

# Simulation-Based Optimal Sensor Scheduling with Application to Observer Trajectory Planning

Sumeetpal Singh<sup>a</sup>, Nikolas Kantas<sup>a</sup>, Ba-Ngu Vo<sup>b</sup>, Arnaud Doucet<sup>a</sup> and Robin J. Evans<sup>b</sup>

<sup>a</sup>Signal Processing Group, Dept. of Eng., Univ. of Cambridge, UK

<sup>b</sup>Dept. of Elec. and Electronic Eng., Univ. of Melbourne, Australia

**Abstract**—Sensor scheduling has been a topic of interest to the target tracking community for some years now. Recently, research into it has enjoyed fresh impetus with the current importance and popularity of applications in Sensor Networks and Robotics. The sensor scheduling problem can be formulated as a controlled Hidden Markov Model. In this paper, we address precisely this problem and consider the case in which the state, observation and action spaces are continuous valued vectors. This general case is important as it is the natural framework for many applications. We present a novel simulation-based method that uses a stochastic gradient algorithm to find optimal actions.<sup>1</sup>

## I. INTRODUCTION

Consider the following continuous state Hidden Markov Model (HMM),  $X_{k+1} = f(X_k, A_{k+1}, W_k) \in \mathbb{R}^{n_x}$ ,  $Y_k = g(X_k, A_k, V_k) \in \mathbb{R}^{n_y}$ , where  $X_k$  is the hidden system state,  $Y_k$  the observation of the state, and  $W_k$  and  $V_k$  are i.i.d. noise terms. Unlike the classical HMM model, the evolution of the state and observation processes depends on an input parameter  $A_k \in \mathbb{R}^{n_a}$ , which is the control or action. In HMM models, one is primarily concerned with the problem of estimating the hidden state, which is achieved by propagating the posterior distribution (or filtering density)  $\pi_k(x)dx = \mathbf{P}(X_k \in dx | A_{1:k}, Y_{1:k})$ . By a judicious choice of control sequence  $\{A_k\}$ , the evolution of the state and observation processes can be ‘steered’ in order to yield filtering densities that give more accurate estimates of the state process. This problem is also known in the literature as the sensor scheduling problem.

Sensor scheduling has been a topic of interest to the target tracking community for the some years now [3], [7], [10], [11], [5]. The classical setting is the problem of tracking a maneuvering target over  $N$  epochs. Here  $X_k$  denotes the state of the target at epoch  $k$ ,  $Y_k$  the observation provided by the sensor, and  $A_k$  some parameter of the sensor that may be adjusted to improve the “quality” of the observation. A measure of tracking performance is the variance of the tracking error over the  $N$  epochs:

$$\mathbf{E} \left\{ (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \right\}, \quad k = 1, \dots, N, \quad (1)$$

where  $\psi : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$  is a suitable *test* function that emphasises the components of interest of the state vector

we wish to track. The aim is to minimise the tracking error variance with respect to the choice of actions  $\{A_1, \dots, A_N\}$ .

When the dynamics of the state and the observation processes are both linear and Gaussian then, the optimal solution to the sensor scheduling problem (1) (when  $\psi$  gives a quadratic cost) can be computed off-line; this is not surprising given that the Kalman filter covariance is also independent of the actual realisation of observations. In the general setting studied in this paper, the dynamics can be both non-linear and non-Gaussian, which means that the filtering density  $\pi_k$ , and integration with respect to it, cannot be evaluated in closed-form. Hence, the variance performance criterion itself does not admit a closed-form expression. To further complicate matters, the actions sought are continuous valued, i.e., vectors in  $\mathbb{R}^{n_a}$ .

To address the complications to do with the non-linear and non-Gaussian dynamics, one could linearise the state and observation model [7]. The majority of works [5], [8], [11] (and references therein), while aim at minimising the tracking error variance, do so approximately by minimising a lower bound to the variance criterion. The bound in question is the Posteriori Cramer-Rao Lower Bound (PCRLB), which is the inverse of the Fisher Information Matrix (FIM). This approach hinges on the ability to propagate recursively the FIM in closed form by a Ricatti-type equation for the non-linear and non-Gaussian filtering problem. Unfortunately, the recursion for the FIM involves evaluating the expectation of certain derivatives of the transition probability density of the state dynamics, as well as the expectation of certain derivatives of the observation likelihood (see (2) and (3) below). As these quantities cannot be evaluated in general except for the linear and Gaussian case, this assumption is either invoked or the authors resort to simulation-based approximations. In addition, the PCRLB bound is not always tight.

The aim of this paper is to solve the sensor scheduling problem with continuous action space directly. We make no assumptions of linearity or Gaussianity for analytic convenience, nor do we discretise the state, observation, or action space. We avoid these restrictive modelling assumptions on the continuous state HMM by recourse to methods based on computer simulation (simulation for short). Under suitable regularity assumptions, one can guarantee convergence to a local optimum of the performance criterion, while it is difficult to make similar assertions about the quality of the

<sup>1</sup>Acknowledgement: S.Singh and A. Doucet were funded by EPSRC, N. Kantas by DIF-DTC, and B. Vo by an ARC large grant.

solutions obtained by other approximate methods proposed in the literature for sensor scheduling.

*Notation:* The notation that is used in the paper is now outlined. The norm of a scalar, vector or matrix is denoted by  $|\cdot|$ . For a vector  $b$ ,  $|b|$  denotes the vector 2-norm  $\sqrt{\sum_i |b(i)|^2}$ . For a matrix  $A$ ,  $|A|$  denotes the matrix 2-norm,  $\max_{b:|b|\neq 0} \frac{|Ab|}{|b|}$ . For convenience, we also denote a vector  $b \in \mathbb{R}^n$  by  $b = [b(i)]_{i=1,\dots,n}$ , or the  $i$ -th component of a vector by  $[b]_i$ . For scalars  $a_{j,i}$ ,  $j = 1, \dots, m$ ,  $i = 1, \dots, n$ , let  $\left[ [a_{j,i}]_{j=1,\dots,m} \right]_{i=1,\dots,n}$  denote the stacked vector  $[a_{1,1}, \dots, a_{m,1}, \dots, a_{1,n}, \dots, a_{m,n}]^T$ . For a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with arguments  $z \in \mathbb{R}^n$ , we denote  $(\partial f / \partial z(i))(z)$  by  $\nabla_{z(i)} f(z)$  and  $\nabla f(z) = [\nabla_{z(1)} f(z), \dots, \nabla_{z(n)} f(z)]^T$ . For the vector valued function  $F = [F_1, \dots, F_n]^T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , let  $\nabla F$  denote the matrix  $[\nabla F_1, \dots, \nabla F_n]$ . For a real-valued integrable functions  $f$  and  $g$ , let  $\langle f, g \rangle$  denote  $\int f(x)g(x)dx$ .

## II. PROBLEM FORMULATION

At time  $k$ , let  $X_k$  and  $Y_k$  be random vectors that model the  $n_x$ -dimensional state and its  $n_y$ -dimensional observation respectively. Suppose that an action  $A_k \in \mathbb{R}^{n_a}$  is applied at time  $k$ . The state  $\{X_k\}_{k \geq 0}$  is an unobserved Markov process with initial distribution and transition law given by

$$X_0 \sim \pi_0, \quad X_{k+1} \sim p(\cdot | X_k, A_{k+1}). \quad (2)$$

(The symbol “ $\sim$ ” means distributed according to.) The observation process  $\{Y_k\}_{k \geq 1}$  is generated according to the state and action dependent probability density

$$Y_k \sim q(\cdot | X_k, A_k). \quad (3)$$

Given the sequence of actions  $a_{1:k} := \{a_1, \dots, a_k\}$  and measurements  $y_{1:k} := \{y_1, \dots, y_k\}$ , the *filtering density* at time  $k$  is denoted by  $\pi_k$ , (or  $\pi_k^{(y_{1:k}, a_{1:k})}$  to emphasise the dependence on  $y_{1:k}$  and  $a_{1:k}$ ), and satisfies the *Bayes* recursion

$$\pi_k(x) = \frac{q(y_k | x, a_k) \int p(x | x', a_k) \pi_{k-1}(x') dx'}{\int \int q(y_k | x, a_k) p(x | x', a_k) \pi_{k-1}(x') dx' dx}. \quad (4)$$

Consider a suitable *test* function  $\psi : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$  where, for example,  $\psi$  could pick out a component of interest of the state vector we wish to estimate. The optimal sensor scheduling problem is to solve

$$\min_{A_{1:N} \in \Theta_A} J(A_{1:N}) = \mathbf{E}_{(\pi_0, A_{1:N})} \left\{ \sum_{k=1}^N \beta^{N-k} (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \right\}, \quad (5)$$

where  $\Theta_A \subset (\mathbb{R}^{n_a})^N$  and for any  $1 \leq k \leq N$  and integrable function  $h : (\mathbb{R}^{n_x})^k \times (\mathbb{R}^{n_a})^k \times (\mathbb{R}^{n_y})^k \rightarrow \mathbb{R}$ ,

$$\mathbf{E}_{(\pi_0, A_{1:k})} \{h(X_{1:k}, A_{1:k}, Y_{1:k})\} := \int h(x_{1:k}, A_{1:k}, y_{1:k}) \times \prod_{i=1}^k q(y_i | x_i, A_i) p(x_i | x_{i-1}, A_i) \pi_0(x_0) dy_{1:k} dx_{0:k}. \quad (6)$$

$\beta \in [0, 1]$  is a discount factor and its inclusion favours better tracking performance in the later epochs.

*Feedback control:* The sensor scheduling problem stated in (5) is an open-loop stochastic control problem. In order to utilise feedback, we will use the *open-loop feedback control* (OLFC) approach, which is described in detail in [1].

*Simulation and gradient based methods:* We do not have a closed-form expression for  $J$  because the filtering density  $\pi_k$  and integration with respect to it cannot be evaluated in closed-form in our general setting. To evaluate  $J(A_{1:N})$ , we could revert to state-space discretisation. One could discretise  $\mathbb{R}^{n_x}$ ,  $\mathbb{R}^{n_y}$  and derive the corresponding state evolution and observation laws, i.e. (2) and (3), for the approximating discrete problem. We may then calculate the approximation to  $J(A_{1:N})$  for any choice of actions. However this approach is computationally prohibitive and we would be limited to a coarse discretisation and a small horizon  $N$  at best. Also, it is not obvious how to choose the grid in  $\mathbb{R}^{n_x}$  and  $\mathbb{R}^{n_y}$ , since, for accuracy of the approximation, the grid should be finer in the regions where density in (6) has more mass. In [11], the HMM is discretised and a close-loop formulation of problem (5) is solved. A close-loop formulation of (5) is known as a Partially Observed Markov Decision Process (POMDP).

We propose to use simulation with Stochastic Approximation (SA) to minimise  $J(A_{1:N})$  when  $\Theta_A$  is an open (i.e. continuous) set without resorting to discretising  $\mathbb{R}^{n_x}$ ,  $\mathbb{R}^{n_y}$  or  $\Theta_A$ . SA is a gradient descent algorithm that only requires noisy estimates of the cost function gradient, i.e.,

$$A_{1:N, k+1} = A_{1:N, k} - \alpha_k (\nabla J(A_{1:N})|_{A_{1:N}=A_{1:N, k}} + \text{noise}), \quad (7)$$

where  $\nabla J(A_{1:N})$  denotes the gradient of  $J$  w.r.t.  $A_{1:N}$ . The step-size  $\alpha_k$  is a non-increasing positive sequence tending to zero. In Section III, we derive the gradient  $\nabla J$ . Once again, we do not have a closed-form expression for  $\nabla J$  for the same reasons as in  $J$ ; the filtering density  $\pi_k$  and integration with respect to it cannot be evaluated in closed-form in our general setting. We will show instead how one may obtain a low variance estimate of  $\nabla J$ , namely  $\widehat{\nabla J}$ . The noise in (7) arises precisely because we use  $\widehat{\nabla J}$  instead of  $\nabla J$ . Under suitable assumptions on the noise in (7), one can guarantee that  $A_{1:N, k}$  eventually converges to a local minimiser of  $J$ .

## III. THE COST GRADIENT AND ITS SIMULATION-BASED APPROXIMATION

In this section, we derive the gradient of the cost function (5) with respect to  $A_{1:N}$ . We then propose a suitable simulation-based approximation for optimising with SA. Keeping in mind that  $(\psi(X_k) - \langle \pi_k, \psi \rangle)^2$  is a function of the form  $h(X_{1:k}, A_{1:k}, Y_{1:k})$ , then (6) implies <sup>2</sup> that  $\mathbf{E}_{A_{1:N}} \left\{ (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \right\} = \mathbf{E}_{A_{1:k}} \left\{ (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \right\}$ . Thus,  $\nabla_{A_l} \mathbf{E}_{A_{1:N}} \left\{ (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \right\} = 0$  for  $l > k$ . For

<sup>2</sup>Because the problem (5) is solved for a fixed initial state distribution  $\pi_0$ , henceforth, we omit reference to  $\pi_0$  in the notation for  $\mathbf{E}_{(\pi_0, A_{1:N})}$  and denote the probability with respect to which this expectation is taken by  $\mathbf{P}_{A_{1:N}}$ .

$l \leq k$ , using (6),

$$\begin{aligned} & \nabla_{A_l} \int \left( \psi(x_k) - \left\langle \pi_k^{(y_{1:k}, A_{1:k})}, \psi \right\rangle \right)^2 \\ & \times \prod_{i=1}^k q(y_i | x_i, A_i) p(x_i | x_{i-1}, A_i) \pi_0(x_0) dx_{0:k} dy_{1:k} \\ & = \int \left( \psi(x_k) - \left\langle \pi_k^{(y_{1:k}, A_{1:k})}, \psi \right\rangle \right)^2 \\ & \times \nabla_{A_l} \left[ \prod_{i=1}^k q(y_i | x_i, A_i) p(x_i | x_{i-1}, A_i) \right] \pi_0(x_0) dx_{0:k} dy_{1:k} \\ & + \int \nabla_{A_l} \left[ \left( \psi(x_k) - \left\langle \pi_k^{(y_{1:k}, A_{1:k})}, \psi \right\rangle \right)^2 \right] \\ & \times \prod_{i=1}^k q(y_i | x_i, A_i) p(x_i | x_{i-1}, A_i) \pi_0(x_0) dx_{0:k} dy_{1:k}. \end{aligned}$$

The first term of the gradient is  $\mathbf{E}_{A_{1:N}} \{ (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \} \times \left[ \frac{\nabla_{A_l} q(Y_l | X_l, A_l)}{q(Y_l | X_l, A_l)} + \frac{\nabla_{A_l} p(X_l | X_{l-1}, A_l)}{p(X_l | X_{l-1}, A_l)} \right]$ . The second term of the gradient is  $\mathbf{E}_{A_{1:N}} \{ \nabla_{A_l} [ (\psi(X_k) - \langle \pi_k^{(Y_{1:k}, A_{1:k})}, \psi \rangle)^2 ] \} = -2 \times \mathbf{E}_{A_{1:N}} \{ (\psi(X_k) - \langle \pi_k^{(Y_{1:k}, A_{1:k})}, \psi \rangle) \nabla_{A_l} \langle \pi_k^{(Y_{1:k}, A_{1:k})}, \psi \rangle \} = 0$ , where the last line follows from conditioning on  $Y_{1:k}$ . It follows from the above derivation that to obtain an unbiased estimator of  $\nabla_{A_l} J(A_{1:N})$  for a given  $A_{1:N}$ , one samples a realisation of states and observations  $(Y_{1:N}, X_{0:N}) \sim \mathbf{P}_{A_{1:N}}$  and forms the following estimate,

$$\begin{aligned} \widehat{\nabla_{A_l} J}(A_{1:N}) &= \sum_{k=l}^N \beta^{N-k} \mathbf{E}_{A_{1:N}} \{ (\psi(X_k) - \langle \pi_k, \psi \rangle)^2 \} \\ & \times \left[ \frac{\nabla_{A_l} q(Y_l | X_l, A_l)}{q(Y_l | X_l, A_l)} + \frac{\nabla_{A_l} p(X_l | X_{l-1}, A_l)}{p(X_l | X_{l-1}, A_l)} \right] | Y_{1:k}, \end{aligned} \quad (8)$$

where we have added the conditioning on  $Y_{1:k}$  as it leads to a lower variance gradient estimate.<sup>3</sup>

We now describe how to implement the gradient estimate (8). In sensor scheduling applications concerning target tracking, the state process  $X_k$  is the state of the target to be tracked and often evolves independently of the action. Henceforth, we assume this independence for simplicity in presentation, i.e.  $p(X_k | X_{k-1})$ , and remark that the work may also be extended to the more general case of state evolution and control dependence.<sup>4</sup> Now, define the real-valued function called the *score* [9],  $S(y, x, a) := q(y | x, a)^{-1} \nabla_a q(y | x, a) \in \mathbb{R}^{n_a}$ . To implement (8), we see that we require both the marginal  $\pi_k$  and the full posterior  $\pi_{0:k}$  for all  $N$  epochs, i.e., for  $1 \leq k \leq N$ . We propose to approximate these quantities using a mixture Dirac delta-masses,

$$\hat{\pi}_{0:k}(x_{0:k}) := \sum_{j=1}^L w_k^{(j)} \delta_{X_{0:k}^{(j)}}(x_{0:k}), \quad (9)$$

<sup>3</sup>The variance is reduced since, for two jointly distributed random variables  $X$  and  $Y$ ,  $\text{var}(E(X|Y)) = \text{var}(X) - E(\text{var}(X|Y))$ , and  $E(\text{var}(X|Y)) > 0$ .

<sup>4</sup>In methods that use the PCRLB [5], [8], [11], even after assuming  $\{X_k\}$  evolves independently of  $\{A_k\}$ , one still needs to evaluate the expectation of derivatives of  $\ln p(X_k | X_{k-1})$  w.r.t.  $X_k$  and  $X_{k-1}$ , while this is not needed in (8).

where  $\delta_{X_{0:k}^{(j)}}$  denotes the Dirac delta-mass located at  $X_{0:k}^{(j)}$  and the *importance weights*  $\{w_k^{(j)}\}_{j=1}^L$  are non-negative scalars that sum to one. The approximation to  $\pi_k$ , namely  $\hat{\pi}_k$ , follows by marginalising  $\hat{\pi}_{0:k}$ . There are a number of ways to define such a point-mass approximation. For example, the simplest scheme would be to sample  $L$  independent state trajectory realisations  $\{X_{0:N}^{(j)}\}_{j=1}^L$  from  $(\prod_{i=1}^N p(x_i | x_{i-1})) \pi_0(x_0)$ . The importance weights would then be  $w_k^{(j)} := \frac{\int \prod_{i=1}^k q(Y_i | X_i^{(j)}, A_i) \pi_0(x_0) dx_{0:k}}{\int \prod_{i=1}^k q(Y_i | X_i^{(j)}, A_i) \pi_0(x_0) dx_{0:k}}$ . For any integrable function  $h$ ,  $\int h(x_{0:k}) \hat{\pi}_{0:k}(x_{0:k}) dx_{0:k}$  converges to  $\int h(x_{0:k}) \pi_{0:k}(x_{0:k}) dx_{0:k}$  as  $L \rightarrow \infty$  (see [2, Ch. 2] for a precise statement of the mode of convergence). Practically though, we would prefer a small sample size  $L$  and this simple scheme of sampling from the state transition model can result in the majority of the importance weights  $w_k^{(j)}$  being very small. There are number of remedies proposed for this in the Sequential Monte Carlo literature [2, Ch. 1.3.2].

Now, for a given  $A_{1:N}$ , one samples a realisation of states and observations  $(X_{0:N}, Y_{1:N}) \sim \mathbf{P}_{A_{1:N}}$  and forms the following biased estimate of  $\nabla_{A_l} J(A_{1:N})$  (8),

$$\begin{aligned} & \sum_{k=l}^N \beta^{N-k} \{ \langle \hat{\pi}_{0:k}, \psi_k^2(\cdot) S(Y_l, \cdot, A_l) \rangle \\ & + \langle \hat{\pi}_k, \psi_k \rangle^2 \langle \hat{\pi}_{0:k}, S(Y_l, \cdot, A_l) \rangle \\ & - 2 \langle \hat{\pi}_k, \psi \rangle \langle \hat{\pi}_{0:k}, \psi_k(\cdot) S(Y_l, \cdot, A_l) \rangle \}. \end{aligned} \quad (10)$$

To prove the convergence of the SA recursion, we would not be able to use standard SA results. Even though (10) is a noisy estimate of  $\nabla_{A_l} J(A_{1:N})$ , the noise is not zero-mean due to the bias of the simulation-based approximation to  $\pi_k$  and  $\pi_{0:k}$ . To assert convergence of (7) to a minima of  $J$ , we would have to gradually increase the number of samples  $L$  to remove the bias. Similar conditions are required for convergence of SA driven by sample averages [9].

Henceforth, we fix the set of  $L$  state trajectory samples  $\{X_{0:N}^{(j)}\}_{j=1}^L$  that is used to approximate  $\pi_{0:k}$  for all  $k$  and form the following approximation to  $J^5$ ,

$$\hat{J}(A_{1:N}) = \sum_{k=1}^N \beta^{N-k} \mathbf{E}_{(\pi_0, A_{1:N})} \left\{ \langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2 \right\}. \quad (11)$$

We will then derive an unbiased estimate of the gradient of  $\hat{J}$  in a similar manner to  $J$  above and minimise  $\hat{J}$  via SA. This approach is easier to analyse and we show that, under suitable assumptions, SA converges to a local minimum of  $\hat{J}$  almost surely.<sup>6</sup>

In the same way as gradient of  $J$  was derived in (8), we

<sup>5</sup>Note that by a conditioning argument,  $J(A_{1:N})$  can be written as  $\sum_{k=1}^N \beta^{N-k} \mathbf{E}_{A_{1:N}} \left\{ \langle \pi_k, \psi^2 \rangle - \langle \pi_k, \psi \rangle^2 \right\}$ .

<sup>6</sup>Since the error in the approximation  $\hat{\pi}_{0:k}$  diminishes as the sample size  $L$  increases,  $\hat{J}$  will be a good approximation to  $J$  for sufficiently large  $L$ .

have

$$\begin{aligned} \nabla_{A_l} \hat{J}(A_{1:N}) = & \\ \sum_{k=1}^N \beta^{N-k} \nabla_{A_l} \mathbf{E}_{A_{1:N}} \left\{ \langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2 \right\} = & \\ \mathbf{E}_{A_{1:N}} \left\{ \sum_{k=l}^N \beta^{N-k} (\langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2) S(Y_l, X_l, A_l) \right\} + & \\ & (12) \end{aligned}$$

$$\mathbf{E}_{A_{1:N}} \left\{ \sum_{k=l}^N \beta^{N-k} (\nabla_{A_l} \langle \hat{\pi}_k, \psi^2 \rangle - 2 \langle \hat{\pi}_k, \psi \rangle \nabla_{A_l} \langle \hat{\pi}_k, \psi \rangle) \right\}, \quad (13)$$

where<sup>7</sup>

$$\begin{aligned} \nabla_{A_l} \langle \hat{\pi}_k, \psi \rangle = & \langle \hat{\pi}_{0:k}, \psi S(Y_l, \cdot, A_l) \rangle \\ & - \langle \hat{\pi}_k, \psi \rangle \langle \hat{\pi}_{0:k}, S(Y_l, \cdot, A_l) \rangle. \end{aligned} \quad (14)$$

It is now straightforward to obtain a simulation-based approximation of  $\nabla \hat{J}(A_{1:N})$ . For a given  $A_{1:N}$ , one samples a realisation of states and observations  $(Y_{1:N}, X_{0:N}) \sim \mathbf{P}_{A_{1:N}}$  and forms the following unbiased estimate of  $\nabla_{A_l} \hat{J}(A_{1:N})$ : For  $l = 1, \dots, N$

$$\begin{aligned} S(Y_l, X_l, A_l) \sum_{k=l}^N \beta^{N-k} \left( \langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2 \right) \\ + \sum_{k=l}^N \beta^{N-k} \left( \nabla_{A_l} \langle \hat{\pi}_k, \psi^2 \rangle - 2 \langle \hat{\pi}_k, \psi \rangle \nabla_{A_l} \langle \hat{\pi}_k, \psi \rangle \right). \end{aligned} \quad (15)$$

#### A. Variance Reduction by Control Variates

The variance of the gradient approximation (15) (or (10)) is quite large, because we are approximating high dimensional integrals using simulation and more so, with moderate sample sizes. Naturally, it would be possible to reduce the variance by simply increasing the number of samples. We do not wish to do so, as our aim is to extract the most accurate estimates of the quantities of interest for a given set of samples.

Given a random variable  $W$  and a zero-mean random variable  $Z$  (control variate or CV) correlated with  $W$ , to estimate  $\mathbf{E}(W)$  we use  $W - bZ$  where  $b$  is a constant (CV constant). The estimator  $W - bZ$  is also unbiased. Furthermore, the function of  $b$ ,  $\mathbf{var}(W - bZ) = \mathbf{var}(W) - 2b\mathbf{cov}(W, Z) + b^2\mathbf{var}(Z)$ , is convex and is minimised at  $b^* = \mathbf{cov}(W, Z)/\mathbf{var}(Z)$ , which implies the variance of the estimate  $W - b^*Z$  of  $\mathbf{E}(W)$  is less than the variance of the estimate  $W$ . In the context of the gradient estimate in (15), we found in implementation that reducing the variance of the estimate of (12) was sufficient.

The score in (12) is zero-mean, i.e.  $\mathbf{E}_{(\pi_0, A_{1:N})} \{S(Y_l, X_l, A_l)\} = 0$ , and we use it as the

<sup>7</sup>It is possible to compute  $\nabla_{A_l} \langle \hat{\pi}_k, \psi \rangle$  when re-sampling is employed in  $\hat{\pi}_k$ ; resampling is commonly used in the Sequential Monte Carlo literature to yield approximations to  $\hat{\pi}_k$

CV. Doing so yields the following unbiased estimator of  $\nabla_{A_l} \hat{J}$  instead of (15),

$$\begin{aligned} \text{diag}(S(Y_l, X_l, A_l)) \left( \sum_{k=l}^N \beta^{N-k} (\langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2) \mathbf{1} - b_l \right) \\ + \sum_{k=l}^N \beta^{N-k} \left( \nabla_{A_l} \langle \hat{\pi}_k, \psi^2 \rangle - 2 \langle \hat{\pi}_k, \psi \rangle \nabla_{A_l} \langle \hat{\pi}_k, \psi \rangle \right), \end{aligned} \quad (16)$$

where  $\mathbf{1} \in \mathbb{R}^{n_a}$  and the CV constant (vector)  $b_l \in \mathbb{R}^{n_a}$  is to be determined in order to minimise the variance of the estimate. Noting that the optimal CV constant is a solution of the minimisation problem (III-A), we may employ the following SA algorithm to converge to it,

$$\begin{aligned} b_l \leftarrow b_l - \beta \text{diag}(S(Y_l, X_l, A_l)) \left( \text{diag}(S(Y_l, X_l, A_l)) b_l \right. \\ \left. - \sum_{k=l}^N \beta^{N-k} (\langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2) \mathbf{1} \right), \end{aligned} \quad (17)$$

where  $\beta$  is the step-size. Under suitable assumptions stated in Section IV, we will have  $b_l$  converging to

$$\begin{aligned} \mathbf{E}_{(\pi_0, A_{1:N})} \left\{ \text{diag}(S(Y_l, X_l, A_l))^2 \right\}^{-1} \\ \times \mathbf{E}_{(\pi_0, A_{1:N})} \left\{ \text{diag}(S(Y_l, X_l, A_l))^2 \right. \\ \left. \times \sum_{k=l}^N \beta^{N-k} (\langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2) \mathbf{1} \right\}. \end{aligned} \quad (18)$$

The same approach applies when minimising the variance of the gradient estimate (10) with control variates.

#### IV. SOLVING THE SENSOR SCHEDULING PROBLEM WITH STOCHASTIC APPROXIMATION

We now state the simulation-based algorithm that will be used to solve the sensor scheduling problem. It is a two time-scale SA algorithm to minimise  $\hat{J}$  using the reduced variance estimate of  $\nabla \hat{J}$  given by (16) and (17). We do so for the (general) case with action path subject to the equality constraint as specified in  $A_{1:N} = F(U_{1:N})$ ,  $U_{1:N} \in (\mathbb{R}^{n_u})^N$ , i.e.  $\Theta_A$  now corresponds to the range of the function  $F$ . Note that although we could also apply algorithm (21)-(24) below with the gradient estimate (10), we minimise  $\hat{J}$  instead as its convergence is easier to study; see comments immediately following (10).

We introduce the following functions to make the presentation of the main algorithm concise. For each  $A_{1:N}$ , define the functions  $h_{i, A_{1:N}} : (\mathbb{R}^{n_x})^{N+1} \times (\mathbb{R}^{n_y})^N \rightarrow (\mathbb{R}^{n_a})^N$ ,  $i = 1, 2$ , as follows:

$$\begin{aligned} h_{1, A_{1:N}}(X_{0:N}, Y_{1:N}) := \\ [S(Y_l, X_l, A_l) \sum_{k=l}^N \beta^{N-k} (\langle \hat{\pi}_k, \psi^2 \rangle - \langle \hat{\pi}_k, \psi \rangle^2)]_{l=1, \dots, N}, \end{aligned} \quad (19)$$

$$\begin{aligned} h_{2, A_{1:N}}(X_{0:N}, Y_{1:N}) := \\ \left[ \sum_{k=l}^N \beta^{N-k} (\nabla_{A_l} \langle \hat{\pi}_k, \psi^2 \rangle - 2 \langle \hat{\pi}_k, \psi \rangle \nabla_{A_l} \langle \hat{\pi}_k, \psi \rangle) \right]_{l=1, \dots, N}. \end{aligned} \quad (20)$$

Note that  $\nabla \hat{J}(A_{1:N}) = \mathbf{E}_{A_{1:N}} \{h_{1,A_{1:N}}(X_{0:N}, Y_{1:N}) + h_{2,A_{1:N}}(X_{0:N}, Y_{1:N})\}$ ,  $\in (\mathbb{R}^{n_a})^N$ .

For technical reasons concerning the convergence of the two time-scale SA algorithm below, we introduce the positive scalar valued function  $\Gamma : (\mathbb{R}^{n_a})^N \rightarrow (0, \infty)$ ,  $\Gamma(b) := (1 + |b|)^{-1}C$ , where  $C$  is a positive constant. The function  $\Gamma$  is needed to ensure that the CV constants remain bounded almost surely (details in [10]).

**The two time-scale SA algorithm for solving the sensor scheduling problem:** For conciseness, let

$$\theta = U_{1:N}, \tilde{\theta} = A_{1:N} \quad (= F(\theta)), \omega = (X_{0:N}, Y_{1:N}).$$

$$\begin{aligned} \theta_{k+1} = & \theta_k - \alpha_{k+1} \Gamma(b_k) \nabla F(\theta_k) \\ & \times (h_{1, \tilde{\theta}_k}(\omega_{k+1}) + h_{2, \tilde{\theta}_k}(\omega_{k+1}) - S_{\tilde{\theta}_k}(\omega_{k+1}) b_k), \end{aligned} \quad (21)$$

$$\begin{aligned} b_{k+1} = & b_k - \beta_{k+1} S_{\tilde{\theta}_k}^2(\omega_{k+1}) b_k \\ & + \beta_{k+1} S_{\tilde{\theta}_k}(\omega_{k+1}) h_{1, \tilde{\theta}_k}(\omega_{k+1}), \end{aligned} \quad (22)$$

$$\omega_{k+1} \sim \mathbf{P}_{\tilde{\theta}_k}, \quad (23)$$

$$\tilde{\theta}_k = F(\theta_k), \quad k \geq 0, \quad (24)$$

where

$$\begin{aligned} S_{A_{1:N,k}}(X_{0:N,k+1}, Y_{1:N,k+1}) = \\ \text{diag} \left( [S(Y_{l,k+1}, X_{l,k+1}, A_{l,k})]_{l=1, \dots, N} \right). \end{aligned} \quad (25)$$

(Note that  $U_{1:N,k} = \theta_k$ ,  $A_{1:N,k} = \tilde{\theta}_k$ ,  $(X_{0:N,k+1}, Y_{1:N,k+1}) = \omega_{k+1}$ .)

*Assumption 1:* The step-size sequences  $\{\alpha_k\}$  and  $\{\beta_k\}$  are non-negative, sum to infinity, are squared summable and, for some  $p > 0$  satisfy  $\sum_k \left(\frac{\alpha_k}{\beta_k}\right)^p < \infty$ .

Typically, the step-sizes are  $\alpha_k = k^{-\alpha}$ ,  $\beta_k = k^{-\beta}$ , where  $\alpha > \beta > 0.5$ . Thus,  $\sum_k \left(\frac{\alpha_k}{\beta_k}\right)^p < \infty$  may only be satisfied for a large positive  $p$ . Since  $\alpha_k$  tends to zero more quickly than  $\beta_k$ , the recursion for the action (21) is said to evolve on a slower time-scale than that for the CV constants (22). By having  $U_{1:N,k}$  evolve more slowly than  $b_k$ , we allow  $b_k$  to ‘track’ the optimal CV constants, which depend on the point at which the gradient  $\nabla \hat{J}$  is evaluated (see (18)). In [10], using results from [6], we establish the convergence of algorithm (21)-(24) for the choice of step-sizes in Assumption 1.

## V. APPLICATION TO OBSERVER TRAJECTORY PLANNING

In observer trajectory planning, we wish to track a maneuvering target for  $N$  epochs. At epoch  $k$ , let  $X_k$  denote the state of the target,  $A_k$  the position of the observer and  $Y_k$  the partial observation of the target state, i.e.,  $Y_k = g(X_k, A_k, V_k)$ , where  $V_k$  denotes noise. Typically, the observer has its own motion model and we let  $X_k^o$  denote state of the observer. The aim of OTP is to adaptively maneuver the observer to optimise the tracking performance the target.

We do not need to specify the target model explicitly. Our only concern is that we can sample from the model. In Section VI, we consider a maneuvering target in the examples. We require an observer model of the form  $A_{1:N} = F(U_{1:N})$  where we exert control on the observer positions  $A_{1:N}$  through the variables  $U_{1:N}$ . For instance, the accelerations of the observer could be determined from the input  $U_{1:N}$ , which will in turn determine the observer trajectory. This is precisely the model we adopt for the examples in Section VI. Let the state of the observer be  $X_k^o = [r_{x,k}^o, v_{x,k}^o, r_{y,k}^o, v_{y,k}^o]^T$ , with  $A_k = [r_{x,k}^o, r_{y,k}^o]^T$ . For example, we could assume a kinematic model for the evolution of the state,

$$X_{k+1}^o = G X_k^o + H \times c_1 \times \arctan(c_2 U_{k+1}) \quad (26)$$

where matrices  $G$  and  $H$  follow the standard definition [7]. As in the standard kinematic model, the initial state  $X_0^o$  is fixed and  $U_{k+1} \in \mathbb{R}^2$  determines the acceleration in the  $x$  and  $y$  directions. The constants  $c_1, c_2$  alter the saturation behaviour of the acceleration. The observer trajectory is completely determined once  $X_0^o$  and  $U_{1:N}$  are given. The function  $F$  is now implicitly defined by (26).

In the bearings-only model, the observation process  $\{Y_k\}_{k \geq 0} (\subset \mathbb{R})$  is generated according to  $Y_k = \arctan\left(\frac{r_{x,k} - A_k(1)}{r_{y,k} - A_k(2)}\right) + V_k$ , where  $V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_Y^2)$ . In our simulation-based framework, we require that the observation process density is known and is differentiable w.r.t.  $A_k$ . The bearings-only case is one such example. At this point we will assume that the  $x$  and  $y$  position of the target correspond to the first and third component of the state descriptor  $X_k$ , i.e.  $X_k = [r_{x,k}, \cdot, r_{y,k}, \dots]^T$ , which is usual convention in the literature.

### A. Convergence For Bearings-Only Tracking

*Proposition 2:* If the support of random variables  $X_{0:N}$  and the range of function  $F$  do not intersect then, we have the desired convergence of two time-scale SA for OTP. That is, almost surely,  $\lim_k \left| b_k - \overline{S^2(A_{1:N,k})}^{-1} \overline{S} \times \overline{h_1(A_{1:N,k})} \right| = 0$  and  $\liminf_k \left| \nabla(\hat{J} \circ F)(U_{1:N,k}) \right| = 0$ .

*Proof:* See [10].  $\blacksquare$

It is interesting to note that the scenario in which the support of  $X_{0:N}$  and the range of  $F$  do not intersect is a standard setting studied by previous works on OTP for bearings-only observations (see references in the Introduction), and hence the conditions of Proposition 2 are not restrictive for the application.

## VI. NUMERICAL EXAMPLE

In all examples below,  $\psi(X_k) = w_1 X_k(1) + w_2 X_k(3)$ . Weights  $w_1, w_2 \in [0, 1]$  are selected to trade-off accuracy in tracking the  $x$ - and  $y$ - coordinates. We solve for the optimal open-loop observer trajectory under a variety of tracking scenarios, namely, with a fast observer, a slow observer and two cooperating observers. Open loop feedback control is implemented for the two observers case. All examples consider a maneuvering target.

*Fast Observers:* The setting for this example is a maneuvering target that is to be tracked by a single fast observer and two cooperating fast observers. The term fast is because in the subsequent example the observer is significantly more constrained in its motion. The observer motion model is given in (26), with a fast or slow observer defined by setting constant  $c_1$ . In Figure 1(a) the optimal open-loop trajectory of the observer is plotted for a horizon 7 problem. The maneuvering target trajectory is also shown. The cloud of particles surrounding the maneuvering target are target trajectory samples drawn from the target dynamic model (2). Note that the target maneuver was intentionally chosen to be far more drastic than is predicted by its model. This was done to contrast the constructed open-loop and open loop feedback control trajectories. In Figure 1(b), we show the difference in the optimal open-loop trajectory obtained when there are two fast observers. Figure 1(c) shows the OLFC obtained for the same two fast observers. The cloud of particles shown is now the particle filter tracking the maneuvering target. We note that the OLFC trajectory performs more maneuvers than the equivalent open loop one.

*Slow Observers:* Figure 2(a) shows the optimal open-loop trajectory of one slow observer, and Figure 2(b) that of two cooperating slow observers. Note that a single slow observer is obliged to do more maneuvers to improve the tracking performance since it is significantly more constrained in motion. Figure 2(c) shows the OLFC obtained. Note that the two observers maneuver in Figure 2(c) much more than in Figure 2(b) as they are responding to the target maneuver.

## VII. CONCLUSION

In this paper, we proposed a novel simulation-based method to solve the sensor scheduling problem for the case in which the state, observation and action spaces are continuous valued vectors. This general continuous state-space case is important as it is the natural framework for many applications, like observer trajectory planning. We presented a novel simulation-based method that used two timescale stochastic approximation to find optimal actions and stated convergence results for the bearings-only tracking problem.

## REFERENCES

- [1] D.P. Bertsekas, *Dynamic programming and optimal control*. Belmont: Athena Scientific, 1995.
- [2] A. Doucet, J.F.G. de Freitas and N.J. Gordon *Sequential Monte Carlo methods in practice*. New York: Springer, 2001.
- [3] R.J. Evans, V. Krishnamurthy and G. Nair, "Sensor adaptive target tracking over variable bandwidth networks," in *Model Identification and Adaptive Control*, G.C. Goodwin (Ed.), Springer-Verlag, London, pp. 115-124, 2001.
- [4] P.W. Glynn and R. Szechtman, "Some new perspectives on the method of control variates," in *Monte Carlo and Quasi-Monte Carlo Methods 2000*, K.T. Fang, F.J. Hickernell and H. Niederreiter, Eds., Springer-Verlag, Berlin, pp. 27-49, 2002.
- [5] M.L. Hernandez, T. Kirubarajan, Y. Bar-Shalom, "Multisensor resource deployment using posterior Cramer-Rao bounds," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 40, no. 2, April, 2004
- [6] V. R. Konda and J. N. Tsitsiklis, "Linear Stochastic Approximation Driven by Slowly Varying Markov Chains", *Systems and Control Letters*, Vol. 50, No. 2, 2003, pp. 95-102.

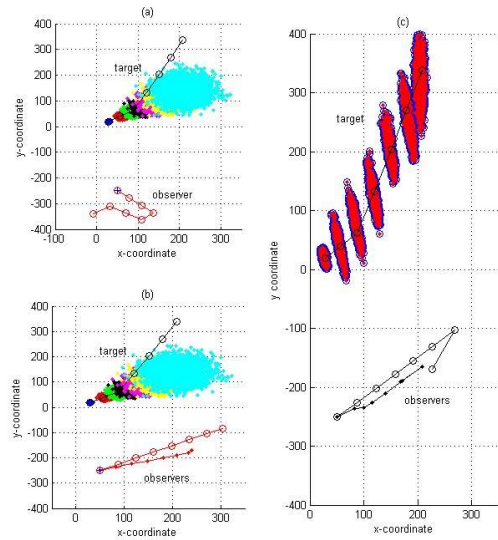


Fig. 1. (a) One observer open loop trajectory with target, (b) Two observers open loop trajectory with target, (c) Two observers OLFC trajectory with target and particle filter. Fast observers commence from (50, -250). Target starts at (0, 0) and moves in the north-west direction.

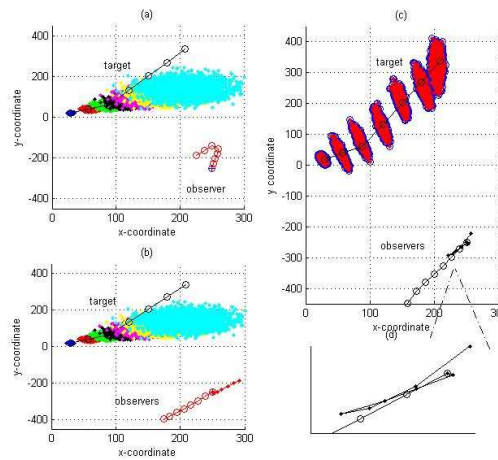


Fig. 2. (a) One observer open loop trajectory with target, (b) Two observers open loop trajectory with target, (c) Two observers OLFC trajectory with target and bootstrap filter, (d) Magnification of the OLFC trajectory of one of the observers. Slow observers commence from coordinate (250, -250).

- [7] A. Logothetis, A. Isaksson and R.J. Evans "An information theoretic approach to observer path design for bearings-only tracking," in *IEEE Conf. Decision Control*, 1997, pp. 3132-3137.
- [8] S. Paris, J.-P. Le Cadre, "Planification for terrain-aided navigation," in *Proceedings of the 5th International Conference on Information Fusion*, pp. 1007-1014, Annapolis, Maryland, 2002.
- [9] G.Ch. Pflug, *Optimization of stochastic models: the interface between simulation and optimization*. Boston: Kluwer, 1996.
- [10] S. Singh, N. Kantas, B. Vo, A. Doucet and R. Evans, "On the convergence of a stochastic optimisation algorithm for optimal observer trajectory planning," Cambridge University Technical Report, 2005, Source: <http://www-sigproc.eng.cam.ac.uk/ss40/papers>.
- [11] O. Tremois and J.-P. LeCadre, "Optimal observer trajectory in bearings-only tracking for manoeuvring sources," *IEE Proc. Radar, Sonar Navig.*, vol. 146, no. 1, pp. 31-39, Feb. 1999.
- [12] H. W.J. Lee, K.L. Teo and E.B. Lim, "Sensor scheduling in continuous time," *Automatica*, vol. 37, pp. 2017-2023, 2001.