{N_\alpha} - \psi_{n,i,\alpha}.\end{aligned}$$

In addition:

$$\begin{aligned}\Xi_{n,\alpha,1}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \left[\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} \right] [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] (\zeta_{2i,n,\varphi} - \zeta_{2i-1,n,\varphi}) \\ \Xi_{n,\alpha,2}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \bar{\psi}_{n,i,\alpha}^{N_\alpha} [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] (H_{2i,n,\alpha,\theta} - H_{2i-1,n,\alpha,\theta}) \\ \Xi_{n,\alpha,3}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] (\zeta_{2i,n,\varphi} - \zeta_{2i-1,n,\varphi}) \\ \Xi_{n,\alpha,4}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \psi_{n,i,\alpha} [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] (H_{2i,n,\alpha,\theta} - H_{2i-1,n,\alpha,\theta}) \\ \Xi_{n,\alpha,5}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \left[\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-2} \right] \times \\ &\quad \eta_{n,\alpha}(\zeta_{2i,n,\varphi} - \zeta_{2i-1,n,\varphi}) [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] (H_{2i,n,\alpha,\theta}) \\ \Xi_{n,\alpha,6}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \eta_{n,\alpha}(H_{2i,n,\alpha,\theta} - H_{2i-1,n,\alpha,\theta}) \bar{\psi}_{n,i,\alpha}^{N_\alpha} \\ \Xi_{n,\alpha,7}^{N_\alpha} &:= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-2} \eta_{n,\alpha}(\zeta_{2i,n,\varphi} - \zeta_{2i-1,n,\varphi}) [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] (H_{2i,n,\alpha,\theta}).\end{aligned}$$

Lemma 1. Assume (A1). Then for every $n \geq 0$, $\varphi : \mathbb{N}_0^d \times X^{n+1} \times \Theta \rightarrow \mathbb{R}$ bounded and $\alpha \in \mathcal{I}$ we have that:

$$\Delta \mathbb{E}_{\pi_{n,\alpha}}^{N_\alpha} [\varphi_\alpha(X_{0:n}, \theta)] - \Delta \mathbb{E}_{\pi_{n,\alpha}} [\varphi_\alpha(X_{0:n}, \theta)] = \sum_{j=1}^7 (-1)^{j+1} \Xi_{n,\alpha,j}^{N_\alpha}.$$

Proof. Follows by standard algebra. (A1) is only used to ensure the existence of all the associated quantities. \square

Proposition 2. Assume (A1)-(A2). Then for every $n \geq 0$, $\varphi : \mathbb{N}_0^d \times X^{n+1} \times \Theta \rightarrow \mathbb{R}$ bounded and $\alpha \in \mathcal{I}$, there exist a $C(\alpha(1 : k_\alpha))$, with $\lim_{\min_{1 \leq i \leq d} \alpha_i \rightarrow +\infty} C(\alpha(1 : k_\alpha)) = 0$, such that for $j \in \{1, \dots, 7\}$, $N_\alpha \geq 1$:

$$\max\{|\mathbb{E}[\Xi_{n,\alpha,j}^{N_\alpha}]|, \mathbb{E}[(\Xi_{n,\alpha,j}^{N_\alpha})^2]\} \leq \frac{C(\alpha(1 : k_\alpha))}{N_\alpha}.$$

Proof. We give the proofs in the case $j = 1$ or $j = 3$. All other cases are essentially the same and omitted. Throughout the proof $C(\alpha(1 : k_\alpha))$ is a constant that depends on $n \geq 0$, $\varphi : \mathbb{N}_0^d \times X^{n+1} \times \Theta \rightarrow \mathbb{R}$, with $\lim_{\min_{1 \leq i \leq d} \alpha_i \rightarrow +\infty} C(\alpha(1 : k_\alpha)) = 0$. The exact value of $C(\alpha(1 : k_\alpha))$ may change from line to line, but the latter property holds.

Set

$$\begin{aligned} \kappa_{n,\alpha,1}(x_{0:n}(1:k_\alpha), s, \theta) &= \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \left[\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} \right] \times \\ &\quad (\zeta_{2i,n,\varphi}(x_{0:n}(1:k_\alpha), s, \theta) - \zeta_{2i-1,n,\varphi}(x_{0:n}(1:k_\alpha), s, \theta)). \end{aligned}$$

Then

$$\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})^2] = \mathbb{E} \left[[\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] \left(\frac{\kappa_{n,\alpha,1}}{\|\kappa_{n,\alpha,1}\|} \right)^2 \|\kappa_{n,\alpha,1}\|^2 \right].$$

Clearly, applying (A2) to the term $\|\kappa_{n,\alpha,1}\|$ and using that

$$[\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] \left(\frac{\kappa_{n,\alpha,1}}{\|\kappa_{n,\alpha,1}\|} \right) \leq 2 \quad (\text{A.1})$$

it follows that

$$\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})^2] \leq C(\alpha(1:k_\alpha)) \mathbb{E} \left[\left(\sum_{i=1}^{k'_\alpha} \left| \eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} \right|^2 \right)^2 \right].$$

Application of the Minkowski inequality yields

$$\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})^2] \leq C(\alpha(1:k_\alpha)) \left(\sum_{i=1}^{k'_\alpha} \mathbb{E} \left[\left(\left[\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} \right]^4 \right)^{1/2} \right]^2 \right)^2.$$

Note that the summand

$$\mathbb{E} \left[\left(\left[\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} \right]^4 \right)^{1/2} \right] = \mathbb{E} \left[\left(\left[\frac{\eta_{n,\alpha}(H_{2i,n,\alpha,\theta}) - \eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})}{\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta}) \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})} \right]^4 \right)^{1/2} \right]$$

then applying (A1)

$$\mathbb{E} \left[\left(\left[\eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta})^{-1} - \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} \right]^4 \right)^{1/2} \right] \leq C \mathbb{E} \left[\left(\eta_{n,\alpha}(H_{2i,n,\alpha,\theta}) - \eta_{n,\alpha}^{N_\alpha}(H_{2i,n,\alpha,\theta}) \right)^4 \right)^{1/2}$$

where $C < +\infty$ is a constant that does not depend upon α . Then applying [28, Proposition 2.9] to the term on the r.h.s. of the above equation, yields that

$$\mathbb{E} \left[\|\kappa_{n,\alpha,1}\|^2 \right] \leq \frac{C(\alpha(1:k_\alpha))}{N_\alpha^2} \quad (\text{A.2})$$

and hence allows one to derive the upper-bound

$$\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})^2] \leq \frac{C(\alpha(1:k_\alpha))}{N_\alpha^2}.$$

For the bias, we have

$$|\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})]| \leq \mathbb{E} \left[\left| [\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha}] \left(\frac{\kappa_{n,\alpha,1}}{\|\kappa_{n,\alpha,1}\|} \right) \right| \|\kappa_{n,\alpha,1}\| \right]$$

it follows by using (A.1)

$$|\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})]| \leq 2 \mathbb{E}[\|\kappa_{n,\alpha,1}\|]$$

then using Jensen's inequality and (A.2), we can conclude that

$$|\mathbb{E}[(\Xi_{n,\alpha,1}^{N_\alpha})]| \leq \frac{C(\alpha(1:k_\alpha))}{N_\alpha}.$$

Set

$$\kappa_{n,\alpha,3}(x_{0:n}(1:k_\alpha), s, \theta) = \sum_{i=1}^{k'_\alpha} \tau_{i,\alpha} \eta_{n,\alpha}(H_{2i,n,\alpha,\theta})^{-1} (\zeta_{2i,n,\varphi}(x_{0:n}(1:k_\alpha), s, \theta) - \zeta_{2i-1,n,\varphi}(x_{0:n}(1:k_\alpha), s, \theta)).$$

Then

$$\mathbb{E}[(\Xi_{n,\alpha,3}^{N_\alpha})^2] = \mathbb{E}[(\eta_{n,\alpha}^{N_\alpha} - \eta_{n,\alpha})(\kappa_{n,\alpha,3})^2].$$

Applying [28, Proposition 2.9] to the term on the r.h.s. yields

$$\mathbb{E}[(\Xi_{n,\alpha,3}^{N_\alpha})^2] \leq \frac{\|\kappa_{n,\alpha,3}\|^2}{N_\alpha}.$$

Application of (A2) gives

$$\mathbb{E}[(\Xi_{n,\alpha,3}^{N_\alpha})^2] \leq \frac{C(\alpha(1:k_\alpha))}{N_\alpha}.$$

For $|\mathbb{E}[(\Xi_{n,\alpha,3}^{N_\alpha})]|$ using a similar decomposition to [13, eq. (A.2)] one can show that

$$|\mathbb{E}[(\Xi_{n,\alpha,3}^{N_\alpha})]| \leq \frac{C(\alpha(1:k_\alpha))}{N_\alpha}.$$

the proof is omitted as it is standard. □

REFERENCES

1. Haji-Ali, A.-L., Nobile, F., and Tempone, R., Multi-index Monte Carlo: when sparsity meets sampling, *Numerische Mathematik*, 132(4):767–806, 2016.
2. Chopin, N., Jacob, P. E., and Papaspiliopoulos, O., SMC²: an efficient algorithm for sequential analysis of state space models, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(3):397–426, 2013.
3. Cappé, O., Ryden, T., and Moulines, E., *Inference in Hidden Markov Models*, Springer: New York, 2005.
4. Andrieu, C., Doucet, A., and Holenstein, R., Particle Markov chain Monte Carlo methods, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342, 2010.
5. Kantas, N., Doucet, A., Singh, S. S., Maciejowski, J., and Chopin, N., On particle methods for parameter estimation in state-space models, *Statistical science*, 30(3):328–351, 2015.
6. Giles, M. B., Multilevel Monte Carlo path simulation, *Operations Research*, 56(3):607–617, 2008.
7. Giles, M. B., Multilevel Monte Carlo methods, *Acta Numerica*, 24:259–328, 2015.
8. Heinrich, S. Multilevel Monte Carlo methods. In *Large-Scale Scientific Computing Methods*, eds. S. Margenov, J. Wasniewski & P. Yalamov. Springer: Berlin, 2001.
9. Jasra, A., Kamatani, K., Law, K. J., and Zhou, Y., A multi-index Markov chain Monte Carlo method, *International Journal for Uncertainty Quantification*, 8(1), 2018.
10. Jasra, A., Kamatani, K., Law, K., and Zhou, Y., Bayesian static parameter estimation for partially observed diffusions via multilevel Monte Carlo, *SIAM Journal on Scientific Computing*, 40(2):A887–A902, 2018.
11. Franks, J., Jasra, A., Law, K., and Vihola, M., Unbiased inference for discretely observed hidden Markov model diffusions, *arXiv preprint arXiv:1807.10259*, 2018.
12. Jasra, A., Law, K., and Suci, C., Advanced multilevel Monte Carlo methods, *International Statistical Review*, to appear <https://doi.org/10.1111/insr.12365>, 2020.
13. Beskos, A., Jasra, A., Law, K., Tempone, R., and Zhou, Y., Multilevel sequential Monte Carlo samplers, *Stochastic Processes and their Applications*, 127(5):1417–1440, 2017.
14. Jasra, A., Kamatani, K., Law, K. J., and Zhou, Y., Multilevel particle filters, *SIAM Journal on Numerical Analysis*, 55(6):3068–3096, 2017.

15. Law, K., Stuart, A., and Zygalakis, K., Data assimilation, *Cham, Switzerland: Springer*, 2015.
16. Bungartz, H.-J. and Griebel, M., Sparse grids, *Acta Numerica*, 13(1):147–269, 2004.
17. Law, K. J. and Stuart, A. M., Evaluating data assimilation algorithms, *Monthly Weather Review*, 140(11):3757–3782, 2012.
18. Stuart, A. M., Inverse problems: a Bayesian perspective, *Acta Numerica*, 19:451–559, 2010.
19. Hesthaven, J. S., Gottlieb, S., and Gottlieb, D., *Spectral methods for time-dependent problems*, Vol. 21, Cambridge University Press, 2007.
20. Del Moral, P. Feynman-Kac formulae. In *Feynman-Kac Formulae*, pp. 47–93. Springer, 2004.
21. Snyder, C., Bengtsson, T., Bickel, P., and Anderson, J., Obstacles to high-dimensional particle filtering, *Monthly Weather Review*, 136(12):4629–4640, 2008.
22. Llopis, F. P., Kantas, N., Beskos, A., and Jasra, A., Particle filtering for stochastic Navier Stokes signal observed with linear additive noise, *SIAM Journal on Scientific Computing*, 40(3):A1544–A1565, 2018.
23. Doucet, A., Pitt, M. K., Deligiannidis, G., and Kohn, R., Efficient implementation of markov chain monte carlo when using an unbiased likelihood estimator, *Biometrika*, 102(2):295–313, 2015.
24. Liu, J. and West, M. Combined parameter and state estimation in simulation-based filtering. In *Sequential Monte Carlo methods in practice*, pp. 197–223. Springer, 2001.
25. Jentzen, A. and Kloeden, P. E., Overcoming the order barrier in the numerical approximation of stochastic partial differential equations with additive space–time noise, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 465(2102):649–667, 2008.
26. Chernov, A., Hoel, H., Law, K., Nobile, F., and Tempone, R., Multilevel ensemble Kalman filtering for spatially extended models, *arXiv preprint arXiv:1608.08558*, 2016.
27. Chopin, N., A sequential particle filter method for static models, *Biometrika*, 89(3):539–552, 2002.
28. Del Moral, P. and Miclo, L. Branching and interacting particle systems approximations of Feynman-Kac formulae with applications to non-linear filtering. In *Seminaire de probabilites XXXIV*, pp. 1–145. Springer, 2000.