

# Sequential Monte Carlo methods for the optimization of a general class of objective functions

Joaquín Míguez\*

Dan Crisan<sup>†</sup>

Petar M. Djurić<sup>‡</sup>

February 2, 2010

## Abstract

We introduce a Monte Carlo method for global optimization of a certain class of (possibly non-convex and possibly non-differentiable) cost functions with respect to a high dimensional signal of interest. The proposed approach involves the transformation of the optimization problem into one of inference in a discrete-time dynamical system in state-space form. In particular, we describe a methodology for constructing an associated state-space model which has the signal of interest as its unobserved dynamic state. The model is matched to the cost function in the sense that the *maximum a posteriori* (MAP) estimate of the system state is also a global minimizer of the cost. The advantage of recasting the optimization problem in an estimation framework is that we can apply the family of sequential Monte Carlo algorithms as an efficient aid for the numerical search of solutions. This kind of techniques produce, in a natural way, a random grid that is dense in the regions where the *a posteriori* probability mass is high (hence, the cost is low) and sparse elsewhere. Simple search techniques can then be applied to locate the best point in the grid with limited extra computations. In the paper, we describe two candidate algorithms, prove that they converge almost surely to a global minimizer of the cost and provide two application examples, including some illustrative numerical results<sup>1</sup>.

## Keywords

Global optimization; MAP estimation; sequential Monte Carlo; convergence of particle filter; state space models; interacting particle systems.

## AMS (MSC2000)

65C35, 65K10, 65C05, 90C56, 60K35, 90C47.

---

\*Department of Signal Theory & Communications, Universidad Carlos III (Spain). E-mail: joaquin.miguez@uc3m.es

<sup>†</sup>Department of Mathematics, Imperial College London (UK). E-mail: d.crisan@imperial.ac.uk

<sup>‡</sup>Department of Electrical & Computer Engineering, Stony Brook University (USA). E-mail: djuric@ece.sunysb.edu

<sup>1</sup>A preliminary version of Proposition 1 in this manuscript was originally presented at the 13th IEEE Digital Signal Processing Workshop (DSP'2009).

# 1 Introduction

Many scientific and engineering problems involve the optimization of general high-dimensional objective functions, not necessarily smooth and possibly with many local extrema. Algorithms designed to be used in such problems are commonly referred to as “global optimization” methods and their application often encounters difficulties related to convergence, numerical stability and computational complexity. We can classify the techniques in the field as either deterministic or stochastic [26]. Deterministic approaches exploit analytical properties of the objective function, typically convexity [18], monotonicity [33] or smoothness [35]. Such techniques can be difficult to apply with general functions, or they can simply be inconvenient to use in the initial stages of the treatment of a real-world problem, when the objective function is being defined and is subject to change. For this reason, many state-of-the-art global optimization algorithms are stochastic, meaning that they involve Monte Carlo simulations and, as a consequence, their output is random even for a fixed input. Examples of stochastic procedures include multistart [27], random search [2], simulated annealing [17] or evolutionary [32] methods.

On the other hand, a body of knowledge on Monte Carlo techniques for inference in broad classes of statistical models has evolved over the past three decades, including Markov-chain Monte Carlo (MCMC) [13] and sequential Monte Carlo [8] methods. The study of the connections between inference and optimization may provide new insights and may possibly lead to the development of some techniques. A well-known example is the interpretation of the simulated annealing method as an MCMC algorithm [29]. More recently, it has been proposed to tackle deterministic optimal control problems by way of Monte Carlo nonlinear filtering methods [23]. The latter approach can be interpreted as an application of the Maslov optimization theory [21].

In this paper, we investigate the transformation of a certain class of optimization problems into equivalent maximum *a posteriori* (MAP) estimation problems in dynamic state-space systems. Specifically, we consider the minimization of a cost function  $C_T(x_{0:T})$ , where  $x_{0:T} = \{x_0, \dots, x_T\}$  is a high-dimensional set of unknowns and where the function can be recursively decomposed into a sequence of costs  $C_{T-1}(x_{0:T-1})$ ,  $C_{T-2}(x_{0:T-2})$ , ...,  $C_0(x_0)$ , with lower-dimensional supports (i.e., the dimension of  $x_{0:t}$  is lower than the dimension of  $x_{0:t+1}$ ). We prove that for costs of this type it is possible to design an associated state-space system such that any MAP estimate of the state of the system at time  $T$  coincides with a global minimum of the cost. The advantage of this reformulation of the problem is that we can draw from a pool of sequential Monte Carlo methods for inference in state-space models. These techniques

produce, in a natural way, a random grid in the space of the unknowns that is dense in the regions where the *a posteriori* probability mass is high (and, equivalently, the cost is low) and sparse elsewhere. A subsequent (simple) search over this grid yields an estimate of the global minimizer. In this work, we start from the standard sequential Monte Carlo technique termed sequential importance resampling [16] (see also [7]) and study its combination with two search procedures. The first one is a direct search over the sample paths in the space of  $x_{0:T}$  generated by the SIR algorithm and has a complexity that grows linearly with the number of samples. The second one performs a trellis search over an extended grid using the Viterbi algorithm [11], as originally suggested in [14], and its complexity grows with the square of the number of samples. Both search procedures can be implemented sequentially and together with the SIR method. Our main contribution is to prove that the two resulting optimization algorithms converge almost surely to a global minimum of the cost (such analysis was not addressed in [14]) and to obtain a lower bound on the the number of samples needed to attain a certain accuracy. These results are new to our best knowledge and start from the derivation of  $L_p$ -bounds for the error of the bootstrap filter in the path space of  $x_{0:T}$ . Such bounds were not explicitly available in the previous literature (see [5] for  $L_2$ -bounds in the path space of  $x_{0:T}$  and [22] for  $L_p$ -bounds on the space of  $x_t$ ). We also provide two application examples, including a typical global optimization problem (the Neumaier 3 problem [1]) and the design of cross-talk cancellation acoustic filters [28]. We further use these two problems to illustrate the numerical performance of the proposed algorithms.

Our approach bears similarities to the work in [23]. However, we do not restrict ourselves to optimal control applications and consider a broader class of minimization problems instead. Indeed, the objective functions studied in [23] are instances of the family of additive costs in Section 5 of this paper. A comparison of our analysis of the asymptotic convergence of the resulting optimization algorithms and that presented in [23] is presented in Remark 3.

The remaining of this paper is organized as follows. After a brief introduction to the notations in the paper, Section 3 describes the class of optimization problems of interest. Their reformulation as inference problems in state-space systems is introduced in Section 4. Two examples, including additive cost functions and minimax problems, are investigated in more detail in Section 5. The sequential Monte Carlo algorithms for global MAP estimation based on state-space models are described and analyzed in Section 6. Section 7 shows some numerical results and, finally, Section 8 is devoted to the conclusions.

## 2 Notation

Random variates and their realizations are represented by the same upper- and lower-case letter, e.g., the random variate  $X$  and its realization  $X = x$ . Random sequences are denoted as  $\{X_t\}_{t \in \mathbb{N}}$ .

Probability density functions (pdf's) are indicated by the letter  $\pi$ . This is an argument-wise notation, hence for the random variates  $X$  and  $Y$ ,  $\pi(x)$  signifies the density of  $X$ , possibly different from  $\pi(y)$ , which represents the pdf of  $Y$ . The integral of a function  $f(x)$  with respect to a measure with density  $\pi(x)$  is denoted by the shorthand  $(f, \pi) \triangleq \int f(x)\pi(x)dx$ .

The letter  $C$  is used throughout the paper to denote costs. It may be an overall cost for a sequence (and we use upper-case letters,  $C_t$  and  $\mathbf{C}_t$ ) or a partial cost for a subsequence (and we write lower-case letters,  $c_t$  and  $\mathbf{c}_t$ ).

## 3 Problem statement

We address the problem of finding the global minima of a certain class of cost functions with recursive structures. Specifically, let  $\{x_t\}_{t \in \mathbb{N} \cup \{0\}}$  and  $\{y_t\}_{t \in \mathbb{N}}$  be discrete-time sequences in  $\mathbb{R}^{d_x}$  and  $\mathbb{R}^{d_y}$ , respectively, where  $d_x$  and  $d_y$  are the (integer) dimensions of the vectors of each sequence. For some arbitrarily large but finite horizon  $T$ , we aim at computing

$$\mathbf{X}_T^c = \arg \min_{x_{0:T}} C_T(x_{0:T}; y_{1:T}), \quad (1)$$

where  $C_T(\cdot; y_{1:T}) : (\mathbb{R}^{d_x})^{T+1} \rightarrow \mathbb{R}^+$  is the real non-negative cost function of interest, the subsequence  $x_{0:T} = \{x_0, x_1, \dots, x_T\}$  denotes the unknowns to be optimized and the subsequence  $y_{1:T} = \{y_1, y_2, \dots, y_T\}$  is known and provides the fixed parameters that determine the specific form of  $C_T$ . The set  $\mathbf{X}_T^c$  contains all the subsequences  $x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}$  for which  $C_T(\cdot; y_{1:T})$  attains its minimum value.

The methods to be introduced in this paper are applicable when the cost function can be constructed recursively, i.e., when there exists a sequence of functions  $C_t(\cdot; y_{1:t}) : (\mathbb{R}^{d_x})^{t+1} \rightarrow \mathbb{R}^+$ ,  $t = 0, 1, \dots, T$ , such that  $C_t(x_{0:t}; y_{1:t})$  can be computed from  $C_{t-1}(x_{0:t-1}; y_{1:t-1})$  by some known update rule. In particular, we assume that  $C_t$  can be decomposed as

$$C_t(x_{0:t}; y_{1:t}) = H(C_{t-1}(x_{0:t-1}; y_{1:t-1}), c_t(x_{0:t}; y_t)), \quad t = 1, 2, \dots, T, \quad (2)$$

where  $H : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$  is the update function and  $c_t(\cdot; y_t) : (\mathbb{R}^{d_x})^{t+1} \rightarrow \mathbb{R}^+$  is termed the partial cost function at time  $t$ . The recursion is initialized with some function  $C_0 : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^+$  which does not formally depend on any element of the sequence  $y_{1:T}$ .

**Example 1** *A toy problem.*

Consider the Becker and Lago problem adapted from [1], which we write as

$$\min_{x_{1:T}} J(x_{1:T}) = \sum_{t=1}^T (|x_t| - 5)^2, \quad (3)$$

subject to  $x_t \in [-10, +10]$  for all  $t \in \{1, 2, \dots, T\}$  (the problem is originally stated for  $T = 2$  in [1]).

There are  $2^T$  obvious solutions at  $x_{1:T}^o \in \{\pm 5, \dots, \pm 5\}$ . For all global minima  $J(x_{1:T}^o) = 0$ .

We represent problem (3) using the notation described at the beginning of this section by defining the partial cost  $c_t(x_{0:t}; y_t) \triangleq (|x_t| - y_t)^2$  and taking an additive update function, i.e.,  $C_t(x_{0:t}; y_{1:t}) = C_{t-1}(x_{0:t-1}; y_{1:t-1}) + c_t(x_{0:t}; y_t)$ . Since the problem does not depend on  $x_0$ , we set  $C_0(x_0) = 0, \forall x_0 \in \mathbb{R}$ , to initialize the recursion. The cost function parameters are all equal in this case, namely  $y_t = 5$  for every  $t = 1, 2, \dots, T$ . The solution set is  $X_T^c = \{\pm 5, \dots, \pm 5\}$ .

Note that, very often, the partial cost  $c_t(x_{0:t}; y_t)$  does not actually depend on the complete sequence  $x_{0:t}$ , but only on some shorter subsequence  $x_{t-k:t}$  ( $k \geq 0$ ). In Example 1 we have the simplest case, in which  $c_t(x_{0:t}; y_t)$  is only a function of  $x_t$  (with fixed parameters  $y_t$ ).

Despite the simplicity of the example above, we may realistically expect that problems of the form of (1) be hard to solve in practical scenarios. Indeed,  $C_T(x_{0:T}; y_{1:T})$  may be analytically intractable and present multiple minima. Also, due to the high dimension of the unknown,  $x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}$ , it may be hard to devise a stable and convergent numerical algorithm with acceptable computational complexity.

In this paper, we propose to recast the optimization problem (1) as one of tracking the state of a dynamic state-space model. The advantage of this transformation is the availability of a pool of sequential Monte Carlo techniques that can be applied to numerically compute dynamic-state estimates. In Section 6 we describe specific algorithms of this class and show how they can asymptotically approximate any element in the solution set  $X_T^c$  with any desired accuracy. Before that, we consider the transformation of (1) into an equivalent estimation problem.

## 4 State-space models

In order to transform problem (1), we consider a state-space model where the unknowns,  $x_{0:T}$ , play the role of the system state and the cost-function parameters,  $y_{1:T}$ , are the associated observations.

To be specific, let  $\{X_t\}_{t \geq 0}$  and  $\{Y_t\}_{t > 0}$  be stochastic processes that take values in  $\mathbb{R}^{d_x}$  and  $\mathbb{R}^{d_y}$ , respectively. We refer to  $\{X_t\}_{t \geq 0}$  and  $\{Y_t\}_{t \geq 0}$  as state and observation processes, respectively. For  $t = 0$ ,

the random variate  $X_0$  has a pdf with respect to the Lebesgue measure in  $\mathbb{R}^{d_x}$ , that we denote as  $\pi(x_0)$ , and, for  $t > 0$ , the process evolves according to the conditional probability law

$$\Pr \{X_t \in A | X_{0:t-1} = x_{0:t-1}\} = \int_A \pi(x_t | x_{0:t-1}) dx_t, \quad (4)$$

where  $\pi(x_t | x_{0:t-1})$  denotes the pdf, with respect to the Lebesgue measure, of  $X_t$  given  $X_{0:t-1} = x_{0:t-1}$  and  $A$  is any Borel subset of  $\mathbb{R}^{d_x}$ . In the sequel, we will use  $\mathcal{B}(\mathbb{R}^d)$  to denote the Borel  $\sigma$ -algebra in  $\mathbb{R}^d$ .

The observation process,  $\{Y_t\}_{t>0}$ , follows the conditional probability law

$$\Pr \{Y_t \in A' | X_{0:t} = x_{0:t}, Y_{1:t-1} = y_{1:t-1}\} = \int_{A'} \pi(y_t | x_{0:t}, y_{1:t-1}) dy_t, \quad (5)$$

where  $\pi(y_t | x_{0:t}, y_{1:t-1})$  denotes the conditional pdf of  $Y_t$  given  $X_{0:t} = x_{0:t}$  and  $Y_{1:t-1} = y_{1:t-1}$ , again with respect to the Lebesgue measure, and  $A' \in \mathcal{B}(\mathbb{R}^{d_y})$ .

We will refer to the densities  $\pi(x_0)$  and  $\pi(x_t | x_{0:t-1})$  as the prior pdf and the transition pdf of the state process, respectively, while for fixed observations  $Y_{1:t} = y_{1:t}$ , the function  $g_t(x_{0:t}) \triangleq \pi(y_t | x_{0:t}, y_{1:t-1})$  is referred to as the likelihood of the state path  $X_{0:t} = x_{0:t}$  at time  $t$ . Together, the densities

$$\pi(x_0), \quad \pi(x_t | x_{0:t-1}) \quad \text{and} \quad \pi(y_t | x_{0:t}, y_{1:t-1}) \quad (6)$$

determine a state-space model. Note that the *a posteriori* pdf of a path  $X_{0:t} = x_{0:t}$  given a sequence of observations  $Y_{1:t} = y_{1:t}$ , denoted  $\pi(x_{0:t} | y_{1:t})$ , can be easily derived from the functions in Eq. (6) using the Bayes' theorem, namely

$$\begin{aligned} \pi(x_{0:t} | y_{1:t}) &\propto \pi(y_t | x_{0:t}, y_{1:t-1}) \pi(x_t | x_{0:t-1}) \pi(x_{0:t-1} | y_{1:t-1}) \\ &= \pi(x_0) \pi(y_1 | x_{0:1}) \pi(x_1 | x_0) \prod_{k=2}^t \pi(y_k | x_{0:k}, y_{1:k-1}) \pi(x_k | x_{0:k-1}). \end{aligned} \quad (7)$$

Since our ultimate interest is to find the values of the subsequence  $x_{0:T}$  that minimize  $C_T(\cdot; y_{1:T})$ , we need to establish a connection between the state-space model (6) and the defined cost function. The relationship is given by way of the posterior pdf in (7), according to the following definition.

**Definition 1** *Let  $y_{1:T}$  be a fixed sequence of observations. A state-space model determined by Eq. (6) is matched to the cost function  $C_T(x_{0:T}; y_{1:T})$  if, and only if,*

$$\arg \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi(x_{0:T} | y_{1:T}) = \arg \min_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} C_T(x_{0:T}; y_{1:T}). \quad (8)$$

For conciseness, we let  $X_t^\pi \triangleq \arg \max_{x_{0:t} \in (\mathbb{R}^{d_x})^{t+1}} \pi(x_{0:t} | y_{1:t})$  and rewrite Eq. (8) as  $X_T^c = X_T^\pi$ .

In many cases, Definition 1 turns out too generic for testing directly whether a state-space model and a cost function are matched. It is more useful to have sufficient conditions in terms of the basic building blocks of the model that allow for  $\mathsf{X}_T^c = \mathsf{X}_T^\pi$ , where the building blocks are, on one side, the densities  $\pi(x_0)$ ,  $\pi(x_t|x_{0:t-1})$  and  $\pi(y_t|x_{0:t}, y_{1:t-1})$  and, on the other side, the partial cost  $c_t(x_{0:t}, y_t)$  and the update function  $H(\cdot, \cdot)$ . Proposition 1 below provides such conditions. Note that, for conciseness of notation, in the sequel we adopt the shorthand

$$\pi_{0:t}(x_{0:t}) \triangleq \pi(x_{0:t}|y_{1:t}), \quad \mathsf{C}_{0:t}(x_{0:t}) \triangleq C_t(x_{0:t}; y_{1:t}) \quad \text{and} \quad \mathsf{c}_t(x_{0:t}) \triangleq c_t(x_{0:t}; y_t). \quad (9)$$

**Proposition 1** *Let  $y_{1:T}$  be an arbitrary, but fixed, sequence of observations and let the pdf's  $\pi(x_0)$ ,  $\pi_{0:t}(x_{0:t})$ ,  $\pi(y_t|x_{0:t}, y_{1:t-1})$ ,  $\pi(x_t|x_{0:t-1})$ , for  $t = 1, \dots, T$ , be proper. If we assume that:*

- (i) *There exists a monotonically decreasing function  $F_0 : \mathbb{R}^+ \rightarrow \mathbb{R}$  such that  $\mathsf{C}_0(x_0) = F_0(\kappa_0\pi(x_0))$ , with a proportionality constant  $\kappa_0$  independent of  $x_0$ .*
- (ii) *There exists a monotonically decreasing function  $F : \mathbb{R}^+ \rightarrow \mathbb{R}$  such that its inverse  $F^{-1}$  factorizes the cost function as*

$$F^{-1}(H(\mathsf{C}_{0:t-1}(x_{0:t-1}), \mathsf{c}_t(x_{0:t}))) = F^{-1}(\mathsf{C}_{0:t-1}(x_{0:t-1})) \times f(\mathsf{C}_{0:t-1}(x_{0:t-1}), \mathsf{c}_t(x_{0:t})), \quad (10)$$

where  $f : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$  is a nonnegative function.

- (iii) *The function  $f$  is such that*

$$f(\mathsf{C}_{t-1}(x_{0:t-1}), \mathsf{c}_t(x_{0:t})) \propto \pi_t(y_t|x_{0:t}, y_{1:t-1})\pi_t(x_t|x_{0:t-1}), \quad (11)$$

with a proportionality constant independent of  $x_{0:t}$ .

Then, for all  $t \in \{1, 2, \dots, T\}$ ,  $\mathsf{C}_{0:t}(x_{0:t}) = F(\kappa_t\pi_{0:t}(x_{0:t}))$ , for some constant  $\kappa_t$  independent of  $x_{0:t}$ , and  $\mathsf{X}_t^c = \mathsf{X}_t^\pi$ .

**Proof:** The proof proceeds by induction in  $t$ . At time  $t = 0$ , by assumption (i) we have  $\mathsf{C}_0(x_0) = F_0(\kappa_0\pi(x_0))$ , hence  $\mathsf{X}_0^c = \mathsf{X}_0^\pi$ . At time  $t - 1$ , we assume that  $\mathsf{C}_{0:t-1}(x_{0:t-1}) = F(\kappa_{t-1}\pi_{0:t-1}(x_{0:t-1}))$  and, as a consequence,  $\pi_{0:t-1}(x_{0:t-1}) \propto F^{-1}(\mathsf{C}_{0:t-1}(x_{0:t-1}))$ .

At time  $t$ , we apply assumption (ii) to obtain

$$\begin{aligned} F^{-1}(\mathsf{C}_{0:t}(x_{0:t})) &= F^{-1}(H(\mathsf{C}_{0:t-1}(x_{0:t-1}), \mathsf{c}_t(x_{0:t}))) \\ &= F^{-1}(\mathsf{C}_{0:t-1}(x_{0:t-1})) \times f(\mathsf{C}_{0:t-1}(x_{0:t-1}), \mathsf{c}_t(x_{0:t})). \end{aligned} \quad (12)$$

By the induction hypothesis,

$$F^{-1}(\mathbf{C}_{0:t-1}(x_{0:t-1})) \propto \pi_{0:t-1}(x_{0:t-1}) \quad (13)$$

and, from assumption (iii),

$$f(\mathbf{C}_{0:t-1}(x_{0:t-1}), \mathbf{c}_t(x_{0:t})) \propto \pi_t(y_t|x_{0:t}, y_{1:t-1})\pi_t(x_t|x_{0:t-1}), \quad (14)$$

hence we can substitute (13) and (14) into Eq. (12) to yield

$$F^{-1}(\mathbf{C}_{0:t}(x_{0:t})) \propto \pi_{0:t-1}(x_{0:t-1})\pi_t(y_t|x_{0:t}, y_{1:t-1})\pi_t(x_t|x_{0:t-1}) \propto \pi_{0:t}(x_{0:t}) \quad (15)$$

and, as a consequence,  $\mathbf{C}_{0:t}(x_{0:t}) = F(\kappa_t \pi_{0:t}(x_{0:t}))$  for some constant  $\kappa_t$  independent of  $x_{0:t}$  (but possibly dependent on  $y_{1:t}$ ), and  $\mathbf{X}_t^\pi = \mathbf{X}_t^c$ .  $\square$

**Remark 1** *Since  $F$  is monotonically decreasing,  $\mathbf{C}_{0:t}(x_{0:t}) = F(\kappa_t \pi_{0:t}(x_{0:t}))$  trivially implies  $\mathbf{X}_t^c = \mathbf{X}_t^\pi$ . The function  $f$  intuitively represents the mechanism that relates the update of the cost function,  $\mathbf{C}_{0:t-1}(x_{0:t-1})$ , with the update of the posterior pdf,  $\pi_{0:t-1}(x_{0:t-1})$ , at time  $t$ , when a new observation,  $y_t$ , is used.*

Most frequently, the functions  $F(\cdot)$  and  $f(\cdot, \cdot)$  in Proposition 1 are of the exponential class and they lead to state-space models that consist of exponential densities as well. This is illustrated by the two families of cost functions studied in Section 5. We conclude the present section by applying Proposition 1 in deriving a state-space model matched to the cost of Example 1.

**Example 2** *A toy problem (continued).*

The cost function of the toy Example 1 has the form  $\mathbf{C}_{0:T}(x_{0:T}) = \sum_{t=1}^T (|x_t| - 5)^2$ . Let  $F(a) = -\log(a)$  and apply the inverse function  $F^{-1}(a) = \exp\{-a\}$  to  $\mathbf{C}_{0:t}(x_{0:t})$ . We obtain

$$\begin{aligned} F^{-1}(\mathbf{C}_{0:t}(x_{0:t})) &= \exp\left\{-\mathbf{C}_{0:t-1}(x_{0:t-1}) - (|x_t| - 5)^2\right\} \\ &= \exp\{-\mathbf{C}_{0:t-1}(x_{0:t-1})\} \times \exp\left\{-(|x_t| - 5)^2\right\}, \end{aligned} \quad (16)$$

hence  $f(\mathbf{C}_{0:t-1}(x_{0:t-1}), \mathbf{c}_t(x_{0:t})) \propto \exp\left\{-(|x_t| - 5)^2\right\}$  and all that remains is to identify  $\pi(y_t|x_{0:t}, y_{1:t-1})$  (for  $y_t = 5$ ) and  $\pi(x_t|x_{0:t-1})$  such that  $\pi(y_t = 5|x_{0:t}, y_{1:t-1})\pi(x_t|x_{0:t-1}) \propto \exp\left\{-(|x_t| - 5)^2\right\}$ . But this is easily achieved by choosing a uniform transition pdf,  $\pi(x_t|x_{0:t-1}) = U(x_t; -10, +10)$ , and a Gaussian conditional density  $\pi(y_t|x_{0:t}, y_{1:t-1}) = \pi(y_t|x_t) = N(y_t; x_t, \frac{1}{2})$ . We complete the description of the matched state-space model by taking a uniform prior,  $\pi(x_0) = U(x_0; -10, +10)$ .



## 5 Examples

In this Section we illustrate the construction of state-space models matched to cost functions by way of two examples, each of them dealing with a class of update functions  $H(\cdot, \cdot)$ . The first one involves a purely additive rule,  $H(a, b) = a + b$ . There is a large number of problems that can be reduced to this form, including the so-called “discounted costs” [30] often applied in finance. Then we study a nonlinear update rule of the form  $H(a, b) = \max(a, b)$ . In the two cases, we explicitly show the functions  $F(\cdot)$  and  $f(\cdot, \cdot)$  that relate the cost  $C_{0:t}(x_{0:t})$  to the posterior pdf,  $\pi_{0:t}(x_{0:t})$ , according to Proposition 1.

### 5.1 Additive cost

Additive costs appear frequently in scientific and engineering problems, e.g., positioning and navigation [31], finance [34] or operational research [3]. Let us consider the generic additive form  $C_{0:t}(x_{0:t}) = C_{0:t-1}(x_{0:t-1}) + c_t(x_{0:t})$ . This cost can be related to a posterior pdf easily by means of the monotonically decreasing functions  $F_0(a) = F(a) = -\log(a)$  and  $f(a, b) = \exp\{-b\}$ , which yield

$$\begin{aligned} F^{-1}(C_{0:t-1}(x_{0:t-1}) + c_t(x_{0:t})) &= F^{-1}(C_{0:t-1}(x_{0:t-1})) \times f(C_{0:t-1}(x_{0:t-1}), c_t(x_{0:t})) \\ &= \exp\{-C_{0:t-1}(x_{0:t-1})\} \times \exp\{-c_t(x_{0:t})\} \\ &= \exp\left\{-C_0(x_0) - \sum_{k=1}^t c_k(x_{0:k})\right\} \propto \pi_{0:t}(x_{0:t}). \end{aligned} \quad (17)$$

For this formal decomposition to be valid, we require integrability of the terms  $F^{-1}(C(x_0))$  and  $f(c_k(x_{0:k}))$ , i.e.,  $\int \exp\{-C(x_0)\} dx_0 < \infty$  and  $\int \exp\{-c_k(x_{0:k})\} dx_k < \infty$ , for each  $k \in \{1, \dots, t\}$ .

**Example 3** *Neumaier 3 problem.*

The Neumaier 3 problem is included in the collection of [1] and consists in the minimization of the cost function

$$J(x_{1:T}) = \sum_{t=1}^T (x_t - 1)^2 - \sum_{t=2}^T x_t x_{t-1}, \quad \text{subject to} \quad -T^2 \leq x_t \leq T^2, \quad t \in \{1, \dots, T < \infty\}. \quad (18)$$

The number of local minima of  $J(x_{1:T})$  is not known, but the global minimum can be expressed as

$$J(x_{1:T}^o) = -\frac{T(T+4)(T-1)}{6}, \quad \text{where} \quad x_t^o = t(T+1-t), \quad t = 1, \dots, T. \quad (19)$$

The cost  $J(x_{1:T})$  has a linear additive form. Specifically, we adapt it to the notation in this paper by defining

$$C_{0:T}(x_{0:T}) = \frac{1}{\sigma^2} \left[ \sum_{t=1}^T (x_t - y_t)^2 - \sum_{t=2}^T x_t x_{t-1} \right], \quad (20)$$

where  $y_t = 1$  for all  $t \geq 1$  and  $\sigma^2 > 0$  is an arbitrary scale parameter. Note that, subject to  $-T^2 \leq x_t \leq T^2$ ,

$$\arg \min_{x_{1:T}} J(x_{1:T}) = \arg \min_{x_{1:T}} C_{0:T}(x_{0:T}), \quad (21)$$

and  $x_0$  is a dummy unknown included only for notational compatibility. The functions  $C_{0:t}$ ,  $t = 2, \dots, T$ , admit the recursive decomposition

$$C_{0:t}(x_{0:t}) = C_{0:t-1}(x_{0:t-1}) + \frac{1}{\sigma^2} [(x_t - y_t)^2 - x_t x_{t-1}]. \quad (22)$$

Therefore the posterior pdf at time  $t \geq 2$  for the matched state-space model has the form

$$\begin{aligned} \pi_{0:t}(x_{0:t}) &\propto \exp\{-C_{0:t}(x_{0:t})\} \\ &= \exp\{-C_{0:t-1}(x_{0:t-1})\} \exp\left\{-\frac{1}{\sigma^2}(x_t - y_t)^2\right\} \exp\left\{\frac{1}{\sigma^2}x_t x_{t-1}\right\}, \end{aligned} \quad (23)$$

while  $\pi(x_{0:1}|y_1) \propto \exp\{-\frac{1}{\sigma^2}(x_1 - y_1)^2\}$  and  $\pi(x_0) = U(x_0; -T^2, +T^2)$  (the value of  $x_0$  does not affect the cost or the pdf, hence the uniform distribution). The resulting likelihood and transition density at time  $t$  are

$$\pi(y_t|x_{0:t}, y_{1:t-1}) = \pi(y_t|x_t) \propto \exp\left\{-\frac{1}{\sigma^2}(x_t - y_t)^2\right\}, \quad t \geq 1 \quad (24)$$

$$\pi(x_t|x_{0:t-1}) = \pi(x_t|x_{t-1}) \propto \exp\left\{\frac{1}{\sigma^2}x_t x_{t-1}\right\}, \quad t \geq 2, \quad (25)$$

respectively, while  $\pi(x_1|x_0) = U(x_1; -T^2, T^2)$ . Note that  $\pi(y_t|x_t) = N(y_t; x_t, \frac{\sigma^2}{2})$  and the normalization constant for  $\pi(x_t|x_{t-1})$  ( $t > 1$ ) is

$$\kappa_t = \frac{x_{t-1}}{\sigma^2} \left( \exp\left\{\frac{T^2 x_{t-1}}{\sigma^2}\right\} - \exp\left\{-\frac{T^2 x_{t-1}}{\sigma^2}\right\} \right)^{-1}. \quad (26)$$

## 5.2 Minimax problems

Optimization problems that consist in the minimization of the maximum value of a certain function abound in engineering, finance and other disciplines (see, e.g., [10, 28, 25]). Let  $a \vee b$  denote the maximum of  $a$  and  $b$ . In this example, we study cost functions of the form  $C_t(x_{0:t}) = C_{t-1}(x_{0:t-1}) \vee c_t(x_{0:t})$ . It turns out that this kind of cost can also be factorized by means of the usual update rule  $F_0(a) = F(a) = -\log(a)$ . Indeed,

$$\begin{aligned} F^{-1}(C_{t-1}(x_{0:t-1}) \vee c_t(x_{0:t})) &= \exp\{- (C_{t-1}(x_{0:t-1}) \vee c_t(x_{0:t}))\} \\ &= \frac{\exp\{-C_{t-1}(x_{0:t-1})\} \times \exp\{-c_t(x_{0:t})\}}{\exp\{-C_{t-1}(x_{0:t-1})\} \vee \exp\{-c_t(x_{0:t})\}}. \end{aligned} \quad (27)$$

By inspection of (27) we quickly notice that

$$f(C_{t-1}(x_{0:t-1}), c_t(x_{t-1:t})) = \frac{\exp\{-c_t(x_{0:t})\}}{\exp\{-C_{t-1}(x_{0:t-1})\} \vee \exp\{-c_t(x_{0:t})\}}. \quad (28)$$

Again, we assume that  $\int \exp\{-C_{k-1}(x_{0:k-1})\} dx_{0:k-1} < \infty$  and  $\int \exp\{-c_k(x_{0:k})\} dx_k < \infty$  for all  $k$ .

**Example 4** *Cross-talk cancellation.*

As shown in [28], the problem of designing an acoustic filter for cross-talk cancellation in a 3D audio system can be stated as a minimax problem. Indeed, let  $h_a(n)$ ,  $n \in \mathbb{Z}$ , be a sequence that represents the combined effect of the acoustic impulse responses between the sound sources (loudspeakers) and (say) the listener's left ear and let  $h_f(n)$ ,  $n \in \mathbb{Z}$ , be the cross-talk cancellation filter that should let the desired source signal pass while mitigating all other signals coming from different sources (see [28] for details). The impulse response  $h_a(n)$  is causal with length  $2M - 1$ , i.e.,  $h_a(n) = 0$  for all  $n < 0$  and  $n \geq 2M - 1$ , while the filter  $h_f(n)$  is assumed causal with length  $K$ , i.e.,  $h_f(n) = 0$  for all  $n < 0$  and  $n \geq K$ .

The goal is to find the response  $h_f(n)$  such that the convolution  $c(n) = h_a(n) * h_f(n) = \sum_{k=0}^{2M-2} h_a(k)h_f(n-k)$  is the closest to a desired response

$$d(n) = \begin{cases} 1 & \text{if } n = 0 \\ 0 & \text{otherwise} \end{cases}, \quad (29)$$

i.e., the filter  $h_f(n)$  is selected to invert the combined acoustic response  $h_a(n)$ . Perfect inversion is not possible, since  $h_f(n)$  has a finite length, hence we seek to solve the equations  $d(n) - c(n) = 0$ ,  $n = 0, \dots, K + 2M - 3$ , approximately instead. Let us collect the complete set of filter coefficients into the vector  $h_f = [h_f(0), \dots, h_f(K-1)]^\top \in \mathbb{R}^K$ . In [28] it is proposed to compute  $h_f$  as the solution of the minimax problem.

$$\hat{h}_f = \arg \min_{h_f \in \mathbb{R}^K} \left\{ J(h_f) = \max_{n \in \{0, \dots, K+2M-3\}} \left| d(n) - \sum_{k=0}^{2M-2} h_a(k)h_f(n-k) \right| \right\}. \quad (30)$$

We can easily rewrite problem (30) using our notation. For the unknowns, we let  $x_t = h_f(t-1) \in \mathbb{R}$  for  $t = 0, 1, \dots, K$  (hence  $x_0 = 0$  and  $x_{t>K} = 0$ ). The desired sequence  $d(n)$  plays the role of the observations, hence  $y_t = d(t-1)$ ,  $t = 1, 2, \dots, K + 2M - 2$ . We define the partial cost at time  $t \geq 1$  as  $c_t(x_{0:t}) = |y_t - \sum_{k=0}^{2M-2} h_a(k)x_{t-k+1}|$  while, trivially,  $C_0(x_0) = 0$ . The overall cost at time  $t$  then becomes  $C_{0:t}(x_{0:t}) = C_{0:t-1}(x_{0:t-1}) \vee c_t(x_{0:t})$ . The time horizon is  $T = K + 2M - 2$  and  $C_{0:T}(x_{0:T}) = J(h_f)$ .

To construct the state-space model matched to  $C_{0:T}(x_{0:T})$  we use Eqs. (27) and (28). Specifically, we

have to choose  $\pi(y_t|x_{0:t}, y_{1:t-1})$  and  $\pi(x_t|x_{0:t-1})$  such that

$$f(\mathbf{C}_{0:t-1}(x_{0:t-1}), \mathbf{c}_t(x_{0:t})) = \frac{\exp\{-\mathbf{c}_t(x_{0:t})\}}{\exp\{-\mathbf{C}_{0:t-1}(x_{0:t-1})\} \vee \exp\{-\mathbf{c}_t(x_{0:t})\}} \propto \pi(y_t|x_{0:t}, y_{1:t-1})\pi(x_t|x_{0:t-1}). \quad (31)$$

The transition pdf  $\pi(x_t|x_{0:t-1})$  is selected to be uniform over the set where  $X_t$  is allowed to take values, i.e.,  $\pi(x_t|x_{0:t-1}) = \pi(x_t) = U(x_t; -a, +a)$  for some prescribed bound  $a > 0$ . The system definition is completed with  $\pi(y_t|x_{0:t}, y_{1:t-1}) \propto f(\mathbf{C}_{0:t-1}(x_{0:t-1}), \mathbf{c}_t(x_{0:t}))$ , while  $\pi_0(x_0) = 0$  with probability 1.

## 6 Algorithms

We have recast the minimization of  $\mathbf{C}_{0:T}(x_{0:T})$  into a problem of MAP estimation for a matched state-space model. However, this is also intractable for most models of practical interest (linear Gaussian systems being the exception) and we need to resort to numerical techniques in order to find the solutions. We propose the use of sequential Monte Carlo (SMC) methods to build a particle approximation of the *a posteriori* smoothing probability measure, from which MAP estimates can be computed. Specifically, we can employ the standard SIR algorithm [16, 9] to obtain a discretization of the path space  $(\mathbb{R}^{d_x})^{T+1}$  consisting of a set of  $N$  paths  $\{x_{0:T}^{(n)}\}_{n=1, \dots, N}$ , and then choose the realization with the highest posterior density. This method can be inefficient in some problems, though, as will be shown in Section 7.1. A more efficient MAP estimation technique can be obtained by combining the SIR and Viterbi [11] algorithms, as suggested in [14]. We prove that both methods guarantee almost sure asymptotic convergence and obtain a lower bound for the necessary number of particles in the discretization of the state-space as a function of the desired accuracy of the approximation, the dimension  $d_x$  and the time horizon  $T$ .

### 6.1 Discretization of the state-space: sequential importance resampling algorithm

We aim at numerically computing solutions of the MAP estimation problem

$$\mathbf{X}_T^\pi = \arg \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T}). \quad (32)$$

Even if the posterior pdf  $\pi_{0:T}(x_{0:T})$  can be evaluated up to a proportionality constant (using the factorization of Eq. (7)) this is, in general, a difficult optimization problem in a high dimensional space, possibly with multiple global and/or local extrema. In this paper we propose to tackle these difficulties by using a SMC method in order to obtain a suitable discretization of the path space  $(\mathbb{R}^{d_x})^{T+1}$ . Different

search methods can subsequently be applied to find the point of the discretized space with the highest density.

SMC algorithms [19, 8] are aimed at recursively computing approximations of the sequence of posterior probability laws

$$\Pr \{A|Y_{1:t} = y_{1:t}\} = \int_A \pi_{0:T}(x_{0:T}) dx_{0:T}, \quad t = 1, 2, \dots, T, \quad (33)$$

where  $A \in \mathcal{B} \left( (\mathbb{R}^{d_x})^{T+1} \right)$  is a Borel set. Specifically, at each time  $t$ , a SMC algorithm generates a random sample  $\Omega_{0:t}^N = \left\{ x_{0:t}^{(n)} \right\}_{n=1, \dots, N}$  such that integrals with respect to the pdf  $\pi_{0:t}(x_{0:t})$  can be approximated by summations [5], i.e.,  $\int f(x_{0:t}) \pi_{0:t}(x_{0:t}) dx_{0:t} \approx \frac{1}{N} \sum_{n=1}^N f(x_{0:t}^{(n)})$ , where  $f : (\mathbb{R}^{d_x})^{t+1} \rightarrow \mathbb{R}$  is a Borel measurable function defined in the path space and integrable with respect to the posterior probability law.

Although various possibilities exist [8], in this paper we consider the standard sequential importance sampling algorithm with resampling at every time step [9], also known as bootstrap filter [16]. We refer to this algorithm as SIR through the the paper. The algorithm is based on the recursive decomposition of  $\pi_{0:T}(x_{0:T})$  given by Eq. (7) and the computational procedure is simple.

- *Initialization.* At time  $t = 0$  we draw  $N$  independent and identically distributed (i.i.d.) samples from the prior probability distribution with density  $\pi_0(x_0)$ . Let us denote this initial sample as  $\Omega_0^N = \{x_0^{(n)}\}_{n=1, \dots, N}$ .
- *Recursive step.* Assume that a random sample  $\Omega_{0:t-1}^N = \{x_{0:t-1}^{(n)}\}_{n=1, \dots, N}$  has been generated up to time  $t - 1$ . Then, at time  $t$ , we take the following steps.
  - i. Draw  $N$  new samples in the state space  $\mathbb{R}^{d_x}$  from the probability distributions with densities  $\pi(x_t|x_{0:t-1}^{(n)})$ ,  $n = 1, \dots, N$ , and denote them as  $\{\bar{x}_t^{(n)}\}_{n=1, \dots, N}$ . Set  $\bar{x}_{0:t}^{(n)} = \{x_{0:t-1}^{(n)}, \bar{x}_t^{(n)}\}$ .
  - ii. Weight each sample according to its likelihood, i.e., compute importance weights  $\tilde{w}_t^{(n)} = \pi(y_t|\bar{x}_{0:t}^{(n)}, y_{1:t-1})$  and normalize them to obtain  $w_t^{(n)} = \tilde{w}_t^{(n)} / \sum_{k=1}^N \tilde{w}_t^{(k)}$ .
  - iii. Resampling: for  $n = 1, \dots, N$ , set  $x_{0:t}^{(n)} = \bar{x}_{0:t}^{(k)}$  with probability  $w_t^{(k)}$ ,  $k \in \{1, \dots, N\}$ . Reset the weights,  $w_t^{(n)} = 1/N$  for  $n = 1, \dots, N$ .

The multinomial resampling procedure in step *iii.* can be substituted by other techniques [4, 6]. We shall use the random grid  $\Omega_{0:T}^N = \{x_{0:T}^{(n)}\}_{n=1, \dots, N}$  as a discrete approximation of the path space  $(\mathbb{R}^{d_x})^{T+1}$  where the random sequence  $X_{0:T}$  takes its values. Note that the SIR algorithm also yields “marginal grids” for each time  $t$ , denoted  $\Omega_t^N = \{x_t^{(n)}\}_{n=1, \dots, N}$ ,  $t = 0, 1, \dots, T$ .

The points of the grid  $\Omega_{0:T}^N$  (often also the points of every  $\Omega_t^N = \{x_t^{(n)}\}_{n=1,\dots,N}$ ) are called particles and the SMC methods that generate them are referred to as particle filters [9] or particle smoothers [15] depending on whether one is interested in the filtering pdf's  $\pi(x_t|y_{1:t})$ ,  $t = 1, 2, \dots$ , or the smoothing pdf's  $\pi_{0:t}(x_{0:t}) = \pi(x_{0:t}|y_{1:t})$ ,  $t = 1, 2, \dots$ , respectively. Using the particles in  $\Omega_{0:T}^N$ , it is straightforward to build a random measure  $\pi_{0:t}^N(dx_{0:t}) = \frac{1}{N} \sum_{n=1}^N \delta_n(dx_{0:t})$ , where  $\delta_n$  is the unit delta measure centered at  $x_{0:T}^{(n)}$ , and use it to approximate integrals of the form

$$(f, \pi_{0:t}) \triangleq \int f(x_{0:t}) \pi_{0:t}(x_{0:t}) dx_{0:t}, \quad (34)$$

where  $f : (\mathbb{R}^{d_x})^{t+1} \rightarrow \mathbb{R}$  is a real function in the space of the paths up to time  $t$ . In this paper, however, we are interested in the densities  $\pi_{0:t}(x_{0:t})$ , rather than the probability measures from which they are derived, hence we find it more convenient to interpret  $\pi_{0:t}^N(x_{0:t})$  as a point-mass approximation of  $\pi_{0:t}(x_{0:t})$ . Thus, we slightly abuse the notation in (34) to write

$$(f, \pi_{0:t}^N) = \int f(x_{0:t}) \pi_{0:t}^N(dx_{0:t}) = \frac{1}{N} \sum_{n=1}^N f_t(x_{0:t}^{(n)}). \quad (35)$$

If the function  $f$  is, for example, bounded, then  $(f, \pi_{0:t}^N)$  is a good approximation of  $(f, \pi_{0:t})$  for  $N$  sufficiently large [5]. We will take advantage of this result in Section 6.3.

## 6.2 MAP estimation algorithms

We propose to use the random grids generated by the SIR algorithm to search for approximate maximizers of the pdf  $\pi_{0:T}(x_{0:T})$ . In particular, we investigate two algorithms. The first one is a straightforward extension of the SIR procedure, while the second one combines it with the Viterbi algorithm as suggested in [14]. We will subsequently refer to them as Algorithm 1 and Algorithm 2, respectively.

### 6.2.1 Algorithm 1

The idea is as simple as to search the element of  $\Omega_{0:T}^N$  with the highest posterior density. For this purpose, note that we can easily extend the SIR algorithm described in Section 6.1 to recursively compute the posterior density of each particle up to a proportionality constant. To be specific, we need to perform the following additional computations.

- At the initialization step, let  $a_0^{(n)} = \log(\pi(x_0^{(n)}))$  for  $n = 1, \dots, N$ .
- At the recursive step, modify steps *ii.* and *iii.* as follows.

- ii. Weight each sample according to its likelihood, i.e., compute importance weights  $\tilde{w}_t^{(n)} = \pi(y_t | \bar{x}_{0:t}^{(n)}, y_{1:t-1})$  and normalize them to obtain  $w_t^{(n)} = \tilde{w}_t^{(n)} / \sum_{k=1}^N \tilde{w}_t^{(k)}$ . Compute  $\bar{a}_t^{(n)} = a_{t-1}^{(n)} + \log\left(\pi(y_t | \bar{x}_{0:t}^{(n)}, y_{1:t-1})\right) + \log\left(\pi(\bar{x}_t^{(n)} | x_{0:t-1}^{(n)})\right)$ .
- iii. Resampling: for  $n = 1, \dots, N$ , set  $x_{0:t}^{(n)} = \bar{x}_{0:t}^{(k)}$  and  $a_t^{(n)} = \bar{a}_t^{(k)}$  with probability  $w_t^{(k)}$ ,  $k \in \{1, \dots, N\}$ . Reset the weights,  $w_t^{(n)} = 1/N$  for  $n = 1, \dots, N$ .

Then we select  $\hat{x}_{0:T}^N = x_{0:T}^{(n_o)}$ , where  $n_o = \arg \max_{n \in \{1, \dots, N\}} a_T^{(n)}$ , as the approximate maximizer of  $\pi_{0:T}(x_{0:T})$ .

### 6.2.2 Algorithm 2

In this Section we briefly describe the MAP estimation algorithm of [14]. Instead of  $\Omega_{0:T}^N$ , we consider now a finer discretization of  $(\mathbb{R}^{d_x})^{T+1}$ , namely the product space  $\bar{\Omega}_{0:T}^N = \Omega_0^N \times \bar{\Omega}_1^N \times \dots \times \bar{\Omega}_T^N$ , where  $\Omega_0^N = \{x_0^{(n)}\}_{n=1, \dots, N}$  and  $\bar{\Omega}_t^N = \{\bar{x}_t^{(n)}\}_{n=1, \dots, N}$  for  $t = 1, 2, \dots, T$ . Specifically note that  $\bar{\Omega}_t^N$  is constructed from the particles available at step *ii.* of the SIR algorithm, i.e., before resampling, to avoid duplicate samples.

Next, assume that  $\pi(x_t | x_{0:t-1}) = \pi(x_t | x_{t-1})$  (i.e., the state-space system is Markovian) and  $\pi(y_t | x_{0:t}, y_{1:t-1}) = \pi(y_t | x_t)$ . Given the random grids  $\Omega_0^N, \bar{\Omega}_t^N, t = 1, \dots, T$ , the Viterbi algorithm outputs a sequence  $(x_0^{(n_o)}, \bar{x}_1^{(n_1)}, \dots, \bar{x}_T^{(n_T)}) \in \bar{\Omega}_{0:T}^N$ , where  $n_i \in \{1, \dots, N\} \forall i$ , with the highest posterior density, i.e., it solves the problem of finding

$$\bar{x}_{0:T}^N \in \arg \max_{\bar{x}_{0:T} \in \bar{\Omega}_{0:T}^N} \pi_{0:T}(\bar{x}_{0:T}) \quad (36)$$

exactly. The procedure is described below.

- *Initialization.* For  $n = 1, \dots, N$ , let  $a_0^{(n)} = \log(\pi(x_0^{(n)}))$ .
- *Recursive step.* At time  $t > 0$ , the random grids  $\bar{\Omega}_{t-1}^N$  and  $\bar{\Omega}_t^N$ , as well as  $\{a_{t-1}^{(n)}\}_{n=1, \dots, N}$ , are available. Then, for  $n = 1, \dots, N$ , compute
  - i.  $a_t^{(n)} = \log(\pi(y_t | \bar{x}_t^{(n)})) + \max_{k \in \{1, \dots, N\}} \left[ a_{t-1}^{(k)} + \log(\pi(\bar{x}_t^{(n)} | \bar{x}_{t-1}^{(k)})) \right]$ ,
  - ii.  $\ell_t^{(n)} = \arg \max_{k \in \{1, \dots, N\}} \left[ a_{t-1}^{(k)} + \log(\pi(\bar{x}_t^{(n)} | \bar{x}_{t-1}^{(k)})) \right]$ .
- *Backtracking.* Computation of an optimal sequence.
  - i. At time  $T$ , let  $j_T = \arg \max_{k \in \{1, \dots, N\}} a_T^{(k)}$  and assign  $\bar{x}_T^N = \bar{x}_T^{(j_T)}$ .
  - iii. For  $t = T-1, T-2, \dots, 0$ , let  $j_t = \ell_{t+1}^{(j_{t+1})}$  and assign  $\bar{x}_t^N = \bar{x}_t^{(j_t)}$ .

The computational complexity of the method is  $\mathcal{O}(N^2)$ . Note that the Viterbi recursion can be run sequentially, together with the SIR algorithm described in Section 6.1. Specifically, we can take a complete recursive step of the Viterbi algorithm right after step *ii.* of the SIR method (i.e., once the random marginal grid  $\bar{\Omega}_t^N$  is obtained). The combination of the SIR and Viterbi methods to compute  $\bar{x}_{0:T}^N$  will be termed Algorithm 2 in the sequel.

### 6.3 Convergence analysis

We now establish the almost sure convergence of the two MAP estimation algorithms described in Section 6.2. In the results that follow, we assume that:

- The sequence  $Y_{1:T} = y_{1:T}$  is fixed (not random).
- The likelihoods  $g_t(x_{0:t}) = \pi(y_t | x_{0:t}, y_{1:t-1})$  are bounded functions of  $x_{0:t}$  for every  $t = 1, 2, \dots, T$ .
- The posterior pdf  $\pi_{0:T}(x_{0:T})$  is uniformly continuous at every point  $\hat{x}_{0:T} \in \mathbf{X}_t^\pi$ .

The first two assumptions are applied to show that the SIR algorithm converges in an adequate way while the third one is used directly in the proof of Theorem 1 below.

Obviously, the convergence of the MAP estimation Algorithms 1 and 2 relies on the convergence of the SIR algorithm. To be precise, given a bounded function  $f : (\mathbb{R}^{d_x})^{T+1}$  our analysis requires the convergence of  $(f, \pi_{0:T}^N)$  toward the actual integral  $(f, \pi_{0:T})$  in the  $L_4$ -norm. Similar, but not directly applicable, results exist in the literature, e.g.,  $L_2$  bounds for the rate of convergence of  $(f, \pi_{0:T}^N)$  to  $(f, \pi_{0:T})$  can be found in [5], while  $L_p$  bounds for the rate of convergence of the corresponding marginals  $(f, \pi_T^N)$  to  $(f, \pi_T)$  in the state space of  $X_t$ , were established in [22] under additional constraints.

In the following Lemma 1, we establish the  $L_p$  bounds for the rate of convergence of  $(f, \pi_{0:T}^N)$  to  $(f, \pi_{0:T})$ . This is required for the subsequent analysis. In Theorem 1, we use the result to prove that Algorithm 1 converges almost surely. More precisely we prove that  $\pi_{0:T}(\hat{x}_{0:T}^N) \rightarrow \pi_{0:T}(\hat{x}_{0:T})$ , with  $\hat{x}_{0:T} \in \mathbf{X}_T^\pi$ . The convergence of Algorithm 2 follows immediately (see Corollary 1). Finally we establish a lower bound on the number of particles  $N$  needed to achieve a certain accuracy.

In the sequel,  $\|\xi\|_p$  denotes the  $L_p$  norm of the random variable  $\xi$  defined as  $\|\xi\|_p = E[\xi^p]^{1/p}$ ,  $E[\cdot]$  denotes the mathematical expectation over all possible realizations of the random measure  $\pi_{0:T}^N$  and  $\|f\|_\infty = \sup_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} |f(x_{0:T})|$  denotes the supremum norm of the function  $f$ .

**Lemma 1** *Let  $f : (\mathbb{R}^{d_x})^{T+1} \rightarrow \mathbb{R}$  be a bounded function defined on the path space. Then there exists a*



constant  $c = c(p, T, y_{1:T})$  independent of  $N$  such that

$$\|(f, \pi_{0:T}^N) - (f, \pi_{0:T})\|_p \leq \frac{c\|f\|_\infty}{\sqrt{N}},$$

for all  $N \geq 1$ .

See the Appendix for a proof.

**Remark 2** Lemma 1 holds true for unbounded functions  $f$  too. In this case, we need to assume that  $|f|^p$  is integrable with respect to the prior distribution<sup>2</sup> and the rate of convergence has the form  $\|(f, \pi_{0:T}^N) - (f, \pi_{0:T})\|_p \leq c/\sqrt{N}$ , where the constant  $c = c(p, T, y_{1:T}, f)$  is independent of  $N$ .

**Theorem 1** Let  $\hat{x}_{0:T}^N$  be the output sequence of Algorithm 1. Then, almost surely,

$$\lim_{N \rightarrow \infty} \pi_{0:T}(\hat{x}_{0:T}^N) = \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T})$$

Moreover any convergent subsequence of  $\hat{x}_{0:T}^N$  has a limit  $\hat{x}_{0:T}$  that belongs to the critical set  $\mathcal{X}_T^\pi$ .

**Proof:** Let  $f : (\mathbb{R}^{d_x})^{T+1} \rightarrow \mathbb{R}$  be a bounded real function of the path  $x_{0:T}$ . From Lemma 1, we obtain

$$\|(f, \pi_{0:T}^N) - (f, \pi_{0:T})\|_p \leq \frac{c\|f\|_\infty}{\sqrt{N}}, \quad (37)$$

The bound (37) for  $p = 4$  implies, using a standard argument, that there exists a positive random variable  $c_T^\epsilon$  such that, for all  $N > 0$ , we have, almost surely,

$$\left( (f, \pi_{0:T}^N) - (f, \pi_{0:T}) \right)^4 \leq \frac{c_T^\epsilon}{N^{1-\epsilon}} \quad (38)$$

for any arbitrarily small  $\epsilon > 0$ . As a consequence, the integral  $(f, \pi_{0:T}^N)$  converges almost surely, i.e.,

$$\lim_{N \rightarrow \infty} |(f, \pi_{0:T}^N) - (f, \pi_{0:T})| = 0. \quad (39)$$

Now, choose any MAP estimate  $\hat{x}_{0:T} \in \mathcal{X}_T^\pi$  and consider the open ball

$$B_k(\hat{x}_{0:T}) = \left\{ z \in (\mathbb{R}^{d_x})^{T+1} : \|z - \hat{x}_{0:T}\| < \frac{1}{k} \right\} \quad (40)$$

where  $k$  is a positive integer and  $\|\cdot\|$  denotes the norm of the Euclidean space  $(\mathbb{R}^{d_x})^{T+1}$ . The indicator function

$$I_{B_k(\hat{x}_{0:T})}(x_{0:T}) = \begin{cases} 1 & \text{if } x_{0:T} \in B_k(\hat{x}_{0:T}) \\ 0 & \text{otherwise} \end{cases} \quad (41)$$

---

<sup>2</sup>Implicitly,  $|f|^p$  is also integrable with respect to the posterior distribution since the likelihood function is bounded. Otherwise, the integrability with respect to the posterior distribution has to be assumed.

is real and bounded, hence

$$\lim_{N \rightarrow \infty} |(I_{B_k(\hat{x}_{0:T})}, \pi_{0:T}^N) - (I_{B_k(\hat{x}_{0:T})}, \pi_{0:T})| = 0 \quad (42)$$

almost surely. Since the posterior pdf  $\pi_{0:T}(x_{0:T})$  is uniformly continuous at  $\hat{x}_{0:T} \in \mathbb{X}_t^\tau$  and  $\pi_{0:T}(\hat{x}_{0:T}) > 0$  it follows that  $\pi_{0:T}(x_{0:T})$  is positive on an open ball around  $\hat{x}_{0:T}$ . In particular, the value  $A_k = (I_{B_k(\hat{x}_{0:T})}, \pi_{0:T})$  is strictly positive. Also

$$(I_{B_k(\hat{x}_{0:T})}, \pi_{0:T}^N) = \frac{m(N, k)}{N}, \quad (43)$$

where  $m(N, k)$  denotes the number of elements of the discretized path space  $\Omega_{0:T}^N$  that belong to the ball  $B_k(\hat{x}_{0:T})$  (equivalently,  $m(N, k) = |\Omega_{0:T}^N \cap B_k(\hat{x}_{0:T})|$  is the number of points in the discrete intersection set  $\Omega_{0:T}^N \cap B_k(\hat{x}_{0:T})$ ). Since  $\lim_{N \rightarrow \infty} \left| \frac{m(N, k)}{N} - A_k \right| = 0$  for any  $k$ , it follows that  $m(N, k) > 0$  almost surely whenever  $N$  is sufficiently large.

Since, for any integer  $k > 0$ , the intersection  $\Omega_{0:T}^N \cap B_k(\hat{x}_{0:T})$  is nonempty for  $N$  sufficiently large (almost surely), we can choose a point  $x_{0:T}^{N, k} \in \Omega_{0:T}^N \cap B_k(\hat{x}_{0:T})$ . Then  $\pi_{0:T}(x_{0:T}^{N, k}) \leq \pi_{0:T}(\hat{x}_{0:T})$  but, by construction of Algorithm 1, we also have that  $\pi_{0:T}(x_{0:T}^{N, k}) \leq \pi_{0:T}(\hat{x}_{0:T}^N)$ . Therefore,  $\pi_{0:T}(x_{0:T}^{N, k}) \leq \pi_{0:T}(\hat{x}_{0:T}^N) \leq \pi_{0:T}(\hat{x}_{0:T})$ . Since  $\pi_{0:T}$  is continuous at  $\hat{x}_{0:T}$  and  $\lim_{k \rightarrow \infty} \hat{x}_{0:T}^{N, k} = \hat{x}_{0:T}$  (as  $\|\hat{x}_{0:T} - \hat{x}_{0:T}^{N, k}\| < 1/k$ ), we deduce that  $\lim_{k \rightarrow \infty} \pi_{0:T}(x_{0:T}^{N, k}) = \pi_{0:T}(\hat{x}_{0:T})$ . Hence, also  $\lim_{k \rightarrow \infty} \pi_{0:T}(\hat{x}_{0:T}^N) = \pi_{0:T}(\hat{x}_{0:T})$ . Moreover, if  $\{\hat{x}_{0:T}^{N_i}\}_{i \in \mathbb{N}}$  is a convergent subsequence of  $\{\hat{x}_{0:T}^N\}_{N \in \mathbb{N}}$  with limit, say,  $\check{x}_{0:T}$  it follows that  $\pi_{0:T}(\check{x}_{0:T}) = \lim_{i \rightarrow \infty} \pi_{0:T}(\hat{x}_{0:T}^{N_i}) = \pi_{0:T}(\hat{x}_{0:T})$ . Therefore  $\check{x}_{0:T} \in \mathbb{X}_T^\tau$  which concludes the proof.  $\square$

**Remark 3** In [23] a different approach is used to prove a similar result to Theorem 1 based on the propagation of chaos property of genealogical tree simulations models (see [20] for details). The basic idea is that a sub-sample from  $\{x_{0:T}^{(i)}\}_{i=1, \dots, N}$  behaves asymptotically as a perfect sample from  $\pi_{0:T}$ . More precisely, using Theorem 8.3.3 in [20] one can show that if  $\pi_{0:T}^{\otimes q}$  is the tensor product of  $q$  copies of the measure  $\pi_{0:T}$ , then

$$\|\text{Law}(x_{0:T}^{(1)}, x_{0:T}^{(2)}, \dots, x_{0:T}^{(q)}) - \pi_{0:T}^{\otimes q}\|_{tv} \leq \frac{q^2}{N} c(T), \quad (44)$$

where  $\|\cdot\|_{tv}$  is the total variation norm between two probability measures and  $c(T)$  is a constant with respect to  $N$ . By choosing  $q=q(N)$  to be of order  $o(N)$  one can show that, for any  $\delta > 0$ ,

$$\Pr \left( \max_{i=1, \dots, q(N)} \pi_{0:T}(x_{0:T}^{(i)}) < \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T}) - \delta \right) \leq c(T) \frac{q(N)^2}{N} + \pi_{0:T}(A(\delta))^{q(N)}, \quad (45)$$

where  $A(\delta)$  is defined to be the set

$$A(\delta) = \left\{ x_{0:T} \in (\mathbb{R}^{d_x})^{T+1} : \pi_{0:T}(x_{0:T}) < \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T}) - \delta \right\}. \quad (46)$$

This, in turn, leads to the convergence in probability (and not almost surely) of the estimator to  $\max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T})$ .

**Corollary 1** Let  $\bar{x}_{0:T}^N$  be the output sequence of Algorithm 2. Then, almost surely,

$$\lim_{N \rightarrow \infty} \pi_{0:T}(\bar{x}_{0:T}^N) = \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T})$$

Moreover any convergent subsequence of  $\bar{x}_{0:T}^N$  has a limit  $\hat{x}_{0:T}$  that belongs to the critical set  $\mathsf{X}_T^\pi$ .

**Proof:** Simply note that  $\Omega_{0:T}^N \subset \bar{\Omega}_{0:T}^N$  and, as a consequence,  $\pi_{0:T}(\hat{x}_{0:T}^N) \leq \pi_{0:T}(\bar{x}_{0:T}^N) \leq \pi_{0:T}(\hat{x}_{0:T})$ .  $\square$

**Remark 4** We emphasize that the sequences  $\hat{x}_{0:T}^N$  and  $\bar{x}_{0:T}^N$  may not necessarily be convergent themselves as they may contain subsequences that converge to different elements of the critical set  $\mathsf{X}_T^\pi$  (we have not assumed uniqueness of the global minimizer). Moreover, if

$$\limsup_{\|x_{0:T}\| \rightarrow \infty} \pi_{0:T}(x_{0:T}) = \max_{x_{0:T} \in (\mathbb{R}^{d_x})^{T+1}} \pi_{0:T}(x_{0:T}) \quad (47)$$

then the sequence may contain subsequences that diverge to infinity or the entire sequence can diverge to infinity. If that is the case we need to restrict the search for a global minimizer to a (sufficiently large) compact set. However, in general,  $\lim_{\|x_{0:T}\| \rightarrow \infty} \pi_{0:T}(x_{0:T}) = 0$  and, therefore, ending up with a sequence divergent to infinity does not occur.

Equation (38) states that the fourth-order approximation error for a real bounded function of  $x_{0:T}$  converges “almost linearly” with the number of particles  $N$  that determines the accuracy of the discretization of the state-space  $\Omega_{0:T}^N$ . This enables us to find how large should the number of particles  $N$  be such that the (random) grids  $\Omega_{0:T}^N$  (respectively,  $\bar{\Omega}_{0:T}^N$ ) contain points at a distance from a true MAP estimate smaller than  $\frac{1}{k}$ , for  $k$  arbitrary but sufficiently large.

**Theorem 2** For sufficiently large  $k$ , the (random) grids  $\Omega_{0:T}^N$  and  $\bar{\Omega}_{0:T}^N$  contain points at a distance from a true MAP estimate smaller than  $\frac{1}{k}$  provided that  $N > ck^{\frac{4d_x(T+1)}{1-\epsilon}}$ , where  $c$  is a positive random variable that depends on the sequence  $y_{0:T}$  and the maximum value of the posterior density  $\pi_{0:T}(\hat{x}_{0:T})$ ,  $\hat{x}_{0:T} \in \mathsf{X}_T^\pi$ .

**Proof:** Eq. (38) implies that there exists a positive random variable  $c_T^\epsilon$  such that, for all  $N > 0$ , we have

$$\left| \frac{m(N, k)}{N} - A_k \right| \leq \frac{(c_T^\epsilon)^{\frac{1}{4}}}{N^{\frac{1-\epsilon}{4}}} \quad (48)$$

almost surely for any arbitrarily small  $\epsilon > 0$ . Recall that  $A_k = \int_{B_k(\hat{x}_{0:T})} \pi_{0:T}(x_{0:T}) dx_{0:T}$ , for some  $\hat{x}_{0:T} \in \mathsf{X}_T^\pi$ . When  $k$  is sufficiently large,  $\pi_{0:T}(x_{0:T})$  is very close to  $\pi_{0:T}(\hat{x}_{0:T})$  for any  $x_{0:T} \in B_k(\hat{x}_{0:T})$ . In

particular, we can assume that  $\frac{1}{2}\pi_{0:T}(\hat{x}_{0:T}) \leq \pi_{0:T}(x_{0:T}) \leq \pi_{0:T}(\hat{x}_{0:T})$  for any  $x_{0:T} \in B_k(\hat{x}_{0:T})$ . Therefore we can deduce that  $A_k \geq \frac{q_{T+1}}{2}\pi_{0:T}(\hat{x}_{0:T}) \left(\frac{1}{k}\right)^{d_x(T+1)}$ , where  $q_{T+1}$  is the volume of the unit ball in  $(\mathbb{R}^{d_x})^{T+1}$ , and from (48) we arrive at

$$\frac{q_{T+1}}{2}\pi_{0:T}(\hat{x}_{0:T}) \left(\frac{1}{k}\right)^{d_x(T+1)} - \frac{(c_T^\epsilon)^{\frac{1}{4}}}{N^{\frac{1-\epsilon}{4}}} \leq \frac{m(N, k)}{N}. \quad (49)$$

By inspection of (49), we realize that  $m(N, k)$  can be guaranteed to be strictly positive if we take  $N$  large enough for the inequality

$$\frac{q_{T+1}}{2}\pi_{0:T}(\hat{x}_{0:T}) \left(\frac{1}{k}\right)^{d_x(T+1)} - \frac{(c_T^\epsilon)^{\frac{1}{4}}}{N^{\frac{1-\epsilon}{4}}} > 0 \quad (50)$$

to hold true. Solving for  $N$ , we obtain  $N > ck^{\frac{4d_x(T+1)}{1-\epsilon}}$  for  $c = \left(\frac{2(c_T^\epsilon)^{\frac{1}{4}}}{q_{T+1}\pi_{0:T}(\hat{x}_{0:T})}\right)^{\frac{4d_x(T+1)}{1-\epsilon}}$ .  $\square$

**Remark 5** Under additional assumptions (for example if the state space is compact), one can deduce<sup>3</sup> a smaller lower bound for the size  $N$  of the sample required to obtain a point at a distance less than, say,  $\frac{1}{k}$ . The basis of this is the following exponential bound (see [22] for details and the required assumptions). One can show that there exist constants  $c_1 = c_1(T, f, \delta)$  and  $c_2 = c_2(T, f, \delta)$  such that

$$\Pr(|(f, \pi_{0:T}^N) - (f, \pi_{0:T})| \geq \delta) \leq c_1 e^{-c_2 N \delta^2} \quad (51)$$

for an arbitrarily small  $\delta > 0$ . Using a standard argument, one can deduce from Eq. (51) that there exist two positive random variables  $c_T^1$  and  $c_T^2$  such that, for all  $N > 0$ , we have

$$\left| \frac{m(N, k)}{N} - A_k \right| \leq c_T^1 e^{-c_T^2 N} \quad (52)$$

which implies that if  $N > c \log k$  for a suitably chosen positive random variable  $c$ , then  $m(N, k)$  is strictly positive and, hence, the (random) grids  $\Omega_{0:T}^N$  and  $\bar{\Omega}_{0:T}^N$  contain points at a distance smaller than  $\frac{1}{k}$ .

## 7 Numerical results

In this section we apply the proposed algorithms to the Examples 3 and 4. In particular, we first address the Neumaier 3 problem and show some numerical simulation results for Algorithms 1 and 2, as well as for the accelerated random search method of [2] for comparison. Then, we tackle the cross-talk cancellation problem of Example 4 using Algorithm 1.

---

<sup>3</sup>This approach was suggested to us by Pierre Del Moral.

## 7.1 Neumaier 3 problem

It is straightforward to apply Algorithm 1 to the Neumaier 3 problem described in Example 3. Specifically note that the likelihood is Gaussian,  $\pi(y_t|x_t) \propto \exp\{-\frac{1}{\sigma^2}(y_t - x_t)^2\}$ , while the transition density has an exponential form over a finite support,  $\pi(x_t|x_{t-1}) \propto \exp\{\frac{1}{\sigma^2}x_t x_{t-1}\}$ , for  $x_t \in [-T^2, +T^2]$ . Therefore, the evaluation of  $\pi(y_t|x_t)$  is straightforward and the generation of random samples from  $\pi(x_t|x_{t-1})$  is easily carried out using the inversion method [12]. It can be easily shown that if  $U$  is a uniform random variable in the interval  $[0, 1]$  and  $X_{t-1} = x_{t-1}$  is given, then

$$X_t = \frac{\sigma^2}{x_{t-1}} \log [\exp\{-n^2 x_{t-1}/\sigma^2\} + U (\exp\{n^2 x_{t-1}/\sigma^2\} - \exp\{-n^2 x_{t-1}/\sigma^2\})] \quad (53)$$

is a random variable with pdf  $\pi(x_t|x_{t-1})$ ,  $x_t \in [-T^2, +T^2]$ . The ability to evaluate  $\pi(y_t|x_t)$  and sample from  $\pi(x_t|x_{t-1})$  is sufficient to apply Algorithms 1 and 2.

In the first experiment, we check the influence of the scale factor  $\sigma^2$  on the solutions generated by the proposed optimization algorithms. Note that, even if the choice of  $\sigma^2 > 0$  is irrelevant from the perspective of the critical set  $X_T^\pi$  (i.e., the solutions of the optimization problem<sup>4</sup>  $\arg \min_{x_{1:T} \in [-T^2, +T^2]^T} \pi(x_{0:T}|y_{1:T})$  do *not* depend on  $\sigma^2$ ) the convergence rate of the numerical algorithms used to approximate the solutions in  $X_T^\pi$  may indeed be affected by this parameter.

Therefore, we have applied Algorithm 1 to the Neumaier 3 problem with a low dimension,  $T = 5$ , using  $N = 10^5$  particles and values of  $\sigma^2$  ranging from  $\sigma^2 = T^2$  to  $\sigma^2 = 500T^2$ . Figure 1 (left) shows the average cost of the solutions generated by Algorithm 1 for the various values of  $\sigma^2$ . Each point in the plot has been obtained by averaging the cost of the solution,  $C_{0:T}(\hat{x}_{0:T}^N)$ , over 100 independent simulation runs. The figure also depicts the true minimum cost for reference. It is observed that a small scale factor yields poorer solutions, while for  $\sigma^2 \geq 50T^2$  the generated solutions are close to optimum and any further increase of the scale parameter does not have an apparent effect on performance. In the sequel, we fix  $\sigma^2 = 150T^2$  for the rest of simulations of this example.

Figure 1 (right) shows the convergence of Algorithm 1 as the number of particles,  $N$ , is increased. For a fixed scale factor  $\sigma^2 = 150T^2$  and  $T = 5$  variables, we have carried out 100 independent simulation trials and averaged the cost of the approximate solution,  $C_{0:T}(\hat{x}_{1:T}^N)$ , for several values of  $N$ . The error reduction as  $N$  grows is apparent, but  $N = 10^7$  particles are needed to achieve a cost that is practically indistinguishable from the true minimum.

---

<sup>4</sup>Recall that the cost does not depend on the variable  $x_0$  for the Neumaier 3 problem.

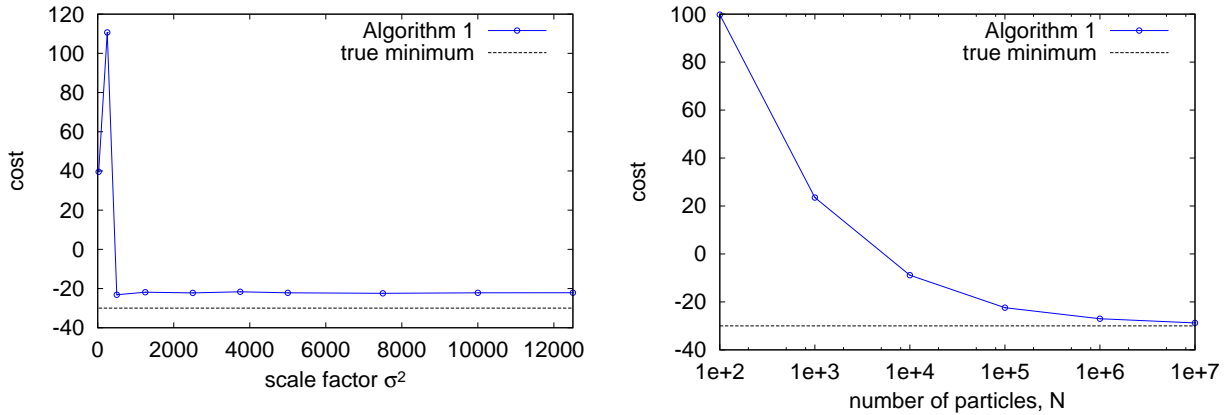


Figure 1: Performance of Algorithm 1 for the Neumaier 3 problem with dimension  $T = 5$ . *Left*: Average cost of the solution  $\hat{x}_{0:T}^N$ , with  $N = 10^5$  particles, for several values of the scale parameter  $\sigma^2$ . *Right*: For fixed  $\sigma^2 = 150T^2$ , average cost of  $\hat{x}_{0:T}^N$  for several values of  $N$ .

In Figure 2 (left) we show that the convergence of Algorithm 2 can be much faster in terms of the number of required particles. For the same  $T = 5$ -dimensional problem, we show the average of the costs  $C_{0:T}(\hat{x}_{1:T}^N)$  (for Algorithm 1) and  $C_{0:T}(\bar{x}_{1:T}^N)$  (for Algorithm 2) over 100 independent simulation runs. This time, the number of particles is increased from  $N = 20$  up to  $N = 500$ . Algorithm 1 attains very poor solutions with this small number of samples, while Algorithm 2 practically achieves the true minimum cost for  $N \geq 50$ .

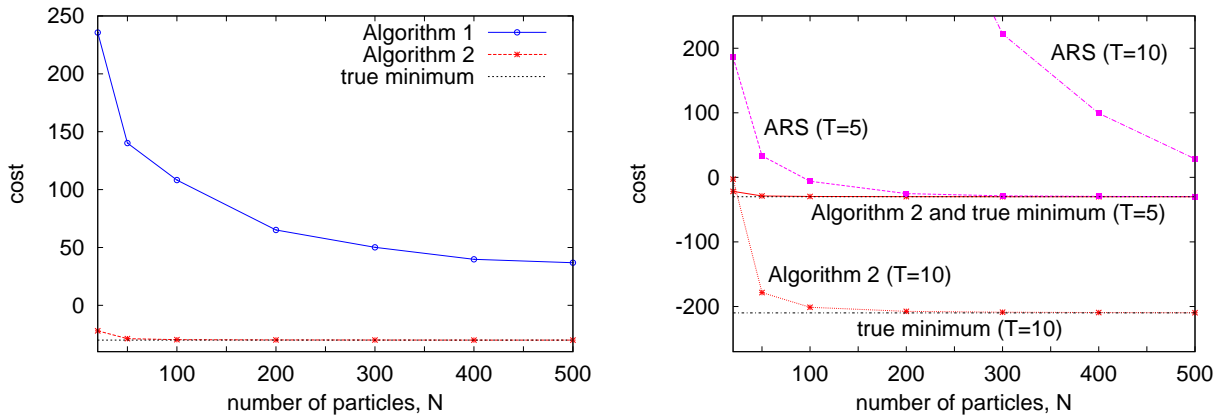


Figure 2: Performance of Algorithm 2 for the Neumaier 3 problem. *Left*: Comparison of Algorithms 1 and 3 in terms of the cost as a function of the number of particles,  $N$ . *Right*: Comparison of Algorithm 2 with the ARS method for  $T = 5$  and  $T = 10$ .

Figure 2 (right) shows a comparison of Algorithm 2 with the accelerated random search (ARS) method [2] when the dimension of the Neumaier 3 problem is  $T = 5$  and  $T = 10$ . The ARS procedure is an iterative

algorithm that seeks solutions within a “contracting-and-expanding” ball centered at the best candidate solution found so far. For the simulations, we set the number of ARS iterations to be equal to the number of particles in Algorithm 2. In this way, both algorithms generate  $N$  complete paths in the space  $[-T^2, +T^2]^T$ . The ARS procedure was run with maximum and minimum radii of  $T^2$  and  $10^{-6}$ , respectively, and a contraction factor of 3 (see [2] for details). Both for  $T = 5$  and  $T = 10$ , Algorithm 2 attains a clearly superior performance. It is worth mentioning that the ARS technique requires sampling uniformly within balls of varying radii. This is done by rejection sampling, but as the dimension of the space increases, the procedure becomes less and less efficient. For  $T = 10$ , an average of  $2 \times 10^5$  actual samples in  $[-T^2, +T^2]^T$  are needed in order to obtain  $N = 500$  effective paths (as a large proportion of them are rejected).

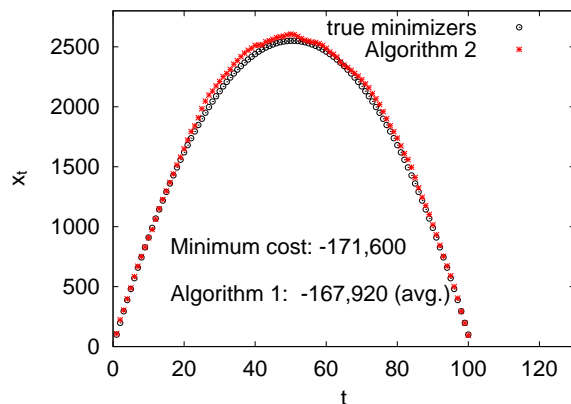


Figure 3: Performance of Algorithm 2 for the Neumaier 3 problem with high dimension,  $T = 100$ , and  $N = 3000$  particles. The figure shows the true minimizers,  $\hat{x}_{1:T}$  (circles), and the approximate values,  $\bar{x}_{1:T}^N$  (stars).

The proposed Algorithm 2 can be used with much higher dimensions. Figure 3 shows the approximate solution  $\bar{x}_{1:T}^N$ , together with the true global minimizer  $\hat{x}_{1:T}$ , for a problem of dimension  $T = 100$  using  $N = 3000$  particles. It can be seen that a close approximation is achieved. The cost of the solution is  $-167,920$ , while the true minimum cost is  $-171,600$ .

## 7.2 Cross-talk cancellation

We turn attention now to the cross-talk cancellation problem of Example 4. Recall that the goal of this problem is to compute an (approximately) inverse filter for a given acoustic response  $h_a(n)$ ,  $n = 0, \dots, 2M - 2$ . In this section, we illustrate the performance of Algorithm 1 in this task. Algorithm

2 can also be applied to this example, although it requires a straightforward adaptation of the procedure presented in Section 6.2.2 motivated by the fact that  $\pi(y_t|x_{0:t}, y_{1:t-1}) = \pi(y_t|x_{t-2M+2:t}) \neq \pi(y_t|x_t)$  in the corresponding state-space model.

For the causal impulse response  $h_a(n)$  to be actually invertible, we must ensure that all the roots of the  $z$ -transform polynomial  $H_a(z^{-1}) = \sum_{n=0}^{2M-2} h_a(n)z^{-n}$  lie inside the unit circle. Formally, let  $z_l^{-1}$ ,  $l = 1, \dots, 2M - 2$ , be the solutions of the equation  $H_a(z^{-1}) = 0$ . If  $|z_l^{-1}| < 1$  for  $l = 1, \dots, 2M - 2$ , then the poles of the inverse system  $H_a^{-1}(z^{-1}) = 1/H_a(z^{-1})$  also lie inside the unit circle and, as a consequence,  $H_a^{-1}(z^{-1})$  is stable. A stable system with stable inverse is termed a “minimum phase” system [24]. A minimum phase acoustic response can be approximately inverted by solving problem (30), if we take sufficiently large  $K$  and allow the cancellation filter coefficients to be, also, large enough, while non-minimum phase systems cannot be inverted because the inverse  $z$ -transform of  $H_a^{-1}(z^{-1})$  diverges.

In the subsequent experiments, the dimension of the problem is given in terms of the time horizon  $T$ . Let us recall that this horizon depends on the length of the acoustic response,  $2M - 1$ , and the length of the cancellation filter,  $K$ . For our simulations, we select the latter as  $K = 7(2M - 1)$ , which yields  $T = K + 2M - 2 = 16M - 9$ . We have set  $M = 7$ , hence  $K = 91$  and  $T = 103$ . Notice, nevertheless, that the number of coefficients to be selected in the sequence  $\hat{x}_{0:T}^N$  is  $K$ , since  $x_{t \leq 0} = x_{t > K} = 0$ . The desired output of the cross-talk cancellation filter is  $d(t) = y_t = \delta_\kappa(t - 1)$  [according to (29)], where  $\delta_\kappa$  denotes the Kronecker delta function. Given a MAP estimate  $\hat{x}_{0:T}^N$ , the actual output of the filter is denoted as  $\hat{y}_t^N = \sum_{k=0}^{2M-2} h_a(k)\hat{x}_{t+1-k}^N$ . Note that  $C_{0:T}(\hat{x}_{0:T}^N) = \max_{t \in \{1, \dots, T\}} |y_t - \hat{y}_t^N|$ .

Figure 4 shows the results obtained by applying Algorithm 1 to three sample acoustic responses. The figure is arranged in columns and each column corresponds to a different realization of  $h_a(n)$ . For each acoustic response, we show the locations of the zeros of the system  $H_a(z^{-1})$  (in the upper plot), the obtained MAP estimate  $\hat{x}_{0:T}^N$  (center plot) and the resulting output  $\hat{y}_{1:T}^N$ , together with the desired signal  $y_t = \delta_\kappa(t - 1)$  (lower plot). In the three simulations, Algorithm 1 was implemented with  $N = 10^4$  particles and state-transition pdf  $\pi(x_t|x_{0:t-1}) = \pi(x_t) = U(x_t; -5, +5)$ .

The first column of Figure 4 depicts the results obtained for an “easy” system  $H_a(z^{-1})$  with all its roots inside the unit circle and close to 0. We see that the impulse response of the cancellation filter is relatively short, i.e., the sequence  $\hat{x}_t^N$  takes residual values for (approximately)  $t > 15$  (hence, we could have inverted  $h_a(n)$  with a much shorter filter). Correspondingly, the difference between the actual output  $\hat{y}_{1:T}^N$  and the desired signal  $y_{1:T}$  is very small.



The second column of Figure 4 shows the results for an acoustic response with larger roots, but still bounded away from the unit circle. We observe that the effective length of the cancellation filter is larger, i.e., the coefficients  $\hat{x}_t^N$  take non-negligible values almost for the whole length,  $K$ , of the filter. Still, the output  $\hat{y}_{1:T}^N$  matches the desired  $y_{1:T}$  closely.

Finally, the third column of Figure 4 illustrates the result of attempting the inversion of an acoustic response with its zeros very close to the unit circle. With the selected filter length,  $K = 91$ , and pdf  $\pi(x_t) = U(x_t; -5, +5)$ , accurate inversion cannot be carried out. As shown in the center plot, the sequence  $\hat{x}_t^N$  does not converge toward 0, as it does in the previous examples, and the output  $\hat{y}_t^N$  departs significantly from  $y_t$ . The accurate inversion of this response would require a longer filter (i.e., greater  $K$ ) and a pdf  $\pi(x_t) = U(x_t; -a, +a)$  with sufficiently large  $a$ .

## 8 Summary

We have shown how a class of optimization problems consisting in the minimization of a cost function with a recursive structure can be transformed into equivalent estimation problems by constructing a suitable state-space dynamic model. The unknowns of the cost function determine the state-space of the dynamic model and both of them are defined to be “matched” when the set of minimizers of the cost coincide with the set of MAP estimates for the dynamic model.

Once recast as an estimation problem, we can take advantage of the SMC methodology for the approximation of probability measures in state-space models in order to numerically find solutions of the original optimization problem. Specifically, the SMC method yields a discretization of the state-space (equivalently, of the space of the unknowns) that is dense where the posterior probability mass is high (equivalently, where the cost is low) and sparse elsewhere. Then, it is possible to approximate the MAP estimates of the model (equivalently, the minimizers of the cost) by searching over this discretized space. We have described two algorithms for this purpose. The first one simply selects the sample path with the highest posterior probability density (equivalently, with the lowest cost) out of those generated by the SMC method. The second one constructs a refined random grid approximation of the state-space by allowing the combination of different sample paths and then searches the best point of this refined grid by means of the Viterbi algorithm. We have shown this scheme to work properly even with the simplest SMC method (the standard SIR algorithm) and have illustrated its performance with two examples borrowed from the fields of global optimization and signal processing. From one of these examples it is clearly seen

that the use of the refined grid can be very advantageous in terms of the accuracy of the solutions.

The approximate MAP estimates (equivalently, cost minimizers) generated by the proposed algorithms depend on the number of samples,  $N$ , allowed to the SMC method in order to discretize the state space. We have proved that, as  $N$  grows, the posterior probability density of the solutions output by Algorithms 1 and 2 converge almost surely to the true maximum *a posteriori* density (equivalently, the cost of the solutions converges to the true minimum). Moreover, we have derived bounds, for the number of samples,  $N$ , needed to attain a certain accuracy that hold almost surely. It is also worth noting that, as an instrument to analyze the proposed optimization algorithms we have derived  $L_p$  bounds for the errors in the integrals of bounded functions in the path space as approximated using the SIR algorithm. This result was not explicitly available in the literature on SMC methods so far.

## Acknowledgements

J. M. acknowledges the support of the Ministry of Science and Technology of Spain (program Consolider-Ingenio 2010 CSD2008-00010 COMONSENS and project DEIPRO TEC2009-14504-C02-01).

Part of this work was done during D. C.'s visit to the Department of Signal Theory & Communications, Universidad Carlos III (Spain), in April 2008. The hospitality of the Department is gratefully acknowledged.

The work of P. M. D. has been supported by the National Science Foundation under Award CCF-0515246 and the Office of Naval Research under Award N00014-06-1-0012. Part of this work was carried out while P. M. D. held a Chair of Excellence of Universidad Carlos III de Madrid-Banco de Santander.

## A Proof of Lemma 1

We proceed by induction in  $T$ . For  $T = 0$ , the random measure  $\pi_{0:0}^N(dx)$  is constructed from an i.i.d. sample of size  $N$  from the distribution with pdf  $\pi_{0:0}$ . Hence, it is straightforward to check that

$$\|(f, \pi_{0:0}^N) - (f, \pi_{0:0})\|_p \leq \frac{c_0^p \|f\|_\infty}{\sqrt{N}}, \quad (54)$$

where  $c_0^p$  is a constant independent of  $N$ .

Now we assume that

$$\|(f, \pi_{0:T}^N) - (f, \pi_{0:T})\|_p \leq \frac{c_T^p \|f\|_\infty}{\sqrt{N}}, \quad (55)$$

for an integer  $T > 0$  and aim at proving the corresponding inequality for  $T + 1$ .

The recursive step of the SIR algorithm, as presented in Section 6.1, consists of three sub-steps. Let  $p_{0:T+1}^N$  be the empirical measure obtained after the first sub-step, i.e.,  $p_{0:T+1}^N(dx) = \frac{1}{N} \sum_{n=1}^N \delta_{\bar{x}_{0:T+1}^{(n)}}(dx)$ , where  $\delta_{\bar{x}_{0:T+1}^{(n)}}$  denotes the unit delta measure centered at  $\bar{x}_{0:T+1}^{(n)}$ . Also let  $\mathcal{G}^{T,N}$  denote the  $\sigma$ -algebra generated by the random variates  $X_{0:T}^{(n)}$ ,  $n = 1, \dots, N$ . Then, for  $f : (\mathbb{R}^{d_x})^{T+2} \rightarrow \mathbb{R}$ , we have

$$E[(f, p_{0:T+1}^N) | \mathcal{G}^{T,N}] = (\bar{f}, \pi_{0:T}^N), \quad (56)$$

where  $\bar{f}$  is obtained from  $f$  by integrating with respect to the measure  $\pi(x_{T+1} | x_{0:T}) dx_{T+1}$ , i.e.,

$$\bar{f}(x_0, x_1, \dots, x_T) \triangleq \int_{\mathbb{R}^{d_x}} f(x_0, x_1, \dots, x_{T+1}) \pi(x_{T+1} | x_{0:T}) dx_{T+1}. \quad (57)$$

Obviously,  $\bar{f}$  is bounded (since  $\|\bar{f}\|_\infty \leq \|f\|_\infty$ ) and, from the induction hypothesis (55), we deduce that

$$\|(\bar{f}, \pi_{0:T}^N) - (\bar{f}, \pi_{0:T})\|_p \leq \frac{c_T^p \|f\|_\infty}{\sqrt{N}}. \quad (58)$$

Moreover, since

$$E\left[\left((f, p_{0:T+1}^N) - E[(f, p_{0:T+1}^N) | \mathcal{G}^{T,N}]\right)^p | \mathcal{G}^{T,N}\right]^{\frac{1}{p}} \leq \frac{\tilde{c}_{T+1}^p \|f\|_\infty}{\sqrt{N}}, \quad (59)$$

where  $\tilde{c}_{T+1}^p$  is a positive random variable independent of  $N$ , it is straightforward to combine (56), (58) and (59) using the triangle inequality to arrive at

$$\|(f, p_{0:T+1}^N) - (\bar{f}, \pi_{0:T})\|_p \leq \frac{\tilde{c}_{T+1}^p \|f\|_\infty}{\sqrt{N}}, \quad (60)$$

where  $\tilde{c}_{T+1}^p = E[\tilde{c}_{T+1}^p] + c_T^p$ .

Consider next the measure  $\bar{\pi}_{0:T+1}^N$  that is obtained after sub-step ii. of the algorithm. This measure can be defined by

$$(f, \bar{\pi}_{0:T+1}^N) = \frac{(f g_{T+1}, p_{0:T+1}^N)}{(g_{T+1}, p_{0:T+1}^N)} \quad (61)$$

(recall that  $g_{T+1}(x_{0:T+1}) = \pi(y_{T+1} | x_{0:T+1}, y_{1:T})$  is the bounded likelihood function). Also let  $p_{0:T+1}(x_{0:T+1}) dx_{0:T+1}$  be the predictive measure that satisfies  $(f, p_{0:T+1}) = (\bar{f}, \pi_{0:T})$  and rewrite (60) as

$$\|(f, p_{0:T+1}^N) - (f, p_{0:T+1})\|_p \leq \frac{\tilde{c}_{T+1}^p \|f\|_\infty}{\sqrt{N}}. \quad (62)$$

Since, from the Bayes' rule,

$$(f, \pi_{0:T+1}) = \frac{(f g_{T+1}, p_{0:T+1})}{(g_{T+1}, p_{0:T+1})}, \quad (63)$$

we can take (61) and (63) together in order to obtain

$$(f, \bar{\pi}_{0:T+1}^N) - (f, \pi_{0:T+1}) = \frac{(f g_{T+1}, p_{0:T+1}^N)}{(g_{T+1}, p_{0:T+1}^N)} - \frac{(f g_{T+1}, p_{0:T+1})}{(g_{T+1}, p_{0:T+1})}. \quad (64)$$

By adding and subtracting the term  $(fg_{T+1}, p_{0:T+1}^N)/(g_{T+1}, p_{0:T+1})$  in the equation above, we easily arrive at

$$\begin{aligned} (f, \bar{\pi}_{0:T+1}^N) - (f, \pi_{0:T+1}) &= \frac{(fg_{T+1}, p_{0:T+1}^N)}{(g_{T+1}, p_{0:T+1}^N)(g_{T+1}, p_{0:T+1})} [(g_{T+1}, p_{0:T+1}) - (g_{T+1}, p_{0:T+1}^N)] \\ &\quad + \frac{1}{(g_{T+1}, p_{0:T+1})} [(fg_{T+1}, p_{0:T+1}^N) - (fg_{T+1}, p_{0:T+1})] \end{aligned} \quad (65)$$

and, since  $|(fg_{T+1}, p_{0:T+1}^N)| \leq \|f\|_\infty (g_{T+1}, p_{0:T+1}^N)$ , it readily follows that

$$\begin{aligned} |(f, \bar{\pi}_{0:T+1}^N) - (f, \pi_{0:T+1})| &\leq \frac{\|f\|_\infty}{(g_{T+1}, p_{0:T+1})} |(g_{T+1}, p_{0:T+1}) - (g_{T+1}, p_{0:T+1}^N)| \\ &\quad + \frac{1}{(g_{T+1}, p_{0:T+1})} |(fg_{T+1}, p_{0:T+1}^N) - (fg_{T+1}, p_{0:T+1})|. \end{aligned} \quad (66)$$

The latter inequality, together with (62) and the assumed boundedness of the likelihood  $g_{T+1}$ , yields

$$\|(f, \bar{\pi}_{0:T+1}^N) - (f, \pi_{0:T+1})\|_p \leq \frac{\mathcal{C}_{T+1}^p \|f\|_\infty}{\sqrt{N}}, \quad (67)$$

where  $\mathcal{C}_{T+1}^p = 2\|g_{T+1}\|_\infty \bar{\mathcal{C}}_{T+1}^p / (g_{T+1}, p_{0:T+1})$  is a constant independent of  $N$ .

In order to analyze the last substep (the resampling), we introduce the  $\sigma$ -algebra generated by the random variates  $\bar{X}_{0:T+1}^{(n)}$ ,  $n = 1, \dots, N$ , and denote it as  $\bar{\mathcal{G}}^{T+1, N}$ . It is straightforward to obtain that  $E[(f, \pi_{0:T+1}^N) | \bar{\mathcal{G}}^{T+1, N}] = (f, \bar{\pi}_{0:T+1}^N)$ , hence the conditional expectation of the error becomes

$$E \left[ ((f, \pi_{0:T+1}^N) - (f, \bar{\pi}_{0:T+1}^N))^p | \bar{\mathcal{G}}^{T+1, N} \right]^{\frac{1}{p}} \leq \frac{\mathcal{C}_{T+1}^p \|f\|_\infty}{\sqrt{N}}, \quad (68)$$

where  $\mathcal{C}_{T+1}^p$  is a positive random variable independent of  $N$ . As a consequence, taking the expectation on  $X_{0:T+1}^{(n)}$ ,  $n = 1, \dots, N$ , yields

$$\|(f, \pi_{0:T+1}^N) - (f, \bar{\pi}_{0:T+1}^N)\|_p \leq \frac{\bar{\mathcal{C}}_{T+1}^p \|f\|_\infty}{\sqrt{N}}, \quad (69)$$

where  $\bar{\mathcal{C}}_{T+1}^p$  is the expected value of  $\mathcal{C}_{T+1}^p$ . Combining (67) and (69) by way of the triangle inequality yields

$$\|(f, \pi_{0:T+1}^N) - (f, \pi_{0:T+1})\|_p \leq \|(f, \pi_{0:T+1}^N) - (f, \bar{\pi}_{0:T+1}^N)\|_p + \|(f, \bar{\pi}_{0:T+1}^N) - (f, \pi_{0:T+1})\|_p \leq \frac{\bar{\mathcal{C}}_{T+1}^p \|f\|_\infty}{\sqrt{N}}, \quad (70)$$

where  $\bar{\mathcal{C}}_{T+1}^p = \bar{\mathcal{C}}_{T+1}^p + \mathcal{C}_{T+1}^p$  is a constant independent of  $N$ .  $\square$

## References

- [1] M. M. Ali, C. Khompatraporn, and Z. B. Zabinsky. A numerical evaluation of several stochastic algorithms on selected continuous global optimization problems. *Journal of Global Optimization*, 31:635–672, 2005.

- [2] M. J. Appel, R. Labarre, and D. Radulovic. On accelerated random search. *SIAM Journal on Optimization*, 14(3):708–730, 2003.
- [3] R. D. Baker. How to correctly calculate discounted healthcare costs and benefits. *The Journal of the Operational Research Society*, 51(7):863–868, July 2000.
- [4] D. Crisan. Particle filters - a theoretical perspective. In A. Doucet, N. de Freitas, and N. Gordon, editors, *Sequential Monte Carlo Methods in Practice*, chapter 2, pages 17–42. Springer, 2001.
- [5] D. Crisan and A. Doucet. Convergence of sequential Monte Carlo methods. Technical Report Cambridge University (CUED/FINFENG/TR381), 2000.
- [6] R. Douc, O. Cappé, and E. Moulines. Comparison of resampling schemes for particle filtering. In *Proceedings of the 4<sup>th</sup> International Symposium on Image and Signal Processing and Analysis*, pages 64–69, September 2005.
- [7] A. Doucet, N. de Freitas, and N. Gordon. An introduction to sequential Monte Carlo methods. In A. Doucet, N. de Freitas, and N. Gordon, editors, *Sequential Monte Carlo Methods in Practice*, chapter 1, pages 4–14. Springer, 2001.
- [8] A. Doucet, N. de Freitas, and N. Gordon, editors. *Sequential Monte Carlo Methods in Practice*. Springer, New York (USA), 2001.
- [9] A. Doucet, S. Godsill, and C. Andrieu. On sequential Monte Carlo Sampling methods for Bayesian filtering. *Statistics and Computing*, 10(3):197–208, 2000.
- [10] D. Du and P. M. Pardalos, editors. *Minimax and applications*. Kluwer Academic Publishers, 1995.
- [11] G. D. Forney. The Viterbi algorithm. *Proceedings of the IEEE*, 61(3):268–278, 1973.
- [12] J.E. Gentle. *Random Number Generation and Monte Carlo Methods*. Springer, 2003.
- [13] W.R. Gilks, S. Richardson, and D. Spiegelhalter. *Markov Chain Monte Carlo in Practice*. Chapman and Hall/CRC, 1996.
- [14] S. Godsill, A. Doucet, and M. West. Maximum a posteriori sequence estimation using Monte Carlo particle filters. *Annals of the Institute of Statistical Mathematics*, 53(1):82–96, March 2001.
- [15] S. Godsill, A. Doucet, and M. West. Monte carlo smoothing for nonlinear time series. *Journal of the American Statistical Association*, 99(465):156–168, March 2004.
- [16] N. Gordon, D. Salmond, and A. F. M. Smith. Novel approach to nonlinear and non-Gaussian Bayesian state estimation. *IEE Proceedings-F*, 140(2):107–113, 1993.

- [17] A.-R. Hedar and M. Fukushima. Derivative-free filter simulated annealing method for constrained continuous global optimization. *Journal of Global Optimization*, 35:521–549, 2006.
- [18] R. Horst and N. V. Thoai. DC programming: overview. *Journal of Optimization Theory and Applications*, 103(1):1–43, 1999.
- [19] J. S. Liu and R. Chen. Sequential Monte Carlo methods for dynamic systems. *Journal of the American Statistical Association*, 93(443):1032–1044, September 1998.
- [20] P. Del Moral. *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*. Springer, 2004.
- [21] P. Del Moral and M. Doisy. On applications of Maslov optimization theory. *Mathematical Notes*, 69(2):232–244, 2001.
- [22] P. Del Moral and L. Miclo. Branching and interacting particle systems. Approximations of Feynman-Kac formulae with applications to non-linear filtering. *Lecture Notes in Mathematics*, pages 1–145, 2000.
- [23] K. Najim, E. Ikonen, and P. Del Moral. Open-loop regulation and tracking control based on a genealogical decision tree. *Neural Computing and Applications*, 15:339–349, 2006.
- [24] A.V. Oppenheim and R.W. Schaffer. *Discrete-Time Signal Processing*. Prentice-Hall, Englewood Cliffs, New Jersey, 1989.
- [25] A. R. Pankov, E. N. Platonov, and K. V. Semenikhin. Minimax optimization of investment portfolio by quantile criterion. *Autom. Remote Control*, 64(7):1122–1137, 2003.
- [26] P. M. Pardalos, H. E. Romeijn, and H. Tuy. Recent developments and trends in global optimization. *Journal of Computational and Applied Mathematics*, 124(1-2):209–228, December 2000.
- [27] D. J. Ram, T. H. Sreenivas, and K. G. Subramaniam. Parallel simulated annealing algorithms. *Journal of Parallel and Distributed Computing*, 37(2):207–212, 1996.
- [28] H. I. K. Rao, V. J. Mathews, and Y.-C. Park. A minimax approach for the joint design of acoustic cross-talk cancellation filters. *IEEE Transactions on Audio, Speech and Language Processing*, 15(8):2287–2298, November 2007.
- [29] C. P. Robert and G. Casella. *Monte Carlo Statistical Methods*. Springer, 2004.
- [30] S. M. Ross. *Applied Probability Models with Optimization Applications*. Dover Publications, New York, 1992.

- [31] A. H. Sayed, A. Tarighat, and N. Khajehnouri. Network based wireless location. *IEEE Signal Processing Magazine*, 22(4):24–40, July 2005.
- [32] R. Storn and K. Price. Differential evolution a simple and efficient heuristic for global optimization over continuous spaces. *Journal of Global Optimization*, 11:341–359, 1997.
- [33] H. Tuy. Monotonic optimization: Problems and solution approaches. *SIAM Journal on Optimization*, 11(2):464–494, 2000.
- [34] W. T. Ziemba and R. G. Vickson, editors. *Stochastic Optimization Models in Finance*. World Scientific, Singapore, 2006.
- [35] W. Zu and Q. Fu. A sequential convexification method (SCM) for continuous global optimization. *Journal of Global Optimization*, 26:167–182, 2003.

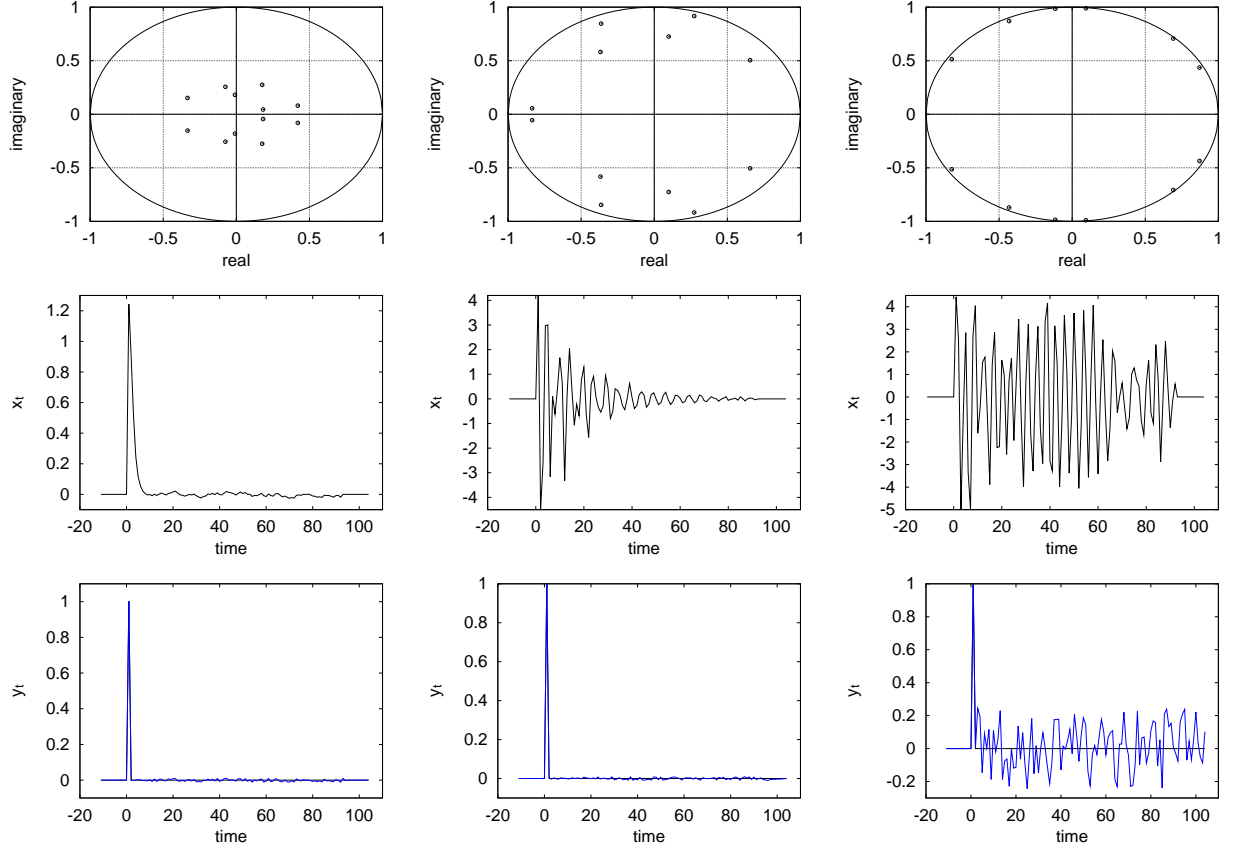


Figure 4: Performance of Algorithm 1 for the cross-talk cancellation problem with  $T = 103$ . **Left column:** Results for a combined response  $h_a(n)$  with all its roots close to 0. The upper plot shows the location of the zeros of  $H_a(z^{-1})$ . The center plot depicts the cancellation filter computed by Algorithm 1,  $\hat{x}_{0:T}^N$ . The lower plot shows the filtered response  $\hat{y}_t^N$  together with the desired response  $y_t = \delta_\kappa(t-1)$ . **Middle column:** Performance of Algorithm 1 as the zeros of  $H_a(z^{-1})$  spread within the unit circle. The upper plot shows the location of the zeros. The center plot shows the cancellation filter coefficients,  $\hat{x}_{0:T}^N$ . The lower plot depicts the filtered response  $\hat{y}_t^N$  together with the desired response  $y_t$ . **Right column:** As the zeros of the acoustic response  $h_a(n)$  move onto the unit circle, the computation of an inverse filter becomes tougher. The upper plot shows the location of the zeros. The center plot shows the cancellation filter coefficients  $\hat{x}_{0:T}^N$ . The lower plot shows the output  $\hat{y}_{1:T}^N$ , which departs from  $y_t$ .