

Adapting the Number of Particles in Sequential Monte Carlo Methods through an Online Scheme for Convergence Assessment

Víctor Elvira, *Member, IEEE*, Joaquín Míguez, Petar M. Djurić, *Fellow, IEEE*

Abstract—Particle filters are broadly used to approximate posterior distributions of hidden states in state-space models by means of sets of weighted particles. While the convergence of the filter is guaranteed when the number of particles tends to infinity, the quality of the approximation is usually unknown but strongly dependent on the number of particles. In this paper, we propose a novel method for assessing the convergence of particle filters in an online manner, as well as a simple scheme for the online adaptation of the number of particles based on the convergence assessment. The method is based on a sequential comparison between the actual observations and their predictive probability distributions approximated by the filter. We provide a rigorous theoretical analysis of the proposed methodology and, as an example of its practical use, we present simulations of a simple algorithm for the dynamic and online adaptation of the number of particles during the operation of a particle filter on a stochastic version of the Lorenz system.

Index Terms—Particle filtering, sequential Monte Carlo, convergence assessment, predictive distribution, convergence analysis, computational complexity, adaptive complexity.

I. INTRODUCTION

A. Background

Many problems in science and engineering can be described by dynamical models where hidden states of the systems change over time and observations that are functions of the states are available. Often, the observations are sequentially acquired and the interest is in making recursive inference on the hidden states. In many applications, the Bayesian approach to the problem is adopted because it allows for optimal inclusion of prior knowledge of the unknown state in the estimation process [1], [2]. In this case, the prior information and the likelihood function that relates the hidden state and the observation are combined yielding a posterior distribution of the state.

Exact Bayesian inference, however, is only possible in a small number of scenarios, including linear Gaussian state-

space models (using the Kalman filter [3], [4]) and finite state-space hidden Markov models (HMM filters [5]). Therefore, in many other practical problems, only approximate inference methods can be used. One class of suboptimal methods is particle filtering, which is also known as sequential Monte Carlo sampling [6], [7], [8], [9], [10]. Since the publication of [11], where the sampling importance resampling (SIR) filter was introduced, particle filtering has received outstanding attention in research and practice. Particle filters approximate posterior distributions of the hidden states sequentially and recursively. They do it by exploiting the principle of importance sampling and by using sets of weighted particles [6], [7], [12].

One key parameter of particle filters is the number of particles. It can be proved that the rate of convergence of the approximate probability distribution towards the true posterior is inversely proportional to the square root of the number of particles used in the filter [12], [13]. This, too, entails that the filter “perfectly” approximates the posterior distribution when the number of particles tends to infinity. However, since the computational cost grows with the number of particles, practitioners must choose a specific number of particles in the design of their filters.

In many applications, the observations arrive sequentially, and there is a strict deadline for processing each new observation. Then, one could argue that the best solution in terms of filter performance is to increase the number of particles as much as possible and keep it fixed. Also, in some hardware implementations, the number of particles is a design parameter that cannot be modified during implementation. Nevertheless, in many other applications where resources are scarce or are shared with a dynamical allocation and/or with energy restrictions, one might be interested in adapting the number of particles in a smart way. One would use enough particles to achieve a certain performance requirement but without wasting resources by using many more particles if they do not translate into a significant improvement of the filter performance.

The selection of the number of particles, however, is often a delicate subject because, (1) the performance of the filter (the quality of the approximation) cannot usually be described in advance as a function of the number of particles, and (2) the mismatch between the approximation provided by the filter and the unknown posterior distribution is obviously also unknown. Therefore, although there is a clear trade-off between performance and computational cost, this

V. Elvira is with Télécom Lille (Institut Mines-Télécom) and CRISAl laboratory (France), and with the Department of Signal Theory and Communications, Universidad Carlos III de Madrid (Spain), e-mail: victor.elvira@telecom-lille.fr. J. Míguez is with the Department of Signal Theory and Communications, Universidad Carlos III de Madrid (Spain), e-mail: joaquin.miguez@uc3m.es. P. M. Djurić is with the Department of Electrical and Computer Engineering, Stony Brook University (USA), e-mail: petar.djuric@stonybrook.edu. This work was partially supported by *Ministerio de Economía y Competitividad* of Spain (TEC2013-41718-R OTOSIS, TEC2012-38883-C02-01 COMPREHENSION, and TEC2015-69868-C2-1-R ADVENTURE), the Office of Naval Research Global (N62909-15-1-2011), and the National Science Foundation (CCF-1320626 and CCF-1618999).

relation is not straightforward; e.g., increasing the number of particles over a certain value may not significantly improve the quality of the approximation while decreasing the number of particles below some other value can dramatically affect the performance of the filter.

Few papers in the wide literature have addressed the problem of online assessment of the filter convergence for the purpose of adapting the number of particles. In [14], the number of particles is selected so that a bound on the approximation error does not exceed a threshold with certain probability. The latter error is defined as the Kullback-Leibler divergence (KLD) between the approximate filter distribution and a grid-discretized version of the true one (which is itself a potentially-costly approximation with an unknown error). In [15], an adaptation of the number of particles is proposed, based on the KLD approach of [14] and an estimate of the variance of the estimators computed via the particle filter, along with an improvement of the proposal distributions. In [16], the adaptation of the number of particles is based on the effective sample size. These methods are heuristic: they do not enjoy any theoretical guarantees (in the assessment of the approximation errors made by the particle filter) and the allocation of particles, therefore, cannot be ensured to be optimal according to any probabilistic criterion. Some techniques based on more solid theoretical ground have been proposed, within the applied probability community, during the last few years. We discuss them below.

Two types of unbiased estimators of the variance in the approximation of integrals using a class of particle filters were analyzed in [17] using the Feynman-Kac framework of [18]. As an application of these results, it was suggested to use these estimators to select the number of particles in the filter. In particular, the scheme proposed in [17] is a batch procedure in which a particle filter is run several times over the whole data sequence, with increasing number of particles, until the variance of the integral of interest is found to fall below a prescribed threshold. This approach cannot be used for online assessment, which is the goal of the present paper. Another batch method (thus, also not applicable for online assessment) for particle allocation has been recently proposed in [19], where an ad hoc autoregressive model is fitted to estimate the variance of the estimators produced by the particle filter.

Papers on so-called *alive* particle filters can also be found in the literature [20], [21], [22]. These articles focus on models where the likelihood function can take zero value for some regions of the state space, in such a way that there is the risk that a collection of zero-weight particles are generated if a standard algorithm is employed. To avoid this limitation, alive particle filters are based on sampling schemes where new particles are generated until a prescribed number of them, M , attain non-zero weights. The computational cost of the algorithm per time step is, therefore, random. Moreover, the number M is chosen a priori and there is no assessment of whether M allows for reaching adequate accuracy of the estimators (the methodology proposed in the present manuscript can be directly applied to alive particle filters in order to adapt M).

In order to guarantee that the particle set yields a sufficiently

good representation, in [23] it is proposed to test whether the particle estimate of the predictive density of the observation at time t given the previous data is sufficiently large, i.e., whether it is above a prescribed (heuristically chosen) threshold. When the particle set does not satisfy this condition, it is discarded and a new collection of particles is generated. The number of particles is not adapted, since all generated sets have the same size. The computational cost of this algorithm is random.

Finally, in [24, Chapter 4] it is proposed to use the coefficient of variation of the weights (or, equivalently, the effective sample size) in order to detect those observations for which there is a large χ^2 -divergence between the proposal distribution used to generate the set of particles and the target distribution. This connection is rigorously established in [24]. The algorithm, however, is computationally costly compared to classical methods: at each time step, a complete set of particles and weights are generated, and the coefficient of variation is computed. If this coefficient is too high, the particles are discarded, the algorithm “rolls back,” and a new, larger set of particles is generated for better representation of the target distribution (this step is termed “refuelling” in [24]). Although the algorithm enjoys theoretical guarantees, it relies on keeping the particle approximation “locked” to the target distribution at all times. It is known that, once the particle filter has lost track of the state distribution, the effective sample size (and, hence, coefficient of variation) becomes uninformative [25] and, therefore, the link with the χ^2 -divergence is lost.

B. Contributions

We introduce a model-independent methodology for the online assessment of the convergence of particle filters and carry out a rigorous analysis that ensures the consistency of the proposed scheme under fairly standard assumptions. The method is an extension of our previous work presented in [26]. In the proposed scheme, the observations are processed one at a time and the filter performance is assessed by measuring the discrepancy between the actual observation at each time step and a number of fictitious data-points drawn from the particle approximation of the predictive probability distribution of the observations. The method can be exploited to adjust the number of particles dynamically when the performance of the filter degrades below a certain desired level. This would allow a practitioner to select the operation point by considering performance-computational cost tradeoffs. Based on the method, we propose a simple and efficient algorithm that adjusts the number of particles in real time. We demonstrate the performance of the algorithm numerically by running it for a stochastic version of the 3-dimensional Lorenz 63 system. As already noted, this paper builds on the method from [26]. However, the main difference here is that the underlying model is not questioned – instead, it is *assumed* to be correct. The connection between [26] and the present work is that they both build upon the ability to compute predictive statistics of the upcoming observations that turn out to be independent of the underlying state space model. In this paper we have rigorous theoretical results regarding the particle approximations of the predictive distribution of the

observations (while this issue was ignored in [26]). Finally, we suggest practical schemes for the online adjustment of the number of particles.

Let us point out that the adaptive procedure for the online selection of the number of particles described herein is only one of many that can exploit the results of the convergence analysis. In other words, our analysis opens the door for development of new families of algorithms for online adaptation of the number of particles by way of online convergence assessment.

C. Organization of the paper

The rest of the paper is organized as follows. In Section II we describe the class of state space Markov models and provide a basic background on the well-known bootstrap particle filter of [11]. The theoretical results that enable the online assessment of particle filters are stated in Section III, with full details and proofs contained in the Supplementary Material. The proposed methodology for online convergence assessment of the particle filter is introduced in Section IV. Furthermore, this section provides a simple algorithm for the dynamic, online adaptation of the number of particles. In Section V, we illustrate the validity of the method by means of computer simulations for a stochastic Lorenz 63 model. Finally, Section VI contains a summary of results and some concluding remarks.

II. PARTICLE FILTERING

In this section we describe the class of state space models of interest and then present the standard particle filter (PF), which is the basic building block for the methods to be introduced later.

A. State space models and stochastic filtering

Let us consider discrete-time, Markov dynamic systems in state-space form described by the triplet¹

$$\mathbf{X}_0 \sim p(\mathbf{x}_0), \quad (1)$$

$$\mathbf{X}_t \sim p(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad (2)$$

$$\mathbf{Y}_t \sim p(\mathbf{y}_t|\mathbf{x}_t), \quad (3)$$

where

- $t \in \mathbb{N}$ denotes discrete time;
- \mathbf{X}_t is the $d_x \times 1$ -dimensional (random) system state at time t , which takes values in the set $\mathcal{X} \subseteq \mathbb{R}^{d_x}$,
- $p(\mathbf{x}_0)$ is the a priori pdf of the state, while
- $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ denotes the conditional density of the state \mathbf{X}_t given $\mathbf{X}_{t-1} = \mathbf{x}_{t-1}$;

¹In most of the paper we abide by a simplified notation where $p(x)$ denotes the probability density function (pdf) of the random variable X . This notation is argument-wise, hence if we have two random variables X and Y , then $p(x)$ and $p(y)$ denote the corresponding density functions, possibly different; $p(x, y)$ denotes the joint pdf and $p(x|y)$ is the conditional pdf of X given $Y = y$. A more accurate notation, which avoids ambiguities, is used for the analysis and the statement of the theoretical results. Besides, vectors are denoted by bold-face letters, e.g., \mathbf{x} , while regular-face is used for scalars, e.g., x .

- \mathbf{Y}_t is the $d_y \times 1$ -dimensional observation vector at time t , which takes values in the set $\mathcal{Y} \subseteq \mathbb{R}^{d_y}$ and is assumed to be conditionally independent of all other observations given the state \mathbf{X}_t ,
- $p(\mathbf{y}_t|\mathbf{x}_t)$ is the conditional pdf of \mathbf{Y}_t given $\mathbf{X}_t = \mathbf{x}_t$. It is often referred to as the *likelihood* of \mathbf{x}_t , when it is viewed as a function of \mathbf{x}_t given \mathbf{y}_t .

The model described by Eqs. (1)–(3) includes a broad class of systems, both linear and nonlinear, with Gaussian or non-Gaussian perturbations. Here we focus on the case where all the model parameters are known. However, the proposed method can also be used for models with unknown parameters for which suitable particle filtering methods are available [27], [28], [29]. We assume that the prior distribution of the state $p(\mathbf{x}_0)$ is also known.

The stochastic filtering problem consists in the computation of the sequence of posterior probability distributions given by the so-called *filtering* densities $p(\mathbf{x}_t|\mathbf{y}_{1:t})$, $t = 1, 2, \dots$. The pdf $p(\mathbf{x}_t|\mathbf{y}_{1:t})$ is closely related to the one-step-ahead predictive state density $p(\mathbf{x}_t|\mathbf{y}_{1:t-1})$, which is of major interest in many applications and can be written down by way of the Chapman-Kolmogorov equation,

$$p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) = \int p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1})d\mathbf{x}_{t-1}. \quad (4)$$

Using Bayes' theorem together with Eq. (4), we obtain the well-known recursive factorization of the filtering pdf

$$p(\mathbf{x}_t|\mathbf{y}_{1:t}) \propto p(\mathbf{y}_t|\mathbf{x}_t) \int p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1})d\mathbf{x}_{t-1}.$$

For conciseness and notational accuracy, we use the measure-theoretic notation

$$\pi_t(d\mathbf{x}_t) := p(\mathbf{x}_t|\mathbf{y}_{1:t})d\mathbf{x}_t, \quad \xi_t(d\mathbf{x}_t) := p(\mathbf{x}_t|\mathbf{y}_{1:t-1})d\mathbf{x}_t$$

to represent the filtering and the predictive posterior probability distributions of the state, respectively. Note that π_t and ξ_t are probability measures, hence, given a Borel set $A \subset \mathcal{X}$, $\pi_t(A) = \int_A \pi_t(d\mathbf{x}_t)$ and $\xi_t(A) = \int_A \xi_t(d\mathbf{x}_t)$ denote the posterior probability of the event $\mathbf{X}_t \in A$ conditional on $\mathbf{Y}_{1:t} = \mathbf{y}_{1:t}$ and $\mathbf{Y}_{1:t-1} = \mathbf{y}_{1:t-1}$, respectively.

However, the object of main interest for the convergence assessment method to be introduced in this paper is the predictive pdf of the observations, namely the function $p(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ and the associated probability measure

$$\mu_t(d\mathbf{y}_t) := p(\mathbf{y}_t|\mathbf{y}_{1:t-1})d\mathbf{y}_t.$$

The density $p(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ is the normalization constant of the filtering density $p(\mathbf{x}_t|\mathbf{y}_{1:t})$, and it is related to the predictive state pdf $p(\mathbf{x}_t|\mathbf{y}_{1:t-1})$ through the integral

$$p(\mathbf{y}_t|\mathbf{y}_{1:t-1}) = \int p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t-1})d\mathbf{x}_t. \quad (5)$$

It also plays a key role in model assessment [26] and model inference problems [28], [29], [30], [31].

B. The standard particle filter

A PF is an algorithm that processes the observations $\{\mathbf{y}_t\}_{t \geq 1}$ sequentially in order to compute Monte Carlo approximations of the sequence of probability measures $\{\pi_t\}_{t \geq 1}$. The simplest algorithm is the so-called *bootstrap particle filter* (BPF) [11] (see also [32]), which consists of a recursive importance sampling procedure and a resampling step. The term “particle” refers to a Monte Carlo sample in the state space \mathcal{X} , which is assigned an importance weight. Below, we outline the BPF algorithm with M particles.

Algorithm 1. Bootstrap particle filter.

- 1) Initialization. At time $t = 0$, draw M i.i.d. samples, $\mathbf{x}_0^{(m)}$, $m = 1, \dots, M$, from the prior $p(\mathbf{x}_0)$.
- 2) Recursive step. Let $\{\mathbf{x}_{t-1}^{(m)}\}_{m=1}^M$ be the particles at time $t - 1$. At time t , proceed with the two steps below.
 - a) For $m = 1, \dots, M$, draw $\bar{\mathbf{x}}_t^{(m)}$ from the model transition pdf $p(\mathbf{x}_t | \mathbf{x}_{t-1}^{(m)})$. Then compute the normalized importance weights

$$w_t^{(m)} = \frac{p(\mathbf{y}_t | \bar{\mathbf{x}}_t^{(m)})}{\sum_{k=1}^M p(\mathbf{y}_t | \bar{\mathbf{x}}_t^{(k)})}, \quad m = 1, \dots, M. \quad (6)$$
 - b) Resample M times with replacement: for $m = 1, \dots, M$, let $\mathbf{x}_t^{(m)} = \bar{\mathbf{x}}_t^{(k)}$ with probability $w_t^{(k)}$, where $k \in \{1, \dots, M\}$.

For the sake of simplicity, in step 2.(b) above we assume that multinomial resampling [7] is carried out for every $t \geq 1$. The results and methods to be presented in subsequent sections remain valid when resampling is carried out periodically and/or using alternative schemes such as residual [6], stratified [33] or minimum-variance [34] resampling (see also [35]).

The simple BPF yields several useful approximations. After sampling at step 2.(a), the predictive state probability measure ξ_t can be approximated as

$$\xi_t^M(d\mathbf{x}_t) = \frac{1}{M} \sum_{m=1}^M \delta_{\bar{\mathbf{x}}_t^{(m)}}(d\mathbf{x}_t),$$

where $\delta_{\mathbf{x}}$ denotes the Dirac delta measure located at $\mathbf{x} \in \mathcal{X}$. The filter measure π_t can be similarly approximated, either using the particles and weights computed at step 2.(a) or the resampled particles after step 2.(b), i.e.,

$$\bar{\pi}_t^M = \sum_{m=1}^M w_t^{(m)} \delta_{\bar{\mathbf{x}}_t^{(m)}} \quad \text{and} \quad \pi_t^M = \frac{1}{M} \sum_{m=1}^M \delta_{\mathbf{x}_t^{(m)}},$$

respectively. In addition, the BPF yields natural approximations of the predictive pdf's of \mathbf{X}_t and \mathbf{Y}_t given the earlier observations $\mathbf{Y}_{1:t-1} = \mathbf{y}_{1:t-1}$. If we specifically denote these functions as $\tilde{p}_t(\mathbf{x}_t) := p(\mathbf{x}_t | \mathbf{y}_{1:t-1})$ and $p_t(\mathbf{y}_t) := p(\mathbf{y}_t | \mathbf{y}_{1:t-1})$, then we readily obtain their respective estimates as mixture distributions with M mixands, or,

$$\tilde{p}_t^M(\mathbf{x}_t) := \sum_{m=1}^M w_{t-1}^M p(\mathbf{x}_t | \mathbf{x}_{t-1}^{(m)}), \quad \text{and}$$

$$p_t^M(\mathbf{y}_t) := \frac{1}{M} \sum_{m=1}^M p(\mathbf{y}_t | \bar{\mathbf{x}}_t^{(m)}),$$

for any $\mathbf{x}_t \in \mathcal{X}$ and $\mathbf{y}_t \in \mathcal{Y}$.

III. A NOVEL ASYMPTOTIC CONVERGENCE RESULT

The convergence of the approximate measures, e.g., ξ_t^M , towards the true ones is usually assessed in terms of the estimates of 1-dimensional statistics of the corresponding probability distribution. To be specific, let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a real integrable function in the state space and denote²

$$(f, \xi_t) := \int f(\mathbf{x}_t) \xi_t(d\mathbf{x}_t).$$

Under mild assumptions on the state space model, it can be proved that

$$\lim_{M \rightarrow \infty} (f, \xi_t^M) = \lim_{M \rightarrow \infty} \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}_t^{(m)}) = (f, \xi_t) \quad (7)$$

almost surely (a.s.) [18], [12].

According to (5), the predictive observation pdf $p_t(\mathbf{y}_t)$ is an integral w.r.t. ξ_t and, as a consequence, Eq. (7) implies that $\lim_{M \rightarrow \infty} p_t^M(\mathbf{y}) = p_t(\mathbf{y})$ a.s. and point-wise for every $\mathbf{y} \in \mathcal{Y}$ under mild assumptions [18]. However, existing theoretical results *do not* ensure that $p_t^M(\mathbf{y})$ can converge *uniformly* on \mathcal{Y} towards $p_t(\mathbf{y})$ and this fact prevents us from claiming that $\lim_{M \rightarrow \infty} \int h(\mathbf{y}) p_t^M(\mathbf{y}) d\mathbf{y} = \int h(\mathbf{y}) p_t(\mathbf{y}) d\mathbf{y} = (h, \mu_t)$ in some proper sense for integrable real functions $h(\mathbf{y})$.

Important contributions of this paper are (a) the proof of a.s. convergence of the random probability measure

$$\mu_t^M(d\mathbf{y}) := p_t^M(\mathbf{y}) d\mathbf{y}$$

towards μ_t (as $M \rightarrow \infty$) under mild regularity assumptions on the state space model, and (b) the provision of explicit error rates. We point out that μ_t^M is not a classical point-mass Monte Carlo approximation of μ_t (as, for example, π_t^M is an approximation of π_t). Instead, the measure μ_t^M is absolutely continuous with respect to the Lebesgue measure (the same as μ_t itself). If a different reference measure were used to define the pdf's $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ and $p(\mathbf{y}_t | \mathbf{x}_t)$, say ν , then both μ_t and μ_t^M would be absolutely continuous with respect to ν . In order to describe how μ_t^M converges to μ_t in a rigorous manner, we need to introduce some notation:

- For each $t \geq 1$, let us define the function $g_t(\mathbf{y}_t, \mathbf{x}_t) := p(\mathbf{y}_t | \mathbf{x}_t)$, i.e., the conditional pdf of \mathbf{y}_t given \mathbf{x}_t . When this function is used as a likelihood, we write $g_t^{\mathbf{y}_t}(\mathbf{x}_t) := g_t(\mathbf{y}_t, \mathbf{x}_t)$ to emphasize that it is a function of \mathbf{x}_t .
- Let $f : \mathcal{Z} \rightarrow \mathbb{R}$ be a real function on some set \mathcal{Z} . We denote the absolute supremum of f as $\|f\|_\infty := \sup_{z \in \mathcal{Z}} |f(z)|$. The set of bounded real functions on \mathcal{Z} is $B(\mathcal{Z}) := \{f : \mathcal{Z} \rightarrow \mathbb{R} \text{ such that } \|f\|_\infty < \infty\}$.
- Let $\mathbf{a} = (a_1, \dots, a_d)$ be a multi-index, where each a_i , $i = 1, 2, \dots, d$, is a non-negative integer. Let $f : \mathcal{Z} \rightarrow \mathbb{R}$ be a real function on a d -dimensional set $\mathcal{Z} \subseteq \mathbb{R}^d$. We use $D^{\mathbf{a}} f(z)$ to denote the partial derivative of f w.r.t. the variable \mathbf{z} determined by the entries of \mathbf{a} , namely,

$$D^{\mathbf{a}} f(z) = \frac{\partial^{a_1} \dots \partial^{a_d} f}{\partial z_1^{a_1} \dots \partial z_d^{a_d}}(z).$$

²Let $(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ be a measurable space, where $\mathcal{Z} \subseteq \mathbb{R}^d$ for some integer $d \geq 1$ and $\mathcal{B}(\mathcal{Z})$ is the Borel σ -algebra of subsets of \mathcal{Z} . If α is a measure on $\mathcal{B}(\mathcal{Z})$ and the function $h : \mathcal{Z} \rightarrow \mathbb{R}$ is integrable with respect to (w.r.t.) α , then we use the shorthand notation $(f, \alpha) := \int f(z) \alpha(dz)$.

The order of the derivative operator D^a is $|a| = \sum_{i=1}^d a_i$.

- The minimum out of two scalar quantities, $a, b \in \mathbb{R}$, is denoted $a \wedge b$.

We make the following assumptions on the likelihood function g_t and the predictive observation measure $\mu_t(dy_t) = p_t(y_t)dy_t$.

- (\mathfrak{L}) For each $t \geq 1$, the function g_t is positive and bounded, i.e., $g_t(\mathbf{y}, \mathbf{x}) > 0$ for any $(\mathbf{y}, \mathbf{x}) \in \mathcal{Y} \times \mathcal{X}$ and $\|g_t\|_\infty = \sup_{(\mathbf{y}, \mathbf{x}) \in \mathcal{Y} \times \mathcal{X}} |g_t(\mathbf{y}, \mathbf{x})| < \infty$.
- (\mathfrak{D}) For each $t \geq 1$, the function $g_t(\mathbf{y}, \mathbf{x})$ is differentiable with respect to \mathbf{y} , with bounded derivatives up to order d_y , hence $D^1 g_t(\mathbf{y}, \mathbf{x}) = \frac{\partial^{d_y} g_t}{\partial y_1 \dots \partial y_{d_y}}(\mathbf{y}, \mathbf{x})$ exists and

$$\|D^1 g_t\|_\infty = \sup_{(\mathbf{y}, \mathbf{x}) \in \mathcal{Y} \times \mathcal{X}} |D^1 g_t(\mathbf{y}, \mathbf{x})| < \infty.$$

- (\mathfrak{C}) For any $0 < \beta < 1$ and any $p \geq 4$, the sequence of hypercubes

$$C_M := \left[-\frac{M^{\frac{\beta}{p}}}{2}, +\frac{M^{\frac{\beta}{p}}}{2}\right] \times \dots \times \left[-\frac{M^{\frac{\beta}{p}}}{2}, +\frac{M^{\frac{\beta}{p}}}{2}\right] \subset \mathbb{R}^{d_y}$$

satisfies the inequality $\mu_t(\overline{C_M}) \leq bM^{-\eta}$ for some constants $b > 0$ and $\eta > 0$ independent of M (yet possibly dependent on β and p), where $\overline{C_M} = \mathbb{R}^{d_y} \setminus C_M$ is the complement of C_M .

Remark 1. Assumptions (\mathfrak{L}) and (\mathfrak{D}) refer to regularity conditions (differentiability and boundedness) that the likelihood function of the state space model should satisfy. Models of observations, for example, of the form $\mathbf{y}_t = f(\mathbf{x}_t) + \mathbf{u}_t$, where f is a (possibly nonlinear) transformation of the state \mathbf{x}_t and \mathbf{u}_t is noise with some differentiable, exponential-type pdf (e.g., Gaussian or mixture-Gaussian), readily satisfy these assumptions. Typical two-sided heavy-tailed distributions, such as Student's t distribution, also satisfy (\mathfrak{L}) and (\mathfrak{D}).

Remark 2. Assumption (\mathfrak{C}) states an explicit bound on the probability under the tails of the pdf $p_t(y_t) = p(\mathbf{y}_t | \mathbf{y}_{1:t-1})$. The bound is polynomial, namely

$$\mu_t(\overline{C_M}) = 1 - \int_{-\frac{1}{2}M^{\frac{\beta}{p}}}^{\frac{1}{2}M^{\frac{\beta}{p}}} \dots \int_{-\frac{1}{2}M^{\frac{\beta}{p}}}^{\frac{1}{2}M^{\frac{\beta}{p}}} p_t(\mathbf{y}) d\mathbf{y} \leq bM^{-\eta},$$

and therefore immediately verified, e.g., by all distributions of the exponential family as well as for many heavy-tailed distributions. For example, when $d_y = 1$ (i.e., the observations are scalars), one can choose the constants b and η such that $bM^{-\eta}$ is an upper bound for the tails of the (heavy-tailed) Pareto, Weibull, Burr or Levy distributions.

It is actually possible to find simple conditions on the conditional pdf of the observations, $g_t(\mathbf{y}_t, \mathbf{x}_t)$, that turn out sufficient for assumption (\mathfrak{C}) to hold true. Let us keep $d_y = 1$, for simplicity, and assume that there exists a sequence of positive constants $\{c_t\}_{t \geq 1}$ such that $g_t(y_t, \mathbf{x}_t)$ has a polynomial upper bound itself, namely

$$\sup_{\mathbf{x}_t \in \mathcal{X}} g_t(y_t, \mathbf{x}_t) \leq c_t |y_t|^{-(1+\epsilon)} \quad (8)$$

for some $\epsilon > 0$ and every y_t such that $|y_t| > \frac{1}{2}$ (note that the smallest set in the sequence C_M is $C_1 = [-\frac{1}{2}, \frac{1}{2}]$). For probability distributions with infinite support and continuous with respect to the Lebesgue measure, the inequality (8) implies that the densities $g_t(y_t, \mathbf{x}_t)$ are integrable for every possible choice of $\mathbf{x}_t \in \mathcal{X}$. Then, the probability below the right tail of $p_t(y)$ is

$$\begin{aligned} \int_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} p_t(y) dy &= \int_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} \int_{\mathcal{X}} g_t(y, \mathbf{x}) \tilde{p}_t(\mathbf{x}) d\mathbf{x} dy \\ &\leq c_t \int_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} y^{-(1+\epsilon)} \int_{\mathcal{X}} \tilde{p}_t(\mathbf{x}) d\mathbf{x} dy, \end{aligned}$$

where the inequality follows from the application of (8). Since $\tilde{p}_t(\mathbf{x})$ is a pdf, we have $\int_{\mathcal{X}} \tilde{p}_t(\mathbf{x}) d\mathbf{x} = 1$ and some elementary calculations yield

$$\begin{aligned} \int_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} p_t(y) dy &\leq c_t \int_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} y^{-(1+\epsilon)} dy \\ &= c_t \left[-\frac{y^{-\epsilon}}{\epsilon} \right]_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} = \frac{2^\epsilon c_t}{\epsilon} M^{-\frac{\epsilon\beta}{p}}. \quad (9) \end{aligned}$$

The same result is easily obtained for the left tail of $p_t(y)$, hence

$$\begin{aligned} \mu_t(\overline{C_M}) &= \int_{\frac{1}{2}M^{\frac{\beta}{p}}}^{\infty} p_t(y) dy + \int_{-\infty}^{-\frac{1}{2}M^{\frac{\beta}{p}}} p_t(y) dy \\ &\leq \frac{2^{1+\epsilon} c_t}{\epsilon} M^{-\frac{\epsilon\beta}{p}}. \quad (10) \end{aligned}$$

By comparing (10) and the inequality $\mu_t(\overline{C_M}) \leq bM^{-\eta}$, we readily see that we can choose $b = \frac{2^{1+\epsilon} c_t}{\epsilon}$ and $\eta = \frac{\epsilon\beta}{p} > 0$ to uphold assumption (\mathfrak{C}). A similar derivation can be carried out when $d_y > 1$.

Theorem 1. Assume that (\mathfrak{L}), (\mathfrak{D}) and (\mathfrak{C}) hold and the observations $\mathbf{y}_{1:t-1}$ are fixed (and otherwise arbitrary). Then, for every $h \in B(\mathcal{Y})$ and any $\epsilon \in (0, \frac{1}{2})$ there exists an a.s. finite r.v. W_t^ϵ , independent of M , such that

$$|(h, \mu_t^M) - (h, \mu_t)| \leq \frac{W_t^\epsilon}{M^{(\frac{1}{2}-\epsilon)\wedge\eta}}.$$

In particular,

$$\lim_{M \rightarrow \infty} (h, \mu_t^M) = (h, \mu_t) \quad \text{a.s.}$$

See the Supplementary Material for a proof.

Note that the r.v. W_t^ϵ in the statement of Theorem 1 depends on the time instant t . It is possible to remove this dependence if the constants b and η in assumption (\mathfrak{C}) are chosen to be independent of t and we impose further constraints on the likelihood function and the Markov kernel of the state space model (similar to the sufficient conditions for uniform convergence in, e.g., [18] or [36]).

IV. ONLINE SELECTION OF THE NUMBER OF PARTICLES

In the sequel we assume scalar observations, hence $d_y = 1$ and $\mathbf{y}_t = y_t$ (while $d_x \geq 1$ is arbitrary). A discussion of how to proceed when $d_y > 1$ is provided in Section IV-E.

Our goal is to evaluate the convergence of the BPF (namely, the accuracy of the approximation $p_t^M(y_t)$) in real time and, based on the convergence assessment, adapt the computational effort of the algorithm, i.e., the number of used particles M .

To that end, we run the BPF in the usual way with a light addition of computations. At each iteration we generate K “fictitious observations”, denoted $\tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$, from the approximate predictive pdf $p_t^M(y_t)$. If the BPF is operating with a small enough level of error, then Theorem 1 states that these fictitious observations come approximately from the same distribution as the acquired observation, i.e., $\mu_t^M(dy_t) \approx \mu_t(dy_t)$. In that case, as we explain in Subsection IV-B, a statistic $a_{K,t}^K$ can be constructed using $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$, which necessarily has an (approximately) uniform distribution independently of the specific form of the state-space model (1)–(3). By collecting a sequence of such statistics, say $a_{t-W+1}^K, \dots, a_t^K$ for some window size W , one can easily test whether their empirical distribution is close to uniform using standard procedures. The better the approximation $\mu_t^M \approx \mu_t$ generated by the BPF, the better fit with the uniform distribution can be expected.

If $K \ll M$ and W is not too large, the cost of the added computations is negligible compared to the cost of running the BPF with M particles and, as we numerically show in Section V, the ability to adapt the number of particles online leads to a very significant reduction of the running times without compromising the estimation accuracy.

Below we describe the method, justify its theoretical validity and discuss its computational complexity as well as its extension to the case of multidimensional y_t ’s.

A. Generation of fictitious observations

The proposed method demands at each time t the generation of K fictitious observations (i.e., Monte Carlo samples), denoted $\{\tilde{y}_t^{(k)}\}_{k=1}^K$, from the approximate predictive observation pdf $p_t^M(y_t) = \frac{1}{M} \sum_{m=1}^M p(y_t | \bar{x}_t^{(m)})$. Since the latter density is a finite mixture, drawing from $p_t^M(y_t)$ is straightforward as long as the conditional density of the observations, $p(y_t | \bar{x}_t)$, is itself amenable to sampling. In order to generate $\tilde{y}_t^{(k)}$, it is enough to draw a sample $j^{(k)}$ from the discrete uniform distribution on $\{1, 2, \dots, M\}$ and then generate $\tilde{y}_t^{(k)} \sim p(y_t | \bar{x}_t^{(j^{(k)})})$.

B. Assessing convergence via invariant statistics

For simplicity, let us assume first that $p_t^M(y_t) = p_t(y_t) = p(y_t | y_{1:t-1})$, i.e., there is no approximation error and, therefore, the fictitious observations $\{\tilde{y}_t^{(k)}\}_{k=1}^K$ have the same distribution as the true observation y_t . We define the set $\mathcal{A}_{K,t} := \{y \in \{\tilde{y}_t^{(k)}\}_{k=1}^K : y < y_t\}$ and the r.v. $A_{K,t} := |\mathcal{A}_{K,t}| \in \{0, 1, \dots, K\}$. Note that $\mathcal{A}_{K,t}$ is the set of fictitious observations which are smaller than the actual one, while $A_{K,t}$ is the number of such observations. If we let \mathbb{Q}_K denote the probability mass function (pmf) of A_K , it is not hard to show that \mathbb{Q}_K is uniform independently of the value and distribution of y_t . This is rigorously given by the Proposition below.

Proposition 1. *If $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ are i.i.d. samples from a common continuous (but otherwise arbitrary) probability distribution, then the pmf of the r.v. $A_{K,t}$ is*

$$\mathbb{Q}_K(n) = \frac{1}{K+1}, \quad n = 0, \dots, K. \quad (11)$$

Proof: Since $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ are i.i.d., all possible orderings of the $K+1$ samples are a priori equally probable, and the value of the r.v. $A_{K,t}$ depends uniquely on the relative position of y_t after the samples are sorted (e.g., if y_t is the smallest sample, then $A_{K,t} = 0$, if there is exactly one $\tilde{y}_t^{(i)} < y_t$ then $A_{K,t} = 1$, etc.). There are $(K+1)!$ different ways in which the samples $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ can be ordered, but $A_{K,t}$ can only take values from 0 to K . In particular, given the relative position of y_t , there are $K!$ different ways in which the remaining samples $\tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ can be arranged. Therefore, $\mathbb{Q}_K(A_K = n) = \frac{K!}{(K+1)!} = \frac{1}{K+1}$ for every $n \in \{0, 1, \dots, K\}$. \square

For the case of interest in this paper, the r.v.’s $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ (the actual and fictitious observations) have a common probability distribution given by the measure μ_t and are generated independently. For the class of state space models described in Section II, and the explicit assumptions in Section III, the measure μ_t is absolutely continuous w.r.t. the Lebesgue measure (with associated density $p_t(y)$) and, therefore, $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ are indeed continuous r.v.’s and the assumptions of Proposition 1 are met. Moreover, it can also be proved that the variables in the sequence $A_{K,t}$ are independent.

Proposition 2. *If the r.v.’s $y_t, \tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ are i.i.d. with common pdf $p_t(y)$, then the r.v.’s in the sequence $\{A_{K,t}\}_{t \geq 1}$ are independent.*

See Appendix A for a proof.

In practice, $p_t^M(y_t)$ is just an approximation of the predictive observation pdf $p_t(y_t)$ and, therefore, the actual and fictitious observations are not i.i.d. However, under the assumptions of Theorem 1, the a.s. convergence of the approximate measure $\mu_t^M(dy_t) = p_t^M(y_t)dy_t$ enables us to obtain an “approximate version” of the uniform distribution in Proposition 1, with the error vanishing as $M \rightarrow \infty$. To be specific, we introduce the set $\mathcal{A}_{K,M,t} := \{y \in \{\tilde{y}_t^{(k)}\}_{k=1}^K : y < y_t\}$, which depends on M because of the mismatch between $p_t^M(y_t)$ and $p_t(y_t)$, and the associated r.v. $A_{K,M,t} = |\mathcal{A}_{K,M,t}|$ with pmf $\mathbb{Q}_{K,M,t}$. We have the following convergence result for $\mathbb{Q}_{K,M,t}$.

Theorem 2. *Let y_t be a sample from $p_t(y_t)$ and let $\{\tilde{y}_t^{(k)}\}_{k=1}^K$ be i.i.d. samples from $p_t^M(y_t)$. If the observations $y_{1:t-1}$ are fixed and Assumptions (L), (D) and (E) hold, then there exists a sequence of non-negative r.v.’s $\{\varepsilon_t^M\}_{M \in \mathbb{N}}$ such that $\lim_{M \rightarrow \infty} \varepsilon_t^M = 0$ a.s. and*

$$\frac{1}{K+1} - \varepsilon_t^M \leq \mathbb{Q}_{K,M,t}(n) \leq \frac{1}{K+1} + \varepsilon_t^M. \quad (12)$$

In particular, $\lim_{M \rightarrow \infty} \mathbb{Q}_{K,M,t}(n) = \mathbb{Q}_K(n) = \frac{1}{K+1}$ a.s.

See Appendix B for a proof. Proposition 1 states that the statistic $A_{K,t}$ is distribution-invariant, since $\mathbb{Q}_K(n) = \frac{1}{K+1}$

independently of t and the state space model. Similarly, Theorem 2 implies that the statistic $A_{K,M,t}$ is asymptotically distribution-invariant (independently of t and the model) since $\mathbb{Q}_{K,M,t}(n) \rightarrow \frac{1}{K+1}$ when $M \rightarrow \infty$, as the BPF converges.³

C. Algorithm with adaptive number of particles

We propose an algorithm that dynamically adjusts the number of particles of the filter based on the transformed r.v. $A_{K,M,t}$. Table II summarizes the proposed algorithm, that is embedded into a standard BPF (see Section II-B) but can be applied to virtually any other particle filter in a straightforward manner. The parameters of the algorithm are shown in Table I.

The BPF is initialized in Step 1(a) with M_0 initial particles. At each recursion, in Step 2(a), the filtered distribution of the current state is approximated. In Step 2(b), K fictitious observations $\{\tilde{y}_t^{(k)}\}_{k=1}^K$ are drawn and the statistic $A_{K,M,t} = a_{K,M,t}$ is computed. In Step 2(c), once a set of W consecutive statistics have been acquired, $\mathcal{S}_t = \{a_{K,M,t-W+1}, a_{K,M,t-W+2}, \dots, a_{K,M,t-1}, a_{K,M,t}\}$, a statistical test is performed for checking whether \mathcal{S}_t is a sequence of samples from the uniform pmf given by Eq. (11).

There are several approaches that can be used to exploit the information contained in \mathcal{S}_t . Here we perform a Pearson's chi-squared test [37], where the χ_t^2 statistic is computed according to Eq. (13) (see Table II). Then, a p-value $p_{K,t}^*$ for testing the hypothesis that the empirical distribution of \mathcal{S}_t is uniform is computed. The value $p_{K,t}^*$ is obtained by comparing the χ_t^2 statistic with the χ^2 distribution with K degrees of freedom. Intuitively, a large $p_{K,t}^*$ suggests a good match of the sequence \mathcal{S}_t with an i.i.d. sample from the uniform distribution on $\{0, 1, \dots, K\}$, while a small $p_{K,t}^*$ indicates a mismatch. Therefore, the p-value $p_{K,t}^*$ is compared with two different significance levels: a low threshold p_ℓ and a high threshold p_h . If $p_{K,t}^* \leq p_\ell$, the number of particles is increased according to the rule $M_t = f_{\text{up}}(M_{t-1})$ whereas, if $p_{K,t}^* \geq p_h$, the number of particles is decreased according to the rule $M_t = f_{\text{down}}(M_{t-1})$. When $p_\ell < p_{K,t}^* < p_h$, the number of particles remains fixed. These two significance levels allow the practitioner to select the operation range by considering a performance-to-computational-cost tradeoff. Note that we set M_{\min} and M_{\max} , maximum and minimum values for the number of particles, respectively.

A large window W yields a more accurate convergence assessment but increases the latency (or decreases the responsiveness) of the algorithm. If the algorithm must be run online, this latency can be critical for detecting a malfunction of the filter and adapting consequently the number of particles. Therefore there is a tradeoff between the accuracy of the convergence assessment procedure and latency of the algorithm.

D. Computational cost

Compared to the BPF, the additional computational cost of the method is mainly driven by the generation of the K

³Specifically note that, under assumptions (L), (D) and (C), the convergence of the continuous random measure μ_t^M computed via the BPF (which is sufficient to obtain (12); see Appendix B) is guaranteed by Theorem 1.

TABLE I: Parameters of the algorithm

- M_0 , initial number of particles
- M_{\min} , minimum number of particles
- M_{\max} , maximum number of particles
- K , number of fictitious samples per iteration
- W , window length
- p_ℓ , lower significance level of p-values
- p_h , higher significance level of p-values
- $f_{\text{up}}(\cdot)$, rule for increasing M
- $f_{\text{down}}(\cdot)$, rule for decreasing M

TABLE II: Algorithm for adapting the number of particles

1) [Initialization]

- a) Initialize the particles and the weights of the filter as

$$\mathbf{x}_0^{(m)} \sim p(\mathbf{x}_0), \quad m = 1, \dots, M_0,$$

$$w_0^{(m)} = 1/M_0, \quad m = 1, \dots, M_0,$$

and set $n = 1$.

2) [For $t = 1 : T$]

- a) **Bootstrap particle filter:**

- Resample M_n samples of $\bar{\mathbf{x}}_{t-1}^{(m)}$ with weights $w_{t-1}^{(m)}$ to obtain $\mathbf{x}_{t-1}^{(m)}$.
- Propagate $\bar{\mathbf{x}}_t^{(m)} \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}^{(m)})$, $m = 1, \dots, M_n$.
- Compute the non-normalized weights $\bar{w}_t^{(m)} = p(y_t | \bar{\mathbf{x}}_t^{(m)})$, $m = 1, \dots, M_n$.
- Normalize the weights $\bar{w}_t^{(m)}$ to obtain $w_t^{(m)}$, $m = 1, \dots, M_n$.

- b) **Fictitious observations:**

- Draw $\tilde{y}_t^{(k)} \sim p^M(y_t | y_{t-1})$, $k = 1, \dots, K$.
- Compute $a_{K,M,t} = A_{K,M,t}$, i.e., the position of y_t within the set of ordered fictitious observations $\{\tilde{y}_t^{(k)}\}_{k=1}^K$.

- c) If $t = nW$, (**assessment of convergence**):

- Compute the χ_t^2 statistic over the empirical distribution of $\mathcal{S}_t = \{a_{K,M,t}, a_{K,M,t-1}, \dots, a_{K,M,t-W+1}\}$ as

$$\chi_t^2 = \sum_{j=0}^K \frac{(O_j - E_j)^2}{E_j}, \quad (13)$$

where O_j is the frequency of the observations in the window being in the j th relative position, i.e., $O_j = |\{a_{K,M,\tau} = j\}|$, and E_j is the expected frequency under the null hypothesis, i.e., $E_j = W \cdot \mathbb{Q}_K(j) = \frac{W}{K+1}$ (see Eq. (11)).

- Calculate the p-value $p_{K,t}^*$ by comparing the statistic χ_t^2 to the χ^2 -distribution with K degrees of freedom.
- If $p_{K,t}^* \leq p_\ell$
increase $M_n = \min\{f_{\text{up}}(M_{n-1}), M_{\max}\}$.
- Else, if $p_{K,t}^* \geq p_h$,
decrease $M_n = \max\{f_{\text{down}}(M_{n-1}), M_{\min}\}$.
- Else,
 $M_n = M_{n-1}$.
- Set $n = n + 1$.

- d) If $t < Wn$, set $t = t + 1$ and go to 2. Otherwise, end.

fictitious observations at each iteration as shown in Subsection IV-A. The generation of these fictitious observations is a two-step procedure, where in the first step, we draw K discrete indices, say j_1, \dots, j_K , from the set $\{1, \dots, M_n\}$ with uniform probabilities, and in the second step, we draw K samples from $p(y_t | \bar{\mathbf{x}}_t^{(j_1)}), \dots, p(y_t | \bar{\mathbf{x}}_t^{(j_K)})$, respectively.

In the proposed algorithm, a Pearson's χ^2 test is performed with a sequence \mathcal{S}_t of W samples, that is, it is carried out only once every W consecutive time steps. Therefore, the computational cost will depend on the parameters K and W . We will show in Section V that the algorithm can work very well with a low number of fictitious observations, which imposes a very light extra computational load.

E. Multidimensional observations

Through this section, we have assumed scalar observations. In the multidimensional case, with $\mathbf{y}_t = [y_{1,t}, \dots, y_{d_y,t}]^\top$, the same assessment scheme can be applied over each marginal $p(y_{i,t}|\mathbf{y}_{1:t-1})$ of the predictive observation pdf. Theoretical guarantees readily follow from the convergence of the marginal measures $\mu_{i,t}^M(dy_{i,t}) = p^M(y_{i,t}|\mathbf{y}_{1:t-1})dy_{i,t}$ under the same assumptions as the joint measure μ_t^M (see the Supplementary Material).

The algorithm proposed in Section IV-C can be extended to the case with multidimensional observations. One of way of doing it is by performing an independent assessment for each marginal pdf $p(y_{i,t}|\mathbf{y}_{1:t-1})$. As a result, d_y p-values $p_{K,t,i}^*$, with $i = 1, \dots, d_y$, become available for deciding whether to increase, decrease or keep fixed the number of particles. A conservative approach is to increase the number of particles whenever at least one p-value $p_{K,t,i}^*$ is below the threshold p_ℓ . Note that the complexity of this approach grows with the dimension of the observations.

Finally, note that the convergence of the marginals does not imply the convergence of the joint approximation μ_t^M . However, it can be reasonably expected that when all the marginals are approximated well over a period of time, the joint distribution is accurately approximated as well.

V. NUMERICAL EXAMPLE

A. The three-dimensional Lorenz system

1) *Model description:* In this section we show computer simulation results that demonstrate the performance of the proposed method. We consider the problem of tracking the state of a three-dimensional Lorenz system [38] with additive dynamical noise, partial observations and additive measurement noise [39]. Namely, we consider a three-dimensional stochastic process $\{X(s)\}_{s \in (0, \infty)}$ taking values on \mathbb{R}^3 , whose dynamics are described by the system of stochastic differential equations

$$\begin{aligned} dX_1 &= -s(X_1 - Y_1) + dW_1, \\ dX_2 &= rX_1 - X_2 - X_1X_3 + dW_2, \\ dX_3 &= X_1X_2 - bX_3 + dW_3, \end{aligned}$$

where $\{W_i(s)\}_{s \in (0, \infty)}$, $i = 1, 2, 3$, are independent one-dimensional Wiener processes and

$$(s, r, b) = \left(10, 28, \frac{8}{3}\right)$$

are static model parameters broadly used in the literature since they lead to a chaotic behavior [38]. Here we use a discrete-time version of the latter system using an Euler-Maruyama

scheme with integration step $\Delta = 10^{-3}$, which yields the model

$$X_{1,n} = X_{1,n-1} - \Delta s(X_{1,n-1} - X_{2,n-1}) + \sqrt{\Delta}U_{1,n}, \quad (14)$$

$$X_{2,n} = X_{2,n-1} + \Delta(rX_{1,n-1} - X_{2,n-1} - X_{1,n-1}X_{3,n-1}) + \sqrt{\Delta}U_{2,n}, \quad (15)$$

$$X_{3,n} = X_{3,n-1} + \Delta(X_{1,n-1}X_{2,n-1} - bX_{3,n-1}) + \sqrt{\Delta}U_{3,n}, \quad (16)$$

where $\{U_{i,n}\}_{n=0,1,\dots}$, $i = 1, 2, 3$, are independent sequences of i.i.d. normal random variables with zero mean and unit variance. The system (14)–(16) is partially observed every 200 discrete-time steps. Specifically, we collect a sequence of scalar observations $\{Y_t\}_{t=1,2,\dots}$, of the form

$$Y_t = X_{1,200t} + V_t, \quad (17)$$

where the observation noise $\{V_t\}_{t=1,2,\dots}$ is a sequence of i.i.d. normal random variables with zero mean and variance $\sigma^2 = \frac{1}{2}$.

Let $\mathbf{X}_n = (X_{1,n}, X_{2,n}, X_{3,n}) \in \mathbb{R}^3$ be the state vector. The dynamic model given by Eqs. (14)–(16) defines the transition kernel $p(\mathbf{x}_n|\mathbf{x}_{n-1})$ and the observation model of Eq. (17) is the likelihood function

$$p(y_t|x_{1,200t}) \propto \exp \left\{ -\frac{1}{2\sigma^2} (y_t - x_{1,200t})^2 \right\}.$$

The goal is on tracking the sequence of joint posterior probability measures π_t , $t = 1, 2, \dots$, for $\{\hat{\mathbf{X}}_t\}_{t=1,\dots}$, where $\hat{\mathbf{X}}_t = \mathbf{X}_{200t}$. Note that one can draw a sample $\hat{\mathbf{X}}_t = \hat{\mathbf{x}}_t$ conditional on $\hat{\mathbf{X}}_{t-1} = \hat{\mathbf{x}}_{t-1}$ by successively simulating

$$\tilde{\mathbf{x}}_n \sim p(\mathbf{x}_n|\tilde{\mathbf{x}}_{n-1}), \quad n = 200(t-1) + 1, \dots, 200t,$$

where $\tilde{\mathbf{x}}_{200(t-1)} = \hat{\mathbf{x}}_{t-1}$ and $\hat{\mathbf{x}}_t = \tilde{\mathbf{x}}_{200t}$. The prior measure for the state variables is normal, namely

$$\mathbf{X}_0 \sim \mathcal{N}(\mathbf{x}_*, v_0^2 \mathcal{I}_3),$$

where $\mathbf{x}_* = (-5.9165; -5.5233; 24.5723)$ is the mean and $v_0^2 \mathcal{I}_3$ is the covariance matrix of \mathbf{X}_0 , with $v_0^2 = 10$ and \mathcal{I}_3 being the three-dimensional identity matrix.

2) *Simulation setup:* With this example, we aim at showing how the proposed algorithm allows to operate the particle filter with a prescribed performance-to-computational-budget tradeoff. With this purpose, we applied a standard BPF for tracking the sequence of posterior probability measures of the system system (14)–(16) generated by the three-dimensional Lorenz model described above. We generated a sequence of $T = 2000$ synthetic observations, $\{y_t; t = 1, \dots, 2000\}$, spread over an interval of 400 seconds (in continuous time), corresponding to 4×10^5 discrete time steps in the Euler-Maruyama scheme (hence, one observation every 200 steps). Since the time scale of the discrete time approximation of Eqs. (14)–(16) is $n = 200t$, a resampling step is taken every 200 steps of the underlying discrete-time system.

We started running the PF with a sufficiently large number of particles, namely $M = 5000$, and then let the proposed algorithm decrease the number of particles to attain a prescribed point in the performance-to-computation-cost range. This point is controlled by the operation range of the p-value, which is in turn driven by the pair of significance

levels $[p_\ell - p_h]$. We tested the algorithm for different ranges of p-values, namely, $p_\ell \in \{0.5, 0.4, 0.3, 0.2, 0.1, 0.05\}$ and $p_h \in \{0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1\}$. When the p-value is below p_ℓ , the algorithm doubles the number of particles $M_{n+1} = f_{\text{up}}(M_n) = 2M_n$, and when the p-value is over p_h , the number of particles is halved, $M_{n+1} = f_{\text{down}}(M_n) = M_n/2$. We used $K = 7$ fictitious observations and a window of size $W = 20$.

In order to assess the approximation errors, we computed the empirical MSEs of the approximation of the posterior mean, $E[\hat{X}_t | Y_{1:t} = y_{1:t}]$, by averaging the MSEs for the whole sequences. Note that, since the actual expectation cannot be computed in closed form for this system, we used the true underlying sequence $\{X_{200t}\}_{t=1,2,\dots}$ as the ground truth.

3) *Numerical results:* Table III shows results of the MSE of the approximation of the posterior mean, the average number of particles

$$\bar{M} = \frac{2}{T} \sum_{k=\frac{T}{2}+1}^T M_k, \quad (18)$$

the p-values of the χ^2 test, and the Hellinger distance [40] between the empirical distribution of \mathcal{S}_t and the uniform distribution. They were obtained by averaging over 100 runs and averaging over time for each run. The initial number of particles $M_0 = 2^{15}$, and the minimum and maximum number of particles are $M_{\min} = 2^5$ and $M_{\max} = 2^{15}$, respectively. The first half of time steps were discarded for obtaining the displayed results in order to test the behavior of the algorithm for different sets of parameters (see Eq. (18)). Regarding the relation between the MSE and \bar{M} and the p-values, it can be seen that selecting a high operation range yields good performance (low MSE) at the cost of using a large number of particles (high \bar{M}). When we decrease the range of p-values, the algorithm decreases the number of particles, increasing also the approximation error. Table III shows that this conclusion holds for any pair of $[p_\ell - p_h]$.

Figure 1 shows the MSE, the number of particles \bar{M} , and the execution time for the different operation ranges (solid blue line) compared to the particle filter with a fixed number of particles $M = 2^{15}$ (dashed red line). It can be seen that with a moderate operation range ($[p_\ell - p_h] = [0.3 - 0.7]$), the algorithm can perform (in terms of MSE) similarly to the case with fixed M , while reducing the execution time approximately by a factor of four. The execution time can be further reduced by decreasing the operation range, although this worsens the performance.

Figure 2 displays the evolution of the number of particles over time (averaged over 100 runs) for $[p_\ell - p_h] = [0.3 - 0.7]$ both when $M_0 = 5000$ and $M_0 = 10$. In this case, the minimum and maximum number of particles are $M_{\min} = 10$ and $M_{\max} = 5000$, respectively. We see that, after some time, the number of particles adjusted by the algorithm does not depend on M_0 .

Figure 3 shows the same behavior for $[p_\ell - p_h] = [0.2 - 0.6]$. After some time, the filter uses less particles than the filter with results in Fig. 2 because the selected range of thresholds employs smaller p-values.

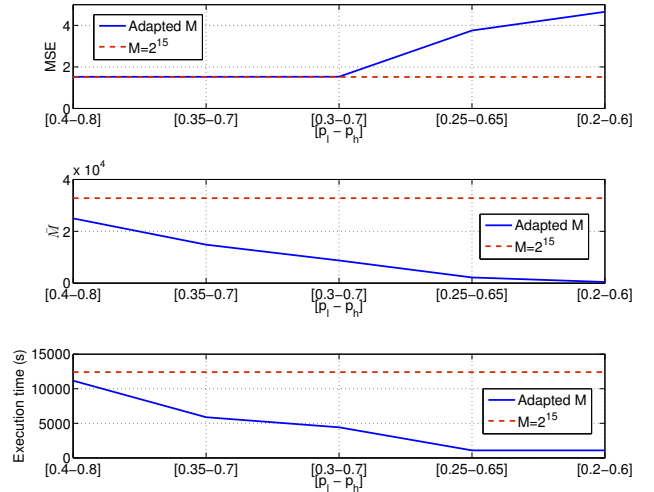


Fig. 1: Lorenz Model (Section V-A). MSE, number of particles \bar{M} and execution time for different pairs of significance levels $[p_\ell - p_h]$ in solid blue line, and with a fixed number of particles $M = 2^{15}$ in dashed red line.

Figure 4 shows histograms of averaged MSE and \bar{M} for simulations performed with two different sets of thresholds: $[p_\ell - p_h] = [0.3 - 0.5]$ and $[p_\ell - p_h] = [0.5 - 0.7]$. In both cases, the initial number of particles is $M_0 = 5000$. It can be seen that a more demanding pair of thresholds ($[p_\ell - p_h] = [0.5 - 0.7]$) leads to better performance and a larger average number of particles. This behavior can also be seen in Figure 5, where the MSE w.r.t. the number of particles is displayed for three different sets of thresholds. Note that a filter with a too relaxed set of thresholds ($[p_\ell - p_h] = [0.05 - 0.4]$) uses very few particles but obtains a poor performance, while a filter with the most stringent set of thresholds ($[p_\ell - p_h] = [0.5 - 0.9]$) consistently yields a low MSE, at the expense of using a larger number of particles.

The numerical results have been computed in a Matlab environment on a computer with an Intel Core i5 processor (2.7 GHz clock frequency) and 12 GB of RAM.

4) *Multidimensional Observations:* Now we consider the case where we have observations also related to the second dimension of the hidden state. In particular, and following the notation of the previous experiment, we collect a sequence of bi-dimensional observations $\{Y_t\}_{t=1,2,\dots}$ with components

$$\begin{aligned} Y_{1,t} &= X_{1,400t} + V_{1,t}, \\ Y_{2,t} &= X_{2,400t} + V_{2,t}, \end{aligned}$$

where the observation noises $\{V_{1,t}\}_{t=1,2,\dots}$ and $\{V_{2,t}\}_{t=1,2,\dots}$ are two sequences of i.i.d. normal random variables with zero mean and variance $\sigma^2 = \frac{1}{2}$. Note that now the state is observed every 400 discrete-time steps in order to make the system more difficult to be tracked.

The implemented algorithm is an extension of the unidimensional case, as suggested in Section IV-E. In particular, we perform the assessment over the marginals

$[p_l - p_h]$	Fixed $M = 2^{15}$	$[0.4 - 0.8]$	$[0.35 - 0.7]$	$[0.3 - 0.7]$	$[0.25 - 0.65]$	$[0.2 - 0.6]$
MSE	1.5193	1.5234	1.5240	1.5287	3.7552	4.6540
\bar{M}	32768	24951	14840	8729	2197	451
p-val	0.5108	0.5089	0.4902	0.4815	0.4872	0.4785
Hell. distance	0.2312	0.2355	0.2493	0.2462	0.2476	0.2521
exec. time (s)	6201	5617	3014	1532	131	67
time ratio	1	1.10	2.1	4.05	47.43	92.36

TABLE III: Lorenz Model (Section V-A): $\Delta = 10^{-3}$, $T_{obs} = 200\Delta$, $\sigma^2 = 0.5$. Algorithm details: $W = 20$, $K = 7$, $M_{\max} = 2^{15}$, $M_{\min} = 2^7$. MSE in the approximation of the posterior mean, averaged number of particles \bar{M} , averaged p-value, and averaged Hellinger distance.

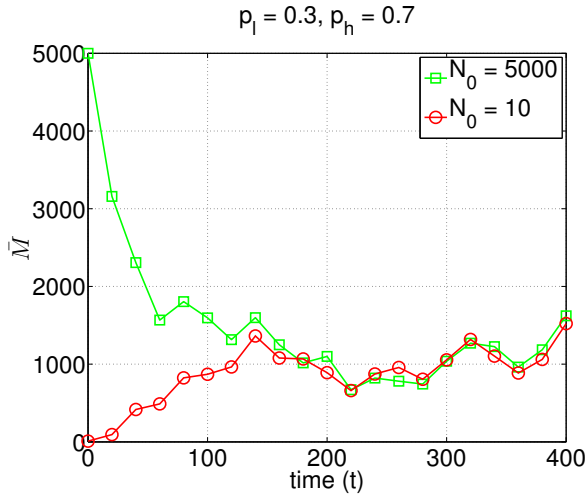


Fig. 2: Lorenz Model (Section V-A). Evolution of the number of particles adapted by the proposed algorithm when the initial number of particles $M_0 \in \{10, 5000\}$. The significance levels were set to $p_\ell = 0.3$ and $p_h = 0.7$.

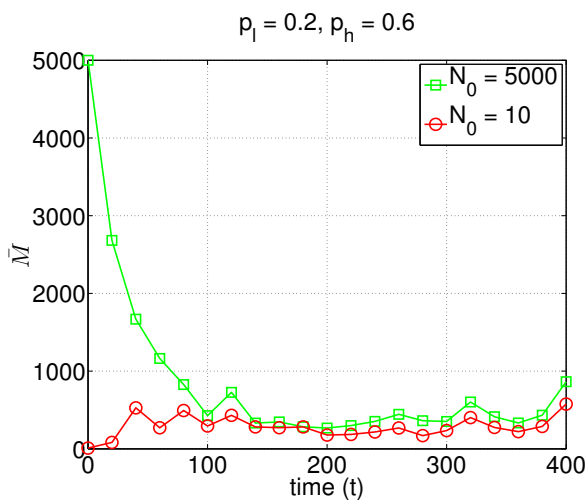


Fig. 3: Lorenz Model (Section V-A). Evolution of the number of particles adapted by the proposed algorithm when the initial number of particles $M_0 \in \{10, 5000\}$. The significance levels were set to $p_\ell = 0.2$ and $p_h = 0.6$.

$p(y_{i,t}|\mathbf{y}_{1:t-1})$, with $i = 1, 2$, and then, with both p-values, we adapt the number of particles as follows: if at least one of the marginals requires more particles, we increase the number of particles; if both marginals indicate no need for change of the number of particles, we keep it fixed; otherwise, we decrease the number.

Table IV shows the MSE in the approximation of the posterior mean, averaged number of particles \bar{M} , averaged p-value (over both dimensions), and the running time. Note that we can extract similar conclusions as in the case with scalar observations.

5) *Discussion:* The assumption (C) of Section III states that the tails of the pdf $p(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ should not be too heavy. Nevertheless, we have shown that the constraint is rather weak, since it is satisfied for all exponential-type distributions as well as for many heavy-tailed distributions. In practice, $p(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ cannot be characterized for most models in a closed form. Here we show the particle approximation of the observation predictive pdf $p^M(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ in the Lorenz 63 model at two different time steps. Figure 6 shows $p^M(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ with $M = 2^{14}$ particles in log-scale at time $t = 9601$. The approximated pdf $p^M(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ is compared with a Gaussian pdf and a Student's t-distribution (with $\nu = 3$), all of them with the same mean and variance. Figure 7 shows the same distributions at a different time step ($t = 10201$). Note that $p^M(\mathbf{y}_t|\mathbf{y}_{1:t-1})$ has very light tails at both time steps, and therefore, the assumption (C) holds in both numerical examples.

B. Non-linear growth model with heavy-tailed observation noise

In this numerical example, we consider the problem of tracking a modified version of the non-linear growth model in [7]. The state and observation equations are given by

$$x_t = \frac{x_{t-1}}{2} + \frac{25x_{t-1}}{1 + x_{t-1}^2} + 8 \cos(\phi t) + u_t, \quad (19)$$

$$y_t = \frac{x_t^2}{20} + v_t, \quad (20)$$

where $\phi = 0.4$ is a frequency parameter (in rad/s), $\{u_t\}_{t \geq 1}$ denotes a sequence of independent zero-mean univariate Gaussian r.v.'s with variance $\sigma_u^2 = 2$, and $\{v_t\}_{t \geq 1}$ is a sequence of independent Student's t-distributed r.v.'s with $\nu = 5$ degrees of freedom. The model is run for $t = 1, 2, \dots, T$, with $T = 5,000$.

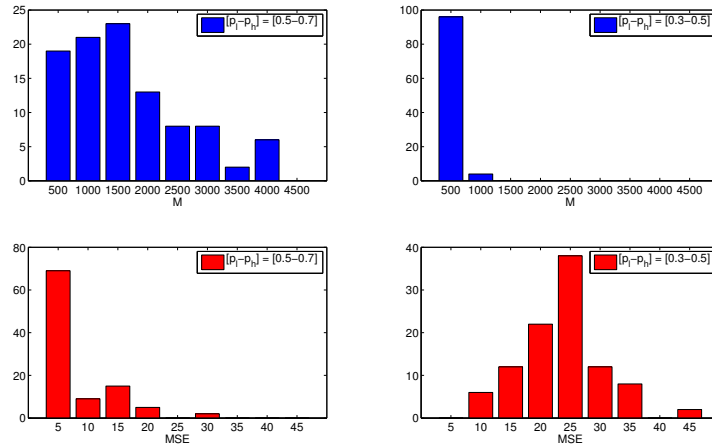


Fig. 4: Lorenz Model (Section V-A). Histograms of averaged MSE and M with $[p_\ell - p_h] = [0.3 - 0.5]$ and $[p_\ell - p_h] = [0.5 - 0.7]$. In both cases, the initial number of particles $M_0 = 5000$.

$[p_l - p_h]$	[0.4 - 0.8]	[0.3 - 0.7]	[0.3 - 0.65]	[0.25 - 0.65]	[0.2 - 0.6]	[0.15 - 0.55]	[0.1 - 0.5]
MSE	2.7151	2.7131	2.8351	3.8862	4.0814	5.4015	7.0323
\bar{M}	26175	19652	15788	7761	3858	539	203
p-val	0.5020	0.4953	0.4858	0.4906	0.4914	0.4820	0.4869
exec. time (s)	2937.9851	2120.0787	1744.3426	772.2125	373.6780	73.1735	38.3487

TABLE IV: Outputs of the particle filter with adaptive M for the Lorenz model (Section V-A) with parameters $\Delta = 10^{-3}$, $T_{obs} = 400\Delta$, $\sigma^2 = 0.5$ and 2-dimensional observations. The algorithm parameters are chosen as $W = 20$, $K = 7$, $M_{\max} = 2^{15}$ and $M_{\min} = 2^7$. We display the MSE in the approximation of the posterior mean, the averaged number of particles \bar{M} , averaged p-value (over both dimensions), and the running time.

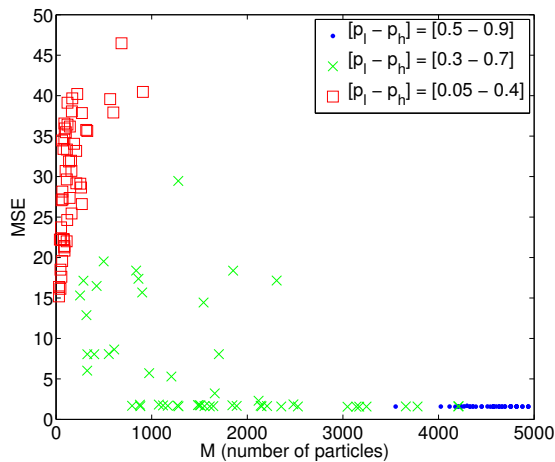


Fig. 5: Lorenz Model (Section V-A). MSE w.r.t. the averaged number of particles M for runs with different sets of thresholds.

First, we have run the standard BPF (with a fixed number of particles) for M in the range between 2 and 2^{14} . Figure 8 shows, for each value of the fixed number of particles M , the MSE of the approximation of the posterior mean, the averaged p-value p^* computed in the algorithm of Table II, and the running time. As expected, the MSE decreases

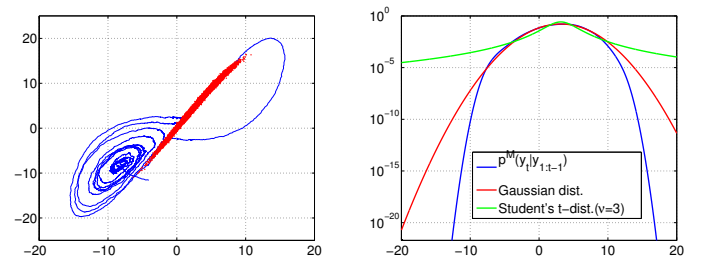


Fig. 6: Approximated observation predictive pdf $p^M(y_t|y_{1:t-1})$, Gaussian distribution, and Student's t-distribution ($\nu = 3$) in log-scale for the stochastic Lorenz 63 example with $M = 2^{14}$ particles. All distributions have the same mean and variance.

with the number of particles, at the expense of increasing the computational complexity of the filter. Note also that, over a certain range of M (namely, $M \geq 2^5$), the performance does not significantly improve. Finally, we see that in this example when the performance is poor, the p-value is very low (in average). This p-value is increased to ≈ 0.5 when the performance of the filter improves.

Then, we have run the particle filter with adaptive number of particles in Table II, with $K = 5$ fictitious observations, window size $W = 15$, p-value thresholds $[p_l - p_h] \in \{[0.4 - 0.68], [0.35 - 0.75], [0.3 - 0.7], [0.3 - 0.65], [0.25 - 0.65], [0.2 - 0.6]\}$, initial number of particles $M_0 = 2^9$,

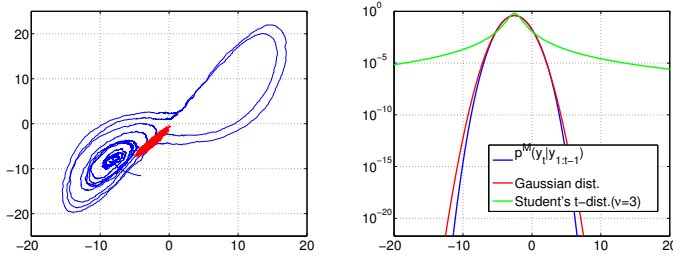


Fig. 7: Approximated observation predictive pdf $p^M(y_t|y_{1:t-1})$, Gaussian distribution, and Student's t-distribution ($\nu = 3$) in log-scale for the stochastic Lorenz 63 example with $M = 2^{14}$ particles. All distributions have the same mean and variance.

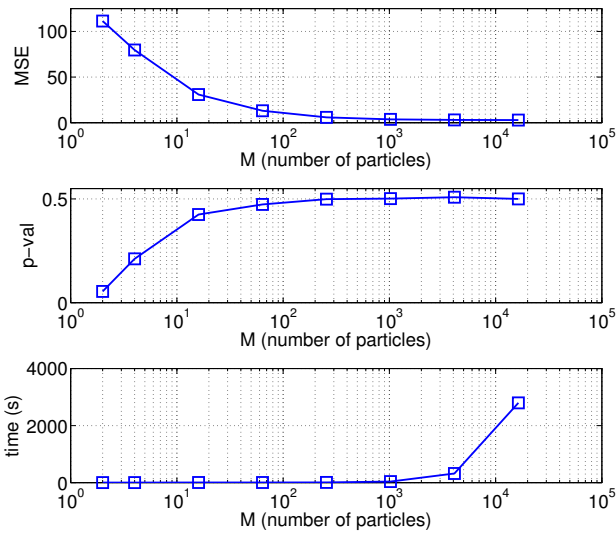


Fig. 8: BPF applied to a stochastic growth model with Student's t-distributed noise, and with fixed number of particles (Section V-B). MSE in the approximation of the posterior mean (top), averaged p-value (middle) and running time (bottom). The results are averaged over 50 independent simulations.

maximum and minimum number of particles $M_{\max} = 2^{14}$ and $M_{\min} = 2^4$, respectively, $f_{\text{up}}(M_{n-1}) = 2M_{n-1}$, and $f_{\text{down}}(M_{n-1}) = M_{n-1}/2$.

Table V displays the MSE of the approximation of the posterior mean, the averaged number of particles, the average p-value, and the running time in seconds, for the different choices of $[p_l - p_h]$. The results are averaged over 50 independent trials. Again, the pair of thresholds $[p_l - p_h]$ allows to operate at different complexity-performance regimes; decreasing the pair of parameters worsens the performance of the filter but enables a reduction in computational load.

VI. CONCLUSIONS

In practice, the number of particles needed in a particle filter is usually determined in an ad hoc manner. Furthermore, this number is typically kept constant throughout tracking. In

this paper, we have proposed a methodology for the online determination of the number of particles needed by the filter. The approach is based on assessing the convergence of the predictive distribution of the observations online. First we have proved, under standard assumptions, a novel convergence result on the approximation of this distribution. Then, we have proposed a method for adapting the number of particles based on the online assessment of the filter convergence. We have illustrated the performance of the suggested algorithm by computer simulations. The proposed procedure is simple but not unique. Namely, with the proposed methodology one can develop a range of algorithms for adapting the number of particles. Furthermore, while the analysis and examples have been presented for the standard bootstrap particle filter for simplicity and clarity, it is straightforward to extend it to more sophisticated algorithms using adaptive proposals [41] or parallelization schemes [42], [43].

APPENDIX A PROOF OF PROPOSITION 2

The sequence of r.v.'s $\{A_{K,t}\}_{t \geq 1}$ are *constructed* to be independent. To see this, let us look into the generation of $A_{K,t}$ and $A_{K,t+1}$. Below, we are using capital letters to denote a r.v. (e.g., Y_t) and lower-case letters for its realisations (e.g., y_t).

At time t , the r.v. $A_{K,t}$ is constructed by means of a nonlinear transformation of the r.v.'s Y_t and $\{\tilde{Y}_t^{(k)}\}_{k=1,\dots,K}$. The latter are referred to as fictitious observations in the paper. Let us denote this many-to-one transformation as ψ , i.e.,

$$A_{K,t} = \psi(Y_t, \tilde{Y}_t^{(1)}, \dots, \tilde{Y}_t^{(K)}). \quad (21)$$

Under the sole assumption that $\{Y_t, \tilde{Y}_t^{(1)}, \dots, \tilde{Y}_t^{(K)}\}$ are i.i.d. continuous r.v.'s, Proposition 1 states that $A_{K,t}$ has a uniform probability distribution. To be precise, $A_{K,t}$ takes values on $\{0, \dots, K\}$, and its probability mass function is $P(A_{K,t} = n) = \frac{1}{K+1}$ for every $n \in \{0, \dots, K\}$.

In our case, the common pdf of the r.v.'s $\{Y_t, \tilde{Y}_t^{(1)}, \dots, \tilde{Y}_t^{(K)}\}$ is $p_t(y_t) = \int g_t(y_t, \mathbf{x}) \xi_t(d\mathbf{x}) = p(y_t|y_{1:t-1})$. However, the actual form of $p(y_t|y_{1:t-1})$ plays no role whatsoever in Proposition 1. In other words, $A_{K,t}$ is uniform as long as $\{Y_t, \tilde{Y}_t^{(1)}, \dots, \tilde{Y}_t^{(K)}\}$ are i.i.d. and this results holds independently of the actual sequence $y_{1:t-1}$ (which determines the form of $p(y_t|y_{1:t-1})$).

We move on to time $t+1$. The r.v. $A_{K,t+1}$ is obtained as a nonlinear transformation of $\{Y_{t+1}, \tilde{Y}_{t+1}^{(1)}, \dots, \tilde{Y}_{t+1}^{(K)}\}$, namely,

$$A_{K,t+1} = \psi(Y_{t+1}, \tilde{Y}_{t+1}^{(1)}, \dots, \tilde{Y}_{t+1}^{(K)}).$$

From Proposition 1, if $\{Y_{t+1}, \tilde{Y}_{t+1}^{(1)}, \dots, \tilde{Y}_{t+1}^{(K)}\}$ are i.i.d. then $A_{K,t+1}$ has a uniform distribution, i.e., $P(A_{K,t+1} = n) = \frac{1}{K+1}$ for every $n \in \{0, \dots, K\}$. As before, this is true independently of the specific common pdf of the r.v.'s $\{Y_{t+1}, \tilde{Y}_{t+1}^{(1)}, \dots, \tilde{Y}_{t+1}^{(K)}\}$. This common pdf is $p_{t+1}(y_{t+1}) = (g_{t+1}^{y_{t+1}}, \xi_{t+1}) = p(y_{t+1}|y_{1:t})$ and, therefore, $A_{K,t+1}$ is uniform without regard to the sequence $y_{1:t}$ (which determines the form of $p(y_{t+1}|y_{1:t})$) and, in particular, without regard to the observed realisation $Y_t = y_t$.

$[p_l - p_h]$	[0.4 – 0.8]	[0.35 – 0.75]	[0.3 – 0.7]	[0.3 – 0.65]	[0.25 – 0.65]	[0.2 – 0.6]
MSE	2.8707	3.4945	4.7687	9.0465	10.5826	17.6967
\bar{M}	9739	7478	6251	3168	2087	232
p-val	0.4976	0.4950	0.4893	0.4837	0.4730	0.4604
exec. time (s)	3613	2515	1427	561	234	21

TABLE V: Output of the algorithm in Table II for a stochastic growth model with Student's t-distributed observation noise, with adaptive M (Section V-B). The algorithm parameters are chosen as $W = 15$, $K = 1$, $M_{\max} = 2^{14}$, $M_{\min} = 2^6$. We display the MSE in the approximation of the posterior mean, the average number of particles \bar{M} , the average p-value, and the running time.

Now, since $A_{K,t+1}$ is uniform for any $Y_t = y_t$ (and, obviously, for any $\tilde{Y}_t^{(k)} = \tilde{y}_t^{(k)}$, $k = 1, \dots, K$), and $A_{K,t}$ is obtained as a transformation of $\{Y_t, \tilde{Y}_t^{(1)}, \dots, \tilde{Y}_t^{(K)}\}$ (see (21) above), then it follows that $A_{K,t+1}$ has a uniform distribution for every possible realisation $A_{K,t} = n$. This implies that the conditional distribution of $A_{K,t+1}$ given $A_{K,t}$ is uniform, i.e.,

$$P(A_{K,t+1} = n | A_{K,t} = m) = \frac{1}{K+1}, \quad (22)$$

$\forall (n, m) \in \{0, \dots, K\} \times \{0, \dots, K\}$. However, Eq. (22) readily entails independence. If we let $P(A_{K,t+1}, A_{K,t})$ denote the joint probability mass function of $A_{K,t+1}$ and $A_{K,t}$, then from the definition of conditional probability we have

$$\begin{aligned} P(A_{K,t+1} = n, A_{K,t} = m) &= \\ P(A_{K,t+1} = n | A_{K,t} = m) P(A_{K,t} = m) &= \\ \frac{1}{K+1} \times \frac{1}{K+1} &= \\ P(A_{K,t+1} = n) P(A_{K,t} = m), \end{aligned} \quad (23)$$

for any n and m within the set $\{0, \dots, K\}$.

APPENDIX B PROOF OF THEOREM 2

Let Y_t denote the (random) observation at time t . Assume, without loss of generality, that $\mathcal{Y} = \mathbb{R}$. The probability measure associated to $Y_t | Y_{1:t-1} = y_{1:t-1}$ is $\mu_t(dy)$ and, therefore, we can write the cumulative distribution function of $Y_t | Y_{1:t-1} = y_{1:t-1}$ as $F_t(z) = (I_{(-\infty, z]}, \mu_t)$, where

$$I_A(y) = \begin{cases} 1, & \text{if } y \in A \\ 0, & \text{otherwise} \end{cases}$$

is the indicator function. Obviously, $\|I_A\|_\infty = 1 < \infty$ independently of the set A and, therefore, Theorem 1 yields

$$\lim_{M \rightarrow \infty} F_t^M(z) = F_t(z) \quad \text{a.s.}$$

for any $z \in \mathbb{R}$, where $F_t^M(z) = (I_{(-\infty, z]}, \mu_t^M)$ is the approximation of the cdf of $Y_t | Y_{1:t-1} = y_{1:t-1}$ provided by the BPF.

Assume the actual observation is $Y_t = y_t$ and we draw K i.i.d. fictitious observations $\tilde{y}_t^{(1)}, \dots, \tilde{y}_t^{(K)}$ from the distribution with cdf F_t^M . Given $Y_t = y_t$ is fixed, the probability that exactly n out of K of these samples are lesser than y_t coincides with the probability to have n successes out of K trials for a binomial r.v. with parameter (i.e., success probability) $F_t^M(y_t)$, which can be written as

$$h_n^M(y_t) = \binom{K}{n} (F_t^M(y_t))^n (1 - F_t^M(y_t))^{K-n}.$$

By integrating $h_n^M(y_t)$ over the predictive distribution of Y_t , we obtain the probability to have exactly n fictitious observations, out of K , which are less than the r.v. Y_t , i.e., the probability that $A_{K,M,t} = n$ is

$$\mathbb{Q}_{K,M,t}(n) = (h_n^M, \mu_t). \quad (24)$$

However, Theorem 1 yields $\lim_{M \rightarrow \infty} (h_n^M, \mu_t^M) = (h_n^M, \mu_t)$ a.s.⁴ and, in particular, there exists a sequence of non-negative r.v.'s $\{\varepsilon_M\}_{M \geq 1}$ such that $\lim_{M \rightarrow \infty} \varepsilon_M = 0$ a.s. and

$$(h_n^M, \mu_t^M) - \varepsilon_M \leq (h_n^M, \mu_t) \leq (h_n^M, \mu_t^M) + \varepsilon_M \quad (25)$$

for each M . Moreover, it is apparent that $(h_n^M, \mu_t^M) = \frac{1}{K+1}$ (see Proposition 1) which, together with (24) and (25) yields the desired relationship

$$\frac{1}{K+1} - \varepsilon_M \leq \mathbb{Q}_{K,M,t}(n) \leq \frac{1}{K+1} + \varepsilon_M$$

for every $n \in \{0, \dots, K\}$. \square

REFERENCES

- [1] M. West and J. Harrison, *Bayesian Forecasting*, 2nd ed., Springer-Verlag, New York, 1996.
- [2] B. Ristic, S. Arulampalam, and N. Gordon, *Beyond the Kalman Filter*, Artech House, Boston, 2004.
- [3] R. E. Kalman, "A new approach to linear filtering and prediction problems," *Journal of Basic Engineering*, vol. 82, pp. 35–45, 1960.
- [4] B. D. O. Anderson and J. B. Moore, *Optimal Filtering*, Englewood Cliffs, 1979.
- [5] L. R. Rabiner and B. H. Juang, "An introduction to hidden markov models," *ASSP Magazine, IEEE*, vol. 3, no. 1, pp. 4–16, 1986.
- [6] J. S. Liu, R. Chen, and W. H. Wong, "Rejection control and sequential importance sampling," *Journal of the American Statistical Association*, vol. 93, no. 443, pp. 1022–1031, September 1998.
- [7] A. Doucet, S. Godsill, and C. Andrieu, "On sequential Monte Carlo Sampling methods for Bayesian filtering," *Statistics and Computing*, vol. 10, no. 3, pp. 197–208, 2000.
- [8] A. Doucet, N. de Freitas, and N. Gordon, Eds., *Sequential Monte Carlo Methods in Practice*, Springer, New York (USA), 2001.
- [9] P. M. Djurić, J. H. Kotecha, J. Zhang, Y. Huang, T. Ghirmai, M. F. Bugallo, and J. Míguez, "Particle filtering," *IEEE Signal Processing Magazine*, vol. 20, no. 5, pp. 19–38, September 2003.
- [10] O. Cappé, S. J. Godsill, and E. Moulines, "An overview of existing methods and recent advances in sequential Monte Carlo," *Proceedings of the IEEE*, vol. 95, no. 5, pp. 899–924, 2007.

⁴Note that $\|h_n^M\|_\infty = 1$ independently of n and M . If we recall the proof of Theorem 1, namely inequality (??), we observe that the error rates for the approximation errors of the form $|(h, \mu_t^M) - (h, \mu_t)|$ depend on the test function h only through its supremum $\|h\|_\infty$, i.e., the r.v. W_t^ε in (??) only depends on the observations $y_{1:t-1}$ and the model (specifically the likelihood functions). Therefore, Theorem 1 (the same as, e.g., Lemmas ?? and ??) also holds for any test function that depends on M (even a random one) as long as its supremum is deterministic and independent of M . This is the case of function $h_n^M(y)$.

- [11] N. Gordon, D. Salmond, and A. F. M. Smith, "Novel approach to nonlinear and non-Gaussian Bayesian state estimation," *IEE Proceedings-F Radar and Signal Processing*, vol. 140, pp. 107–113, 1993.
- [12] A. Bain and D. Crisan, *Fundamentals of Stochastic Filtering*, Springer, 2008.
- [13] P. Del Moral and L. Miclo, "Branching and interacting particle systems. Approximations of Feynman-Kac formulae with applications to nonlinear filtering," *Lecture Notes in Mathematics*, pp. 1–145, 2000.
- [14] D. Fox, "Adapting the sample size in particle filters through KLD-sampling," *The International Journal of Robotics Research*, vol. 22, no. 12, pp. 985–1003, 2003.
- [15] A. Soto, "Self adaptive particle filter," in *IJCAI*, 2005, pp. 1398–1406.
- [16] O. Straka and M. Šimandl, "Particle filter adaptation based on efficient sample size," in *14th IFAC Symposium on System Identification*, 2006.
- [17] A. Lee and N. Whiteley, "Variance estimation and allocation in the particle filter," *arXiv:1509.00394v1 [stat.CO]*, 2015.
- [18] P. Del Moral, *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*, Springer, 2004.
- [19] A. Bhadra and E. L. Ionides, "Adaptive particle allocation in iterated sequential Monte Carlo via approximating meta-models," *Statistics and Computing*, vol. 26, no. 1-2, pp. 393–407, 2016.
- [20] F. LeGland and N. Oudjane, "A sequential particle algorithm that keeps the particle system alive," in *13th European Signal Processing Conference*. IEEE, 2005, pp. 1–4.
- [21] A. Jasra, A. Lee, C. Yau, and X. Zhang, "The alive particle filter," *arXiv:1304.0151*, 2013.
- [22] P. Del Moral, A. Jasra, A. Lee, C. Yau, and X. Zhang, "The alive particle filter and its use in particle Markov chain Monte Carlo," *Stochastic Analysis and Applications*, vol. 33, no. 6, pp. 943–974, 2015.
- [23] X.L. Hu, T.B. Schon, and L. Ljung, "A basic convergence result for particle filtering," *IEEE Transactions on Signal Processing*, vol. 56, no. 4, pp. 1337–1348, 2008.
- [24] J. Cornebise, *Adaptive Sequential Monte Carlo Methods*, Ph.D. thesis, PhD thesis, Télécom ParisTech, 2010. 38, 49, 2009.
- [25] A. Beskos, D. Crisan, and A. Jasra, "On the stability of sequential monte carlo methods in high dimensions," *The Annals of Applied Probability*, vol. 24, no. 4, pp. 1396–1445, 2014.
- [26] P. M. Djurić and J. Míguez, "Assessment of nonlinear dynamic models by Kolmogorov-Smirnov statistics," *IEEE Transactions on Signal Processing*, vol. 58, no. 10, pp. 5069–5079, 2010.
- [27] R. Chen, X. Wang, and J. S. Liu, "Adaptive joint detection and decoding in flat-fading channels via mixture kalman filtering," *IEEE Transactions on Information Theory*, vol. 46, no. 6, pp. 2079–2094, September 2000.
- [28] N. Chopin, P. E. Jacob, and O. Papaspiliopoulos, "SMC2: An efficient algorithm for sequential analysis of state space models," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2012.
- [29] D. Crisan and J. Míguez, "Nested particle filters for online parameter estimation in discrete-time state-space markov models," *to appear in Bernoulli (arXiv: 1308.1883v1 [stat.CO])*, 2016.
- [30] C. Andrieu, A. Doucet, and R. Holenstein, "Particle Markov chain Monte Carlo methods," *Journal of the Royal Statistical Society B*, vol. 72, no. 3, pp. 269–342, 2010.
- [31] E. Koblenz and J. Míguez, "A population monte carlo scheme with transformed weights and its application to stochastic kinetic models," *Statistics and Computing*, vol. 25, no. 2, pp. 407–425, 2015.
- [32] A. Doucet, N. de Freitas, and N. Gordon, "An introduction to sequential Monte Carlo methods," in *Sequential Monte Carlo Methods in Practice*, A. Doucet, N. de Freitas, and N. Gordon, Eds., chapter 1, pp. 4–14. Springer, 2001.
- [33] J. Carpenter, P. Clifford, and P. Fearnhead, "Improved particle filter for nonlinear problems," *IEE Proceedings - Radar, Sonar and Navigation*, vol. 146, no. 1, pp. 2–7, February 1999.
- [34] D. Crisan, "Particle filters - a theoretical perspective," in *Sequential Monte Carlo Methods in Practice*, A. Doucet, N. de Freitas, and N. Gordon, Eds., chapter 2, pp. 17–42. Springer, 2001.
- [35] T. Li, M. Bolić, and P. M. Djurić, "Resampling methods for particle filtering," *IEEE Signal Processing Magazine*, vol. 32, no. 3, pp. 70–86, May 2015.
- [36] K. Heine and D. Crisan, "Uniform approximations of discrete-time filters," *Advances in Applied Probability*, vol. 40, no. 4, pp. 979–1001, 2008.
- [37] R. L. Plackett, "Karl pearson and the chi-squared test," *International Statistical Review/Revue Internationale de Statistique*, pp. 59–72, 1983.
- [38] E. N. Lorenz, "Deterministic nonperiodic flow," *Journal of Atmospheric Sciences*, vol. 20, no. 2, pp. 130–141, 1963.
- [39] A. J. Chorin and P. Krause, "Dimensional reduction for a Bayesian filter," *PNAS*, vol. 101, no. 42, pp. 15013–15017, October 2004.
- [40] M. S. Nikulin, "Hellinger distance," *Encyclopedia of Mathematics*, 2001.
- [41] J. Cornebise, E. Moulines, and J. Olsson, "Adaptive methods for sequential importance sampling with application to state space models," *Statistics and Computing*, vol. 18, no. 4, pp. 461–480, 2008.
- [42] N. Whiteley, A. Lee, and K. Heine, "On the role of interaction in sequential monte carlo algorithms," *Bernoulli*, vol. 22, no. 1, pp. 494–529, 2016.
- [43] B. Paige, F. Wood, A. Doucet, and Y. W. Teh, "Asynchronous anytime sequential monte carlo," in *Advances in Neural Information Processing Systems*, 2014, pp. 3410–3418.



Víctor Elvira received the M.Sc. and Ph.D. degrees in electrical engineering from Universidad de Cantabria (Spain) in 2008 and 2011, respectively. In 2012, he joined Universidad Carlos III de Madrid (Spain) as postdoctoral researcher, and later as an Assistant Professor. In 2016, he joined Télécom Lille (Institut Mines-Télécom) where he is currently an Associate Professor. He also belongs to the CRISTAL laboratory (UMR CNRS 9189). He has been a visiting scholar at the IHP Leibniz Institute (Frankfurt Oder, Germany), University of Helsinki (Finland), Stony Brook University of New York (USA), Federal University of Rio de Janeiro (Brazil), and Paris-Dauphine University (France). His research interests include computational statistics, statistical signal processing, Bayesian analysis, and biomedical signal processing. He has co-authored over 40 journal and peer-reviewed conference papers.



Joaquín Míguez received the M.Sc. and Ph.D. degrees in computer engineering from the University of A Coruña (A Coruña, Spain) in 1997 and 2000, respectively. He has held permanent positions at the Department of Electronics and Systems, University of A Coruña (2000–03), the School of Mathematical Sciences, Queen Mary University of London (2015–2016), and the Department of Signal Theory & Communications, Universidad Carlos III de Madrid (2004–15 and 2016–present). He has also held visiting positions in the Department of Electrical & Computer Engineering of the State University of New York at Stony Brook (2001) and the Department of Mathematics of Imperial College London (2013–14). His research interests are in the fields of applied probability, statistical signal processing, Bayesian analysis, dynamical systems and the theory and applications of Monte Carlo methods. Dr. Míguez has co-authored over 50 international journal papers in the fields of signal processing, communications, mathematical physics, probability and statistics. He has delivered lectures and seminars on various European universities and research centres. He is a co-recipient of the IEEE Signal Processing Magazine Best Paper Award 2007.



Petar M. Djurić (M'90–SM'99–F'06) received the B.S. and M.S. degrees in electrical engineering from the University of Belgrade, Belgrade, in 1981 and 1986, respectively, and the Ph.D. degree in electrical engineering from the University of Rhode Island, Kingston, RI, in 1990. Since 1990, he has been a Professor with the Department of Electrical and Computer Engineering, Stony Brook University, Stony Brook, NY. From 1981 to 1986, he was a Research Associate with the Vinča Institute of Nuclear Sciences, Belgrade. His research interests

include the area of signal and information processing with primary interests in the theory of signal modeling, detection, and estimation; Monte Carlo-based methods; signal and information processing over networks; and applications in a wide range of disciplines. He has been invited to lecture at many universities in the United States and overseas. He received the IEEE Signal Processing Magazine Best Paper Award in 2007 and the EURASIP Technical Achievement Award in 2012. In 2008, he was the Chair of Excellence of Universidad Carlos III de Madrid-Banco de Santander. From 2008 to 2009, he was a Distinguished Lecturer of the IEEE Signal Processing Society. He has been on numerous committees of the IEEE Signal Processing Society and of many professional conferences and workshops. He is a Fellow of EURASIP and the Editor-in-Chief of the IEEE TRANSACTIONS ON SIGNAL AND INFORMATION PROCESSING OVER NETWORKS.