

Improving Population Monte Carlo: Alternative Weighting and Resampling Schemes

Víctor Elvira

Dep. of Signal Theory and Communic., Universidad Carlos III de Madrid, Leganés (Spain).

Luca Martino

Dep. of Mathematics and Statistics, University of Helsinki, Helsinki (Finland).

David Luengo

Dep. of Signal Theory and Communic., Universidad Politécnica de Madrid, Madrid (Spain).

Mónica F. Bugallo

Dep. of Electrical and Computer Eng., Stony Brook University, NY (USA).

Abstract

Population Monte Carlo (PMC) sampling methods are powerful tools for approximating distributions of static unknowns given a set of observations. These methods are iterative in nature: at each step they generate samples from a proposal distribution and assign them weights according to the importance sampling principle. Critical issues in applying PMC methods are the choice of the generating functions for the samples and the avoidance of the sample degeneracy. In this paper, we propose three new schemes that considerably improve the performance of the original PMC formulation by allowing for better exploration of the space of unknowns and by selecting more adequately the surviving samples. A theoretical analysis is performed, proving the superiority of the novel schemes in terms of variance of the associated estimators and preservation of the sample diversity. Furthermore, we show that they outperform other state of the art algorithms (both in terms of mean square error and robustness w.r.t. initialization) through extensive numerical simulations.

Keywords: Population Monte Carlo, adaptive importance sampling, proposal distribution, resampling.

1. Introduction

Bayesian signal processing, which has become very popular over the last years in statistical signal processing, requires computing distributions of unknowns conditioned on observations (and moments of them). Unfortunately, these distributions are often impossible to obtain analytically in many real-world challenging problems. An alternative is then to resort to Monte Carlo (MC) methods, which approximate the target distributions with random measures composed of samples and associated weights [1].

Preprint submitted to Signal Processing

July 12, 2016

A well-known class of MC methods are those based on the adaptive importance sampling (AIS) mechanism, such as Population Monte Carlo (PMC) algorithms [2, 3], which have been used in missing data, tracking, and biological applications, among others [4, 5, 6, 7, 8]. In these methods, a population of probability density functions (pdfs) is adapted for approximating a target distribution through an iterative importance sampling procedure. AIS is often preferred to other MC schemes, such as Markov Chain Monte Carlo (MCMC), since they present several advantages. On the one hand, all the generated samples are employed in the estimation (e.g., there is no “burn-in” period). On the other hand, the corresponding adaptive schemes are more flexible, since they present fewer theoretical issues than adaptive MCMC algorithms. Namely, the convergence of AIS methods can usually be guaranteed under mild assumptions regarding the tails of the distributions and the stability of the adaptive process, whereas adaptive MCMC schemes must be designed very carefully, since the adaptation procedure can easily jeopardize the ergodicity of the chain (e.g., see [9] or [1, Section 7.6.3]).

The most characteristic feature in PMC [3] is arguably the use of resampling procedures for adapting the proposal pdfs (see for instance [10] for a review of resampling methods in particle filtering). The resampling step is a fast, often dimensionality-free, and easy way of adapting the proposal pdfs by using information about the target. However, resampling schemes present some important drawbacks, such as the sample impoverishment. At the resampling step, the proposal pdfs with poor performance (i.e., with low associated weights) are likely to be removed, thus yielding a reduction of diversity. Since the publication of the standard PMC [3], several variants have been considered, partly in an attempt to mitigate this issue. In the D-kernel algorithm [11, 12], the PMC kernel is a mixture of different kernels and the weights of the mixture are iteratively adapted in an implicit expectation-maximization (EM) algorithm. This procedure is refined through a double Rao-Blackwellization in [13]. The mixture population Monte Carlo algorithm (M-PMC) proposed in [14] also adapts a mixture of proposal pdfs (weights and parameters of the kernels). The M-PMC belongs to the family of AIS methods, since it iteratively draws the samples from the mixture that is updated at every iteration without an explicit resampling step. Since drawing from the mixture can be interpreted as an implicit multinomial resampling, this method retains some similarities with the standard PMC scheme. A nonlinear transformation of the importance weights in the PMC framework has also been proposed in [15]. Other sophisticated AIS schemes, such as the AMIS [16] and the APIS [17] algorithms, have been recently proposed in the literature.

In this paper, we study three novel PMC schemes that improve the performance of standard PMC approach by allowing a better exploration of the space of unknowns and by reducing the variance of the estimators. These alternatives can be applied within some other sophisticated AIS approaches as well, such as the SMC samplers [18]. For this reason, we mainly compare them with the standard PMC [3], since the novel schemes could be automatically combined with the more sophisticated AIS techniques.

First of all, we introduce an alternative form of the importance weights, using a mixture of the proposal pdfs in the denominator of the weight ratio. We provide an exhaustive theoretical analysis, proving the unbiasedness and consistency of the resulting estimator, and showing the reduction in the variance of the estimator w.r.t. the estimator obtained using the standard weights. We also prove that the use of this mixture decreases the averaged mismatch between the numerator (target) and the function in the denominator of the IS weight in terms of L_2 distance. Moreover, we test this alterna-

tive scheme in different numerical simulations, including an illustrative toy example in Section 5.1, showing its practical benefit.

In the second proposed scheme, we generate several samples from every proposal pdf (not only one, as in PMC) and then we resample them jointly (all the samples at once, keeping fixed the total number of proposal pdfs). In the third proposed scheme, we consider again the generation of several samples from every proposal pdf, but the resampling is performed separately on the set of samples coming from each proposal, therefore guaranteeing that there will be exactly one representative from each of the individual mixture components in the random measure.

We show, through extensive computer simulations in several different scenarios, that the three newly proposed variants provide a substantial improvement compared to the standard PMC. In addition, we test the proposed variants on a standard implementation of the SMC samplers [18], showing also an improvement of the performance. On the one hand, they yield unbiased estimators with a reduced variance, as also proved theoretically. On the other hand, they outperform the standard PMC in terms of preservation of sample diversity and robustness w.r.t initialization and parameter choice.

2. Problem Statement

Let us consider the variable of interest, $\mathbf{x} \in \mathbb{R}^{D_x}$, and let $\mathbf{y} \in \mathbb{R}^{D_y}$ be the observed data. In a Bayesian framework, the posterior probability density function (pdf), here referred as *target*, contains all the information about the parameters of interest and is defined as

$$\tilde{\pi}(\mathbf{x}|\mathbf{y}) = \frac{\ell(\mathbf{y}|\mathbf{x})p_0(\mathbf{x})}{Z(\mathbf{y})} \propto \pi(\mathbf{x}) = \ell(\mathbf{y}|\mathbf{x})p_0(\mathbf{x}), \quad (1)$$

where $\ell(\mathbf{y}|\mathbf{x})$ is the likelihood function, $p_0(\mathbf{x})$ is the prior pdf, and $Z(\mathbf{y})$ is the model evidence or partition function (useful in model selection).¹ The goal is to compute some moment of \mathbf{x} , i.e., an integral measure w.r.t. the target pdf,

$$I = \frac{1}{Z} \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}, \quad (2)$$

where f can be any square integrable function of \mathbf{x} w.r.t. $\pi(\mathbf{x})$, and $Z = \int \pi(\mathbf{x})d\mathbf{x}$.²

In many practical applications, both the integral (2) and Z cannot be obtained in closed form and must be approximated. Importance sampling methods allow for the approximation of both quantities by a set of properly weighted samples.

3. Population Monte Carlo (PMC)

3.1. Description of the original PMC algorithm

The PMC method [3] is a well-known iterative adaptive importance sampling technique. At each iteration it generates a set of N samples $\{\mathbf{x}_i^{(t)}\}_{i=1}^N$, where t denotes the iteration number and i denotes the sample index. In order to obtain the samples, the

¹From now on, we remove the dependence on \mathbf{y} in order to simplify the notation.

²Let us recall that $f(\mathbf{x})$ is square integrable w.r.t. $\pi(\mathbf{x})$ if $f(\mathbf{x}) \in L^2_\pi$, i.e., if $\int_{\mathcal{X}} f(\mathbf{x})^2 \pi(\mathbf{x}) d\mathbf{x} < \infty$.

original PMC algorithm makes use of a collection of proposal densities $\{q_i^{(t)}(\mathbf{x})\}_{i=1}^N$, with each sample being drawn from a different proposal, $\mathbf{x}_i^{(t)} \sim q_i^{(t)}(\mathbf{x})$ for $i = 1, \dots, N$. Then, they are assigned an importance weight, computed as $w_i^{(t)} = \frac{\pi(\mathbf{x}_i^{(t)})}{q_i^{(t)}(\mathbf{x}_i^{(t)})}$, i.e., the weight of a particular sample represents the ratio between the evaluation, at the sample value, of the target distribution and the evaluation at the sample value of the proposal used to generate it. The method proceeds iteratively (up to the maximum iteration step considered, T), building a global importance sampling estimator using different proposals at every iteration. The new proposals are obtained by updating the set of proposals in the previous iteration.

There are two key issues in the application of PMC methods: the adaptation of the proposals from iteration to iteration and the way resampling is applied. The latter is critical to avoid the degeneracy of the random measure, i.e., to avoid a few particles having extremely large weights and the rest negligible ones [1, 19]. Through the resampling procedure one selects the most promising streams of samples from the first iteration up to the current one. Several resampling procedures have been proposed in the literature [20, 21]. In the standard PMC [3], multinomial resampling is the method of choice, and consists of sampling N times from the discrete probability mass defined by the normalized weights. As a result of this procedure, the new set of parameters used to adapt the proposals for the generation of samples in the next iteration is selected. In summary, the standard PMC technique consists of the steps shown in Table 1.

3.2. Estimators and consistency

All the generated samples can be used to build a global approximation of the target. This can be done by first normalizing all the weights from all the iterations,

$$\bar{\rho}_i^{(t)} = \frac{w_i^{(t)}}{\sum_{\tau=1}^t \sum_{j=1}^N w_j^{(\tau)}}, \quad t = 1, \dots, T, \quad i = 1, \dots, N, \quad (8)$$

and then providing the pairs $\{\mathbf{x}_i^{(t)}, \bar{\rho}_i^{(t)}\}$ for $i = 1, \dots, N$ and $t = 1, \dots, T$. This procedure to compute the weights is equivalent to applying a static importance sampling technique that considers NT different proposal pdfs and all the corresponding samples. If the normalizing constant Z is known, the integral in Eq. (2) is approximated by the unbiased estimator

$$\hat{I}_t = \frac{1}{tN} \frac{1}{Z} \sum_{\tau=1}^t \sum_{j=1}^N w_j^{(\tau)} f(\mathbf{x}_j^{(\tau)}). \quad (9)$$

When the normalizing constant is unknown, the unbiased estimate of Z is substituted in Eq. (9), yielding the self-normalized estimator

$$\tilde{I}_t = \sum_{\tau=1}^t \sum_{j=1}^N \bar{\rho}_j^{(\tau)} f(\mathbf{x}_j^{(\tau)}) = \frac{1}{tN} \frac{1}{\hat{Z}_t} \sum_{\tau=1}^t \sum_{j=1}^N w_j^{(\tau)} f(\mathbf{x}_j^{(\tau)}), \quad (10)$$

where

$$\hat{Z}_t = \frac{1}{tN} \sum_{\tau=1}^t \sum_{j=1}^N w_j^{(\tau)}, \quad (11)$$

is the unbiased estimate of the normalizing constant.

Table 1: **Standard PMC algorithm [3].**

<p>1. [Initialization]: Select the parameters defining the N proposals:</p> <ul style="list-style-type: none"> • The adaptive parameters $\mathcal{P}^{(1)} = \{\boldsymbol{\mu}_1^{(1)}, \dots, \boldsymbol{\mu}_N^{(1)}\}$. • The set of static parameters, $\{\mathbf{C}_i\}_{i=1}^N$. <p>E.g., if the proposals were Gaussian distributions one could select the adapting parameters in $\mathcal{P}^{(1)}$ as the means of the proposals (that would be updated through the iterations) and the static parameters $\{\mathbf{C}_i\}_{i=1}^N$ as their covariances [3].</p> <p>2. [For $t = 1$ to T]:</p> <p>(a) Draw one sample from each proposal pdf,</p> $\mathbf{x}_i^{(t)} \sim q_i^{(t)}(\mathbf{x} \boldsymbol{\mu}_i^{(t)}, \mathbf{C}_i), \quad i = 1, \dots, N. \quad (3)$ <p>(b) Compute the importance weights,</p> $w_i^{(t)} = \frac{\pi(\mathbf{x}_i^{(t)})}{q_i^{(t)}(\mathbf{x}_i^{(t)} \boldsymbol{\mu}_i^{(t)}, \mathbf{C}_i)}, \quad i = 1, \dots, N, \quad (4)$ <p>and normalize them,</p> $\bar{w}_i^{(t)} = \frac{w_i^{(t)}}{\sum_{j=1}^N w_j^{(t)}}, \quad i = 1, \dots, N. \quad (5)$ <p>(c) Perform multinomial resampling by drawing N independent parameters $\boldsymbol{\mu}_i^{(t+1)}$ from the discrete probability random measure,</p> $\hat{\pi}_t^N(\mathbf{x}) = \sum_{i=1}^N \bar{w}_i^{(t)} \delta(\mathbf{x} - \mathbf{x}_i^{(t)}). \quad (6)$ <p>The new set of adaptive parameters defining the next population of proposals becomes</p> $\mathcal{P}^{(t+1)} = \{\boldsymbol{\mu}_1^{(t+1)}, \dots, \boldsymbol{\mu}_N^{(t+1)}\}. \quad (7)$ <p>3. [Output, $t = T$]: Return the pairs $\{\mathbf{x}_i^{(t)}, \bar{\rho}_i^{(t)}\}$, with $\bar{\rho}_i^{(t)}$ given by Eq. (8), for $i = 1, \dots, N$ and $t = 1, \dots, T$.</p>

4. Improved PMC schemes

In the following, we introduce several alternative strategies that decrease the variance of the estimators by exploiting the mixture perspective, and improve the diversity of the population w.r.t. to the standard PMC. More specifically, we study three different PMC schemes: one related to the strategy for calculating the weights and the other two based on modifying the way in which the resampling step is performed. Although we concentrate on the standard PMC, we remark that these alternative schemes can be directly applied or combined in other more sophisticated PMC algorithms. Moreover, the alternative schemes can be easily implemented in other Monte Carlo methods with resampling steps, such as the Sequential Monte Carlo (SMC) samplers [18], as we show in Sections 5.2 and 5.3.

4.1. Scheme 1: Deterministic mixture PMC (DM-PMC)

The underlying idea of PMC is to perform a good adaptation of the location parameters $\boldsymbol{\mu}_i^{(t)}$, i.e., where the proposals of the next iteration will be centered (e.g., if $q_i^{(t)}$ is a Gaussian pdf, then $\boldsymbol{\mu}_i^{(t)}$ is its mean). These parameters are obtained at each iteration by sampling from $\hat{\pi}_{t-1}^N$ in Eq. (6) (i.e., via resampling), which is a random measure that approximates the target distribution, i.e., $\boldsymbol{\mu}_i^{(t)} \sim \hat{\pi}_{t-1}^N$. As a direct consequence of the strong law of large numbers, $\hat{I}_t \rightarrow I$ almost surely (a.s.) as $N \rightarrow \infty$ under very weak assumptions [22] (the support of the proposal includes the support of the target and $I < \infty$). Furthermore, by setting $f_{\mathbf{z}}(\mathbf{X}) = \mathbb{I}(\mathbf{X} \leq \mathbf{z})$, where $\mathbf{X} = [X_1, \dots, X_{D_x}]$, $\mathbf{z} = [z_1, \dots, z_{D_x}]$, and $\mathbb{I}(\mathbf{X} \leq \mathbf{z})$ is defined as

$$\mathbb{I}(\mathbf{X} \leq \mathbf{z}) = \prod_{d=1}^{D_x} \mathbb{I}(X_d \leq z_d),$$

where $\mathbb{I}(X_d \leq z_d)$ denotes the indicator function for the d -th component ($1 \leq d \leq D_x$) of the variable of interest,

$$\mathbb{I}(X_d \leq z_d) = \begin{cases} 1, & X_d \leq z_d; \\ 0, & X_d > z_d, \end{cases}$$

then $I = I(\mathbf{z})$ becomes the multi-variate cumulative distribution function (cdf) of $\pi(\mathbf{z})$. Consequently, since $\hat{I}_t(\mathbf{z}) \rightarrow I(\mathbf{z})$ a.s. for any value of \mathbf{z} as $N \rightarrow \infty$ [Geweke,1989], $\boldsymbol{\mu}_i^{(t)} \sim \pi(\mathbf{x})$ a.s. as $N \rightarrow \infty$. In short, since the cdf associated to $\hat{\pi}_{t-1}^N(\mathbf{x})$ (which is the pdf used for resampling) converges to the target cdf (i.e., the cdf associated to $\pi(\mathbf{x})$) as $N \rightarrow \infty$, then the outputs of the resampling stage (i.e., the means $\boldsymbol{\mu}_i^{(t)}$) are asymptotically distributed as the target.

Therefore, the equally-weighted mixture of the set of proposals at the t -th iteration, given by

$$\psi^{(t)}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N q_i^{(t)}(\mathbf{x} | \boldsymbol{\mu}_i^{(t)}, \mathbf{C}_i), \quad (12)$$

can be seen as a kernel density approximation of the target pdf, where the proposals, $\{q_i^{(t)}(\mathbf{x} | \boldsymbol{\mu}_i^{(t)}, \mathbf{C}_i)\}_{i=1}^N$, play the role of the kernels [23, Chapter 6]. In general, this estimator has non-zero bias and variance, depending on the choice of q , \mathbf{C}_i , and the number of samples, N . However, for a given value of N , there exists an optimal choice of \mathbf{C}_i^* which provides the minimum Mean Integrated Square Error (MISE) estimator [24]. Using this optimal covariance matrix \mathbf{C}_i^* , it can be proved that

$$\psi^{(t)}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N q_i^{(t)}(\mathbf{x} | \boldsymbol{\mu}_i^{(t)}, \mathbf{C}_i^*) \rightarrow \tilde{\pi}(\mathbf{x}) \quad (13)$$

pointwise as $N \rightarrow \infty$ [24]. Hence, resampling naturally leads to a concentration of the proposals around the modes of the target for large values of N .

Therefore, since the performance of an importance sampling method relies on the discrepancy between the numerator (the target) and the denominator (usually, the proposal

pdf), a reasonable choice for calculating the importance weights is

$$w_i^{(t)} = \frac{\pi(\mathbf{x}_i^{(t)})}{\psi^{(t)}(\mathbf{x}_i^{(t)})} = \frac{\pi(\mathbf{x}_i^{(t)})}{\frac{1}{N} \sum_{j=1}^N q_j^{(t)}(\mathbf{x}_i^{(t)} | \boldsymbol{\mu}_j^{(t)}, \mathbf{C}_j)}, \quad (14)$$

where, as opposed to Eq. (4), the complete mixture of proposals $\psi(\mathbf{x})$ is accounted for in the denominator.

4.1.1. Theoretical justification

The first justification for using these *deterministic mixture* (DM) weights is merely mathematical, since the estimator \hat{I}_t of Eq. (9) with these weights is also unbiased (see the proof in Appendix A.2). The main advantage of this new scheme is that it yields more efficient estimators, i.e. with less variance, combining the deterministic mixture sampling (as in standard PMC) with the weight calculation that accounts for the whole mixture. Namely, the estimator \hat{I}_t in Eq. (9), computed using the DM approach, has less variance than the estimator obtained by the standard PMC, as proved in Appendix A.3 for any target and set of proposal pdfs. These DM weights have been explored in the literature of multiple importance sampling (see for instance the balance heuristic strategy of [25, Section 3.3.] or the deterministic mixture approach of [26, Section 4.3.]).

The intuition behind the variance reduction is clear in a multi-modal scenario, where different proposals have been successfully adapted covering the different modes, and therefore, the whole mixture of proposals has less mismatch w.r.t. the target than each proposal separately. Indeed, it can be easily proved that the mismatch of the whole mixture w.r.t to the target is always less than the average mismatch of each proposal. More precisely, let us consider the L_p functional distance (with $p > 1$) among the target and an arbitrary function $g(\mathbf{x})$,

$$D_p(\tilde{\pi}(\mathbf{x}), g(\mathbf{x})) = \left[\int |\tilde{\pi}(\mathbf{x}) - g(\mathbf{x})|^p d\mathbf{x} \right]^{1/p}, \quad (15)$$

and let us recall Jensen's inequality [27],

$$\varphi \left(\sum_{i=1}^N \alpha_i z_i \right) \leq \sum_{i=1}^N \alpha_i \varphi(z_i), \quad (16)$$

which is valid for any convex function $\varphi(\cdot)$, any set of non-negative weights α_i such that $\sum_{i=1}^N \alpha_i = 1$, and any collection of points $\{z_i\}_{i=1}^N$ in the support of φ . Then, by using Jensen's inequality in Eq. (16) with $\varphi(z(\mathbf{x})) = \left[\int |z(\mathbf{x})|^p d\mathbf{x} \right]^{1/p}$, $\alpha_i = \frac{1}{N}$ and $z_i = \tilde{\pi}(\mathbf{x}) - q_i(\mathbf{x})$, it is straightforward to show that

$$D_p(\tilde{\pi}(\mathbf{x}), \psi(\mathbf{x})) = D_p \left(\tilde{\pi}(\mathbf{x}), \frac{1}{N} \sum_{i=1}^N q_i(\mathbf{x}) \right) \leq \frac{1}{N} \sum_{i=1}^N D_p(\tilde{\pi}(\mathbf{x}), q_i(\mathbf{x})). \quad (17)$$

Indeed, although we have focused on the L_p distance, the proof is valid for any distance function which is based on a norm (i.e., any distance s.t. $\varphi(z(\mathbf{x})) = \|z(\mathbf{x})\|$ for some norm $\|\cdot\|$), since every norm is a convex function.

Another benefit of the DM-PMC scheme is the improvement in the exploratory behavior of the algorithm. Namely, since the weights in DM-PMC take into account all the proposals (i.e., the complete mixture) for their calculation, they temper the overrepresentation of high probability areas of the target. Note that, as a consequence of the variance reduction of the DM weights, the effective sample size in DM-PMC in a specific iteration is larger than with the standard IS weights.³ The expression $\widehat{ESS} = \frac{1}{\sum_n \bar{w}_n^2}$ is widely used as a sample approximation of the effective sample size (see its derivation in [28]). Therefore, if the true underlying effective sample size (the ratio of variances) is larger with the DM weights (than with the standard IS weights), a similar behavior can be considered for \widehat{ESS} . As a consequence, the diversity loss associated to the resampling step is reduced by using with the DM weights. See [29] for a more detailed discussion about effective sample size in static multiple importance sampling schemes.

4.1.2. Computational complexity discussion

In this DM-PMC scheme, the performance is improved at the expense of an increase in the computational cost (in terms of proposal evaluations) in the calculation of the weights. However, it is crucial to note that all the proposed schemes keep the same number of evaluations of the target as in the standard PMC. Hence, if the target evaluation is much more costly than the evaluation of the proposal pdfs (as it often happens in practical applications), the increase in computational cost can be negligible in many scenarios of interest. Note that other adaptive multiple IS algorithms, e.g. [12, 14, 13], also increase the number of proposal evaluations, and they state that the most significant computational cost is associated to the evaluation of the target (see this argument in [14, Section 2.2.]).

Finally, note that the variant *partial*-DM proposed in [30] within the static multiple IS framework, could be easily adapted to the DM-PMC. In this weighting scheme, a partition (forming subsets) of the set of proposals is a priori performed. The weight of each sample only accounts at the denominator for a subset of proposals, i.e., reducing the number of proposal evaluations. This variant achieves an intermediate point in the complexity-performance tradeoff, between the standard weights and the DM weights.

4.1.3. Comparison with other methods

Note that other methods also use a mixture of proposals at the denominator of the weights. For instance, in the D-kernel of [12], each sample is drawn from a mixture of D kernels (proposals), and this same mixture is evaluated at the denominator of each weight. Nevertheless, note that these D kernels are centered at the same position, and the weight of each sample ignores the locations of the $N - 1$ proposals. In the M-PMC of [14], a single mixture is used for sampling and weighting the N samples at each iteration. Note that this method does not use an explicit resampling step, and the mixture is completely adapted (weights, means, and covariances).

In the sequel, we adopt the weights of Eq. (14) for the other two proposed PMC schemes due to their theoretical and practical advantages discussed above.

³The effective sample size is the number of independent samples drawn from the target distribution that are equivalent (in terms of variance of the estimators) to the performance of the N samples used in the importance sampling estimator.

4.2. Scheme 2: Multiple samples per mixand with global resampling (GR-PMC)

We propose to draw K samples per individual proposal or mixand, instead of only one as done in the standard PMC algorithm. Namely,

$$\mathbf{x}_{i,k}^{(t)} \sim q_i(\mathbf{x}|\boldsymbol{\mu}_i^{(t)}, \mathbf{C}_i) \quad (18)$$

for $i = 1, \dots, N$ and $k = 1, \dots, K$. Then, we compute the corresponding DM weights as in (14),

$$w_{i,k}^{(t)} = \frac{\pi(\mathbf{x}_{i,k}^{(t)})}{\frac{1}{N} \sum_{j=1}^N q_j^{(t)}(\mathbf{x}_{i,k}^{(t)}|\boldsymbol{\mu}_j^{(t)}, \mathbf{C}_j)}. \quad (19)$$

Therefore, at each iteration we have a set of KN generated samples, i.e., $\mathcal{X}^{(t)} = \{\mathbf{x}_{1,1}^{(t)}, \dots, \mathbf{x}_{1,K}^{(t)}, \dots, \mathbf{x}_{N,1}^{(t)}, \dots, \mathbf{x}_{N,K}^{(t)}\}$. Resampling is performed in the same way as in standard PMC, although now the objective is to downsample, from KN samples to N samples, according to the normalized weights,

$$\bar{w}_{i,k}^{(t)} = \frac{w_{i,k}^{(t)}}{\sum_{j=1}^N \sum_{\ell=1}^K w_{j,\ell}^{(t)}}. \quad (20)$$

We refer to this type of resampling as *global* resampling, since all the samples, regardless of the proposal used to generate them, are resampled together. After resampling, a new set of adapted parameters for the next iteration, $\mathcal{P}^{(t+1)} = \{\boldsymbol{\mu}_1^{(t+1)}, \dots, \boldsymbol{\mu}_N^{(t+1)}\}$, is obtained. Note that, through this paper, for sake of simplicity in the explanation of the proposed improvements, we use the standard multinomial resampling, but other resampling schemes that reduce the path-degeneracy problem can be considered instead, e.g. the residual or stratified resampling, (see [31, 20]).

The PMC algorithms suffer from sample impoverishment, which is a side effect inherent to adaptive algorithms with resampling steps such as SMC samplers or particle filters (see for instance [32, Section V-C] or [33, Section 2]). In other words, there is a diversity reduction of the samples after the resampling step (in a very adverse scenario, the N resampled samples can be N copies of the sample). The sample impoverishment of the standard PMC is illustrated in Fig. 4, Fig. 5, and Fig. A.7, where the increase of diversity of the algorithms proposed in this paper is shown by numerical simulations. These figures correspond to the example of Section 5.2 and will be properly introduced below. In multimodal scenarios, proposals of the standard PMC that are exploring areas with negligible probability masses are very likely to be removed before they find unexplored relevant areas. If we draw K samples per proposal, the samples of a well-placed proposal will have similarly high weights, but as for the explorative proposals, increasing K also increases their chances of discovering local relevant features of the target $\tilde{\pi}(\mathbf{x})$. Then, the GR-PMC promotes the local exploration of the explorative proposals, increasing the chances of not being removed in the resampling step. Figures 4 and 5 show the reduction of path-degeneracy of GR-PMC in a multimodal scenario, and they will be properly explained in the example of Section 5.2.

Note that using $K > 1$ does not entail an increase in the computational cost w.r.t. the standard PMC or DM-PMC (where $K = 1$) if the number of evaluations of the target is fixed to $L = KNT$. Indeed, since the number of resampling stages is reduced to

$T = L/(KN)$, the computational cost decreases, although at the expense of performing less adaptation steps than for $K = 1$. Therefore, for a fixed budget of target evaluations L and a fixed number of proposals N , one must decide whether to promote the local exploration (possibly reducing the path degeneracy) by increasing K , or performing more adaptation steps T . Thus, there is a trade-off between local and global exploration as the numerical experiments will also show in Section 5. This suggests that, for a fixed computational budget L , there exists an optimal value of samples per proposal and iteration, K^* , which will also depend on the target and cannot be found analytically. This issue can be partially addressed through the use of local resampling, as shown in the following section.

4.3. Scheme 3: Multiple samples per proposal with local resampling (LR-PMC)

Consider again K samples generated from each proposal pdf. In this alternative scheme, the estimators are built as in GR-PMC, i.e., with the weights of Eq. (19). Nevertheless, unlike the previous method, here the resampling step is performed independently for each proposal. Namely, at the t -th iteration, K samples are drawn from each of the N proposal pdfs, and N parallel resampling procedures are independently performed within each subset of K samples (see Fig. 1 for a visual comparison of both resampling schemes). More precisely, the adaptive parameter for the next iteration of the i -th proposal, $\boldsymbol{\mu}_i^{(t+1)}$ for $i = 1, \dots, N$, is resampled from the set

$$\mathcal{X}_i^{(t)} = \{\mathbf{x}_{i,1}^{(t)}, \dots, \mathbf{x}_{i,K}^{(t)}\}, \quad (21)$$

using the multinomial probability mass function with probabilities

$$\bar{w}_{i,k}^{(t)} = \frac{w_{i,k}^{(t)}}{\sum_{\ell=1}^K w_{i,\ell}^{(t)}}, \quad k = 1, \dots, K. \quad (22)$$

where the unnormalized weights $w_{i,k}^{(t)}$ are given by Eq. (19). Note that again we can use any resampling technique, including the standard multinomial or other advanced schemes [31, 20]. In LR-PMC, there is no loss of diversity in the population of proposals, since each proposal at the current iteration yields another proposal in the next iteration. In other words, exactly one particle per proposal survives after the resampling step.

The adaptation scheme of LR-PMC can be intuitively understood as follows. Let us consider for a moment a modified version of LR-PMC where the weights used in the resampling are those of standard PMC of Eq. (5) instead of the DM weights of Eq. (22). This modified scheme is equivalent to N parallel PMC samplers, where the i -th PMC draws K samples from the i -th proposal, applying a resampling step independently from the other $N - 1$ PMC samplers. By using the DM weights in LR-PMC, we incorporate cooperation among the N proposals. When the proposal pdfs are close to each other, the local resampling scheme (with DM weights) adds a “repulsive” interaction: among the K samples of a specific proposal, the resampling promotes the samples in areas that are less covered by the other $N - 1$ proposals (and where, at the same time, the target evaluation is high). Therefore, this scheme performs a cooperative exploration of the state space by the N proposals. Note that, when the proposal pdfs are located far away from each other, the weights of the K samples of a specific proposal are in practice not

affected by other $N - 1$ proposals. In this case, the LR-PMC works as the N parallel PMC samplers described above.

Finally, let us remark that a mixed global-local resampling strategy (e.g., performing local resampling on clusters of proposals) could also be devised in order to obtain the advantages of both global and local resampling.

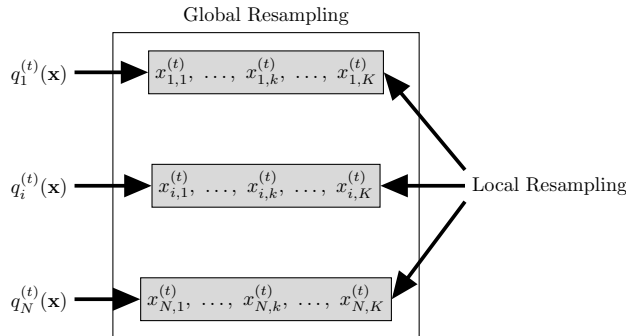


Figure 1: Sketch of the global and local resampling schemes considering N proposal pdfs at the t -th iteration, $q_i^{(t)}$ for $i = 1, \dots, N$ and $t = 1, \dots, T$, and K samples per proposal.

5. Numerical results

5.1. Estimation of the normalizing constant

Let us consider, as a target pdf, a bimodal mixture of Gaussians $\pi(\mathbf{x}) = \frac{1}{2}\mathcal{N}(x; \nu_1, c_1^2) + \frac{1}{2}\mathcal{N}(x; \nu_2, c_2^2)$ with $\nu_1 = -3$ and $\nu_2 = 3$, and $c_1^2 = 1$ and $c_2^2 = 1$. The proposal pdfs are also Gaussians: $q_1(x) = \mathcal{N}(x; \mu_1, \sigma_1)$ and $q_2(x) = \mathcal{N}(x; \mu_2, \sigma_2)$. At this point, we consider two scenarios:

- Scenario 1: In this case, $\mu_1 = \nu_1$, $\mu_2 = \nu_2$, $\sigma_1^2 = c_1^2$, and $\sigma_2^2 = c_2^2$. Then, both proposal pdfs can be seen as a whole mixture that exactly replicates the target, i.e., $\pi(\mathbf{x}) = \frac{1}{2}q_1(\mathbf{x}) + \frac{1}{2}q_2(\mathbf{x})$. This is the desired situation pursued by an adaptive importance sampling algorithm: each proposal is centered at a different mode of the target, and their scale parameters perfectly match the scales of the modes. Fig. 2(a) shows the target pdf in solid black line, and both proposal pdfs in blue and red dashed lines, respectively. Note that the proposals are scaled (each one integrates up to $1/2$ so we can see the perfect matching between the target and the mixture of proposal densities).
- Scenario 2: In this case, $\mu_1 = -2.5$, $\mu_2 = 2.5$, $\sigma_1^2 = 1.2$, and $\sigma_2^2 = 1.2$. Therefore, there is a mismatch between the target and the two proposals. Fig. 3(a) shows the target pdf in solid black line, and both proposal pdfs in blue and red dashed lines, respectively.

The goal is estimating the normalizing constant using the estimator \hat{Z} of Eq. (11) with $N = 2$ samples, one from each proposal, and $t = 1$. We use the standard PMC weights of Eq. (4) (estimator \hat{Z}_{IS}) and the DM-PMC weights of Eq. (14) (estimator \hat{Z}_{DM}).

	Estimator	\hat{Z}_{IS}	\hat{Z}_{DM}		Estimator	\hat{Z}_{IS}	\hat{Z}_{DM}
(Sc. 1)	Max.	35864	1	(Sc. 2)	Max.	77238	1.59
	$Var(\hat{Z})$	7891	0		$Var(\hat{Z})$	6874	0.01

Table 2: **(Ex. of Section. 5.1)** Maximum value of the estimator \hat{Z} in $2 \cdot 10^5$ runs for each scheme, in two different scenarios.

In order to characterize the two estimators, we run $2 \cdot 10^5$ simulations for each method. Note that the true value is $Z = 1$.

Figure 2(b) shows a boxplot of the distribution of the estimator \hat{Z} , obtained with both methods for Scenario 1. The blue lower and upper edges of the box correspond to the 25th and 75th percentiles, respectively, while the red line represents the median. The vertical black dashed whiskers extend to the minimum and maximum obtained values. Since the maxima cannot be appreciated in the figure, they are displayed in Table 2, altogether with the variance of the estimators. Note that even in this extremely simple and idealized scenario (perfect adaptation), the estimator obtained using the standard IS weights (i.e., the estimator used in standard PMC) has a poor performance. In most of the realizations, $\hat{Z}_{IS} \approx 0.5$ because each proposal (which integrates up to one) is adapted to one of the two modes (which contain roughly half of the probability mass).⁴ Since $E[\hat{Z}_{IS}] = Z = 1$, in a few runs the value the \hat{Z}_{IS} is extremely high as shown in Table 2. These huge values occur when a sample drawn from the tail of the proposal falls close to the other mode of the target (where actually the other proposal is placed). On the other hand, note that the DM estimator has a perfect performance (i.e., $\hat{Z}_{DM} = 1$ always, thus implying zero variance). Hence, this simple example shows that a substantial variance reduction can be attained by using the mixture at the denominator.

Figure 3(b) shows an equivalent boxplot for Scenario 2. In this case, the mismatch between proposals and target pdfs worsens both schemes. Note that the estimator \hat{Z}_{DM} now does not perfectly approximates Z , but still largely outperform the estimator \hat{Z}_{IS} . In particular, the median is still around the true value, and its variance is smaller.

5.2. Bi-dimensional example

We first consider a bivariate multimodal target pdf, consisting of a mixture of five Gaussians, i.e.,

$$\pi(\mathbf{x}) = \frac{1}{5} \sum_{i=1}^5 \mathcal{N}(\mathbf{x}; \boldsymbol{\nu}_i, \boldsymbol{\Sigma}_i), \quad \mathbf{x} \in \mathbb{R}^2, \quad (23)$$

⁴In this setup, each proposal approximately covers a different half of the target probability mass, since each one coincides with a different mode of the target. However, in standard PMC, the weight of each sample only accounts for its own proposal, and therefore there is not an exchange of information among the two proposals. Note that, if both proposals were covering the same mode (and therefore missing the other one), the weights would also be $w = 0.5$ in most of the runs; the lack of information exchange between the two samples, makes it impossible to know whether the target mass reported by the weight of each sample is the same and should be accounted “once”, or whether it is from another area and it should be accounted “twice”.

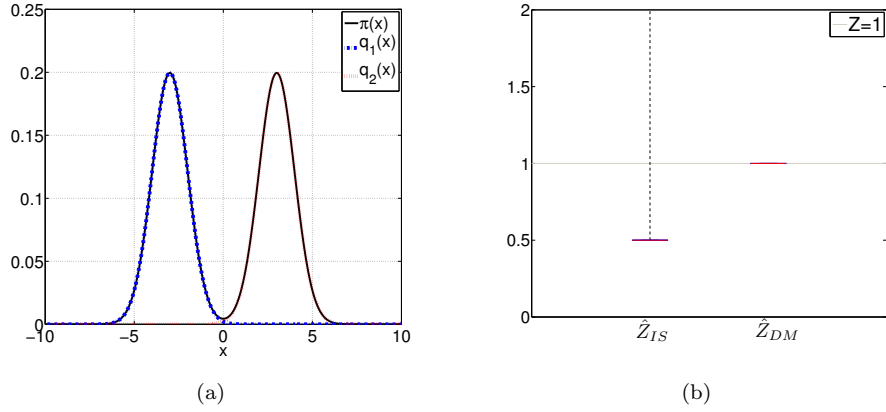


Figure 2: **(Ex. of Section 5.1)** Estimation of the normalizing constant (true value $Z = 1$) in Scenario 1 (perfect matching). **(a)** Target pdf (black solid line) and proposal pdfs (red and blue dashed lines). **(b)** Boxplot showing the 25th and 75th percentiles of the estimators \hat{Z}_{IS} and \hat{Z}_{DM} . The maximum value of \hat{Z}_{IS} is 35864.

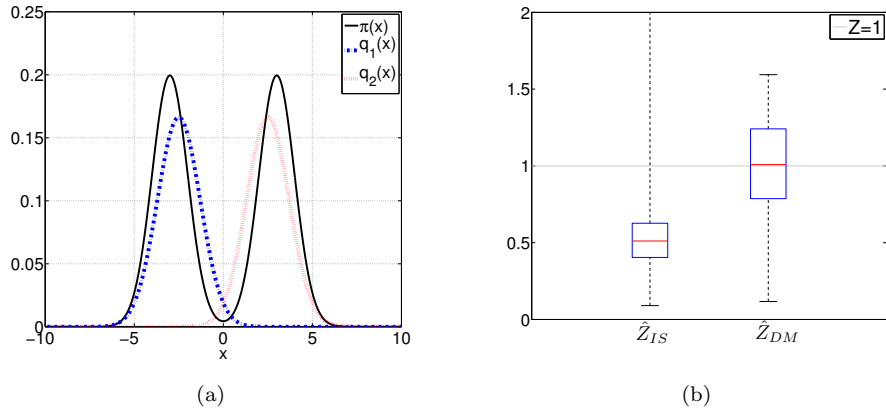


Figure 3: **(Ex. of Section 5.1)** Estimation of the normalizing constant (true value $Z = 1$) in Scenario 2 (proposal-target mismatch). **(a)** Target pdf (black solid line) and proposal pdfs (red and blue dashed lines). **(b)** Boxplot showing the distribution of the estimators \hat{Z}_{IS} and \hat{Z}_{DM} . The maximum value of \hat{Z}_{IS} is 77238.

where $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{C})$ denotes a normalized Gaussian pdf with mean vector $\boldsymbol{\mu}$ and covariance matrix \mathbf{C} , $\boldsymbol{\nu}_1 = [-10, -10]^\top$, $\boldsymbol{\nu}_2 = [0, 16]^\top$, $\boldsymbol{\nu}_3 = [13, 8]^\top$, $\boldsymbol{\nu}_4 = [-9, 7]^\top$, $\boldsymbol{\nu}_5 = [14, -14]^\top$, $\boldsymbol{\Sigma}_1 = [2, 0.6; 0.6, 1]$, $\boldsymbol{\Sigma}_2 = [2, -0.4; -0.4, 2]$, $\boldsymbol{\Sigma}_3 = [2, 0.8; 0.8, 2]$, $\boldsymbol{\Sigma}_4 = [3, 0; 0, 0.5]$, and $\boldsymbol{\Sigma}_5 = [2, -0.1; -0.1, 2]$. In this example, we can analytically compute different moments of the target in (23), and therefore we can easily validate the performance of the different techniques. In particular, we consider the computation of the mean of the target, $E[\mathbf{X}] = [1.6, 1.4]^\top$, and the normalizing constant, $Z = 1$, for $\mathbf{x} \sim \frac{1}{Z}\pi(\mathbf{x})$. We use as figure of merit the Mean Squared Error (MSE) in the estimation of $E[\mathbf{X}]$ (averaged over both components) and Z .

		$L = NKT = 2 \cdot 10^5$					
N	Algorithm	$\sigma = 1$	$\sigma = 2$	$\sigma = 5$	$\sigma = 10$	$\sigma = 20$	$\sigma = 70$
5	Standard PMC [3]	92.80 (85.23-99.57)	38.71 (31.96-47.69)	12.65 (7.10-19.04)	0.38 (0.28-0.53)	0.047 (0.033-0.065)	37.44 (21.01-55.62)
100		75.17 (72.72-78.20)	59.42 (54.78-64.23)	14.24 (12.04-16.57)	0.25 (0.21-0.30)	0.028 (0.023-0.033)	0.18 (0.15-0.22)
$5 \cdot 10^4$		68.29 (66.92-69.19)	37.44 (34.57-41.98)	7.01 (5.72-7.86)	0.25 (0.18-0.34)	0.033 (0.027-0.039)	0.17 (0.14-0.21)
	DM-PMC ($K = 1$)	72.48 (69.79-75.14)	36.21 (33.54-39.26)	5.34 (4.41-6.33)	0.036 (0.030-0.043)	0.029 (0.024-0.034)	0.21 (0.18-0.25)
	GR-PMC ($K = 2$)	69.41 (66.02-72.30)	26.23 (22.26-30.83)	3.09 (1.88-4.69)	0.022 (0.019-0.027)	0.028 (0.022-0.033)	0.17 (0.14-0.21)
	LR-PMC ($K = 2$)	2.68 (1.85-3.54)	0.007 (0.005-0.009)	0.010 (0.008-0.012)	0.018 (0.014-0.022)	0.102 (0.084-0.122)	32.88 (27.89-38.69)
	GR-PMC ($K = 5$)	67.04 (64.26-69.53)	17.44 (14.74-20.55)	0.11 (0.03-0.25)	0.013 (0.011-0.016)	0.023 (0.018-0.027)	0.15 (0.12-0.17)
	LR-PMC ($K = 5$)	8.04 (6.65-9.65)	0.012 (0.007-0.019)	0.008 (0.005-0.012)	0.016 (0.013-0.019)	0.027 (0.021-0.033)	2.00 (1.52-2.60)
	GR-PMC ($K = 20$)	61.58 (56.94-66.03)	15.13 (12.30-18.81)	0.42 (0.03-1.18)	0.012 (0.010-0.014)	0.024 (0.020-0.029)	0.14 (0.12-0.17)
	LR-PMC ($K = 20$)	9.51 (8.49-10.53)	1.16 (0.54-1.89)	0.011 (0.008-0.014)	0.013 (0.011-0.016)	0.023 (0.019-0.028)	0.22 (0.18-0.26)
	GR-PMC ($K = 100$)	64.94 (61.67-67.66)	12.50 (10.65-15.53)	0.08 (0.02-0.20)	0.015 (0.011-0.018)	0.026 (0.021-0.030)	0.18 (0.15-0.21)
	LR-PMC ($K = 100$)	9.60 (8.58-10.66)	1.21 (0.64-1.88)	0.022 (0.016-0.029)	0.015 (0.012-0.018)	0.026 (0.022-0.032)	0.20 (0.16-0.24)
	GR-PMC ($K = 500$)	58.49 (54.10-62.20)	9.63 (7.81-11.45)	0.08 (0.06-0.10)	0.014 (0.011-0.016)	0.024 (0.019-0.030)	0.16 (0.14-0.20)
	LR-PMC ($K = 500$)	14.79 (13.12-16.54)	6.72 (5.30-8.39)	0.10 (0.06-0.14)	0.010 (0.008-0.013)	0.024 (0.018-0.030)	0.20 (0.16-0.25)
	100	M-PMC [14]	71.39 (65.22-77.36)	81.33 (71.59-90.04)	18.14 (13.51-22.90)	0.058 (0.052-0.067)	0.031 (0.016-0.056)
10	SMC [18]	84.14 (73.46-97.81)	81.68 (67.66-95.91)	6.49 (2.58-10.45)	0.76 (0.15-1.71)	0.024 (0.021-0.027)	4.60 (1.64-8.51)
100		77.00 (76.35-77.66)	76.57 (75.60-77.66)	15.98 (15.42-16.59)	0.79 (0.64-0.97)	0.068 (0.065-0.072)	0.86 (0.79-0.93)
$5 \cdot 10^4$		69.08 (68.34-69.91)	51.29 (44.10-57.26)	20.48 (8.86-36.70)	0.22 (0.14-0.31)	0.038 (0.019-0.061)	0.68 (0.39-1.03)
100	DM-SMC ($K = 1$)	70.95 (70.16-71.74)	42.40 (41.49-43.39)	1.91 (1.72-2.15)	0.039 (0.037-0.040)	0.027 (0.026-0.029)	0.19 (0.18-0.19)
	GR-SMC ($K = 5$)	66.64 (65.42-67.84)	41.54 (39.93-43.01)	0.16 (0.15-0.18)	0.015 (0.014-0.016)	0.024 (0.023-0.025)	0.19 (0.19-0.20)
	LR-SMC ($K = 5$)	8.16 (7.68-8.66)	2.32 (1.92-2.71)	0.007 (0.006-0.008)	0.015 (0.014-0.016)	0.027 (0.026-0.028)	2.19 (2.08-2.29)
	GR-SMC ($K = 20$)	65.48 (64.16-66.67)	37.91 (36.21-39.75)	0.10 (0.05-0.18)	0.013 (0.012-0.014)	0.025 (0.024-0.026)	0.19 (0.18-0.20)
	LR-SMC ($K = 20$)	8.88 (8.45-9.32)	4.15 (3.65-4.62)	0.010 (0.008-0.012)	0.014 (0.013-0.014)	0.026 (0.025-0.027)	0.20 (0.19-0.20)

Table 3: (**Ex. of Section 5.2**) MSE in the estimation of $E[\mathbf{X}]$, for several values of σ and K , keeping the total number of evaluations of the target fixed to $L = KNT = 2 \cdot 10^5$ in all algorithms. The best results for each value of σ are highlighted in bold-face.

For simplicity, we assume Gaussian proposal densities for all of the methods compared, and deliberately choose a “bad” initialization of the means in order to test the robustness and the adaptation capabilities. Specifically, the initial adaptive parameters of the individual proposals are selected uniformly within the $[-4, 4] \times [-4, 4]$ square, i.e., $\boldsymbol{\mu}_i^{(1)} \sim \mathcal{U}([-4, 4] \times [-4, 4])$ for $i = 1, \dots, N$. This initialization is considered “bad”, since none of the modes of the target falls within the initialization square. We test all the alternatives using the same isotropic covariance matrices for all the Gaussian proposals, $\mathbf{C}_i = \sigma^2 \mathbf{I}_2$ with $\sigma \in \{1, 2, 5, 10, 20, 70\}$. All the results have been averaged over 500 independent experiments, where the computational cost of the different techniques (in terms of the total number of evaluations of the target distribution) is fixed to $L = KNT$.⁵ We compare the following schemes:

- **Standard PMC [3]:** Standard PMC algorithm described in Table 1 with $N = 100$ proposals and $T = 2000$ iterations. The total number of samples drawn is $L = NT = 2 \cdot 10^5$.
- **M-PMC [14]:** M-PMC algorithm proposed in [14] with $D = 100$ proposals, $N = 100$ samples per iteration, and $T = 2000$ iterations. The total number of samples drawn is $L = NT = 2 \cdot 10^5$.
- **SMC [18]:** We apply a Sequential Monte Carlo (SMC) scheme combining resampling and MCMC steps. Specifically, we consider Metropolis-Hastings (MH) steps as forward reversible kernels. In this example, we do not employ a sequence of tempered target pdfs, i.e., we consider always the true target density. The proposal pdfs for the MH kernels coincide with the Gaussian proposals employed in the propagation resampling steps, with the scale parameters \mathbf{C}_i of the other tested methods. Due to the application of the MH steps, in this case, $L > 2 \cdot 10^5$.
- **K-PMC:** Standard PMC scheme using $N = 100$ proposals, but drawing $K > 1$ samples per proposal at each iteration and performing global resampling (GR). In order to keep the total number of samples constant, the number of iterations of the algorithm is now $T = 2 \cdot 10^5 / (KN)$.
- **DM-PMC:** Standard PMC using the weights of Eq. (14), $N = 100$ proposals, $T = 2000$ iterations, and drawing $K = 1$ samples per proposal.
- **GR-PMC:** PMC scheme with multiple samples per mixand (K), weights computed as Eq. (19), and global resampling (GR). We use $N = 100$ proposals and $T = L / (KN)$ iterations with $L = 2 \cdot 10^5$ (as in the three previous schemes). In particular, we test the values $K \in \{2, 5, 20, 100, 500\}$, and thus $T \in \{1000, 400, 100, 20, 4\}$.
- **LR-PMC:** PMC scheme with multiple samples per mixand (K) and local resampling (LR). All the parameters are selected as in the GR-PMC scheme.
- **Improved SMC:** SMC scheme with the improvements proposed in this paper. In all cases, we use the weights of Eq. (14) (DM-SMC), and we try the GR-SMC and LR-SMC variants. We test $K \in \{5, 20\}$

⁵Note that $L = KNT$ also corresponds to the total number of samples generated in all the schemes.

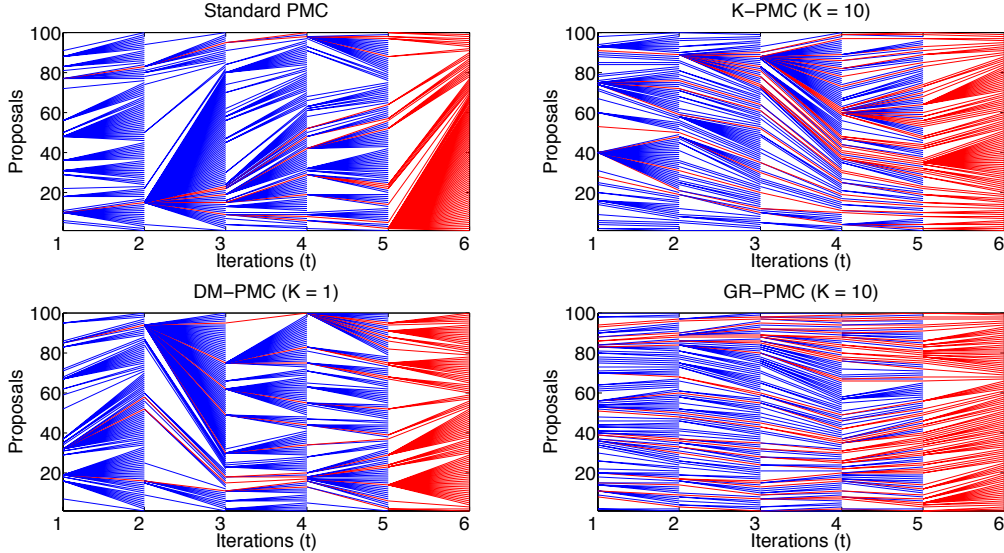


Figure 4: (**Ex. of Section 5.2**) Graphical representation of the indexes of the proposals used to generate the population for the next iteration with different schemes (6 iterations; $N = 100$, $\sigma = 5$). For each pair of iterations, lines link each surviving proposal (“father” proposal) with the next generation. In red, proposals surviving from the 1st to the 6th iteration.

Table 3 shows the MSE in the estimation of $E[\mathbf{X}]$ (averaged over both components) for $\mathbf{x} \sim \frac{1}{Z}\pi(\mathbf{x})$. We can see that all the proposed schemes outperform the standard PMC for any value of σ . In general, the local resampling (LR) works better than the global resampling (GR). Moreover, we note that the optimum value of K depends on the value of σ , the scale parameter of the proposals: for small values of σ (e.g., $\sigma = 1$ or $\sigma = 2$) small values of K lead to better performance, whereas a larger value of K (and thus less iterations T) can be used for larger values of σ (e.g., $\sigma = 10$ or $\sigma = 20$). In addition, the proposed methods also outperform the M-PMC algorithm in this scenario. Note that M-PMC is an adaptive importance sampling algorithm that does not perform the resampling step. Finally, note that the performance of the SMC sampler can be also improved with the proposed modifications.

The large MSE values in Table 3 for some schemes and sets of parameters are due to the fact that they fail at discovering all the modes of the target pdf. In order to clarify this issue, Fig. A.7 shows the evolution of the population of proposals for the first 4 iterations of the standard PMC ($K = 1$), K-PMC (with $K = 10$), and DM-PMC with global resampling (also for $K = 1$ and $K = 10$). Standard PMC tends to concentrate the whole population on one or two modes, very loosely covering the remaining ones and completely missing the mode in the bottom right corner. This issue is partly solved by using $K = 10$ (after 4 iterations the proposals are evenly distributed around 3 out of the 5 modes) or DM-PMC with $K = 1$ (after 4 iterations the proposals are uniformly distributed among 4 out of the 5 modes). Combining both approaches (DM-PMC and $K = 10$) an approximately uniform distribution of the proposals around all the modes

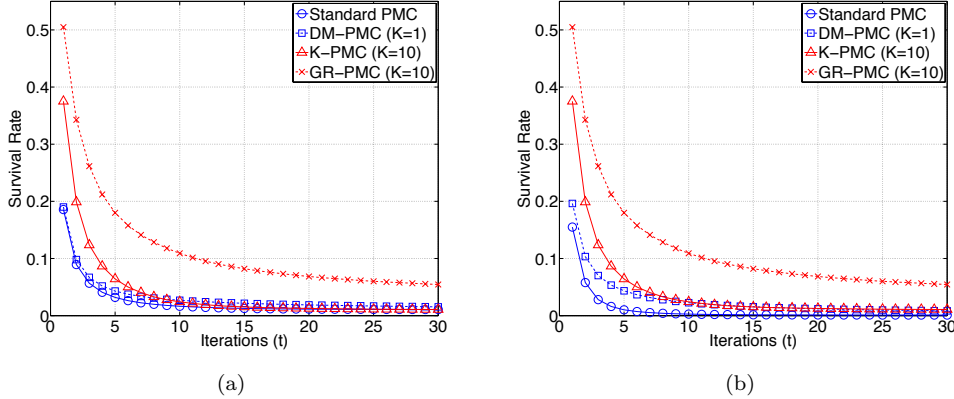


Figure 5: **(Ex. of Section 5.2)** Survival rate (after resampling) of the proposals vs. the distance in iterations among the proposals averaged over 500 runs ($\sigma = 5$). **(a)** $N = 100$ for all methods. **(b)** Different values of N , fixing $NK = 1000$ and thus $N \in \{100, 1000\}$.

of the target is attained.

Finally, in Figs. 4 and 5 we explore a well-known problem of PMC: the survival of proposals as the algorithm evolves. On the one hand, Fig. 4 shows which proposals have been used to generate the starting population for the next iteration. After 6 iterations, all of the $N = 100$ proposals in the population have arisen from only 2 of the proposals in the initial population. This situation hardly improves by using the DM-PMC: now 4 initial proposals have generated all the $N = 100$ proposals in the 6-th iteration. However, by drawing multiple samples per mixand ($K = 10$) the situation improves dramatically both when using the standard IS weights (9 proposals survive until the 6-th iteration) and especially when using the DM-PMC (19 surviving proposals). On the other hand, Fig. 5 shows the evolution in the survival rate of proposals w.r.t. the distance in iterations (or generations). In standard PMC, after very few iterations, most of the ancestors do not survive. This rate falls down as t increases in all cases, but the DM weights and especially the use of multiple samples per mixand help in slowing down this decrease. Therefore, we can conclude that the newly proposed schemes can be very useful in preserving the diversity in the population of proposals.

5.3. High-dimensional example

We consider a target corresponding to a mixture of isotropic Gaussians

$$\pi(\mathbf{x}) = \frac{1}{3} \sum_{k=1}^3 \mathcal{N}(\mathbf{x}; \boldsymbol{\nu}_k, \boldsymbol{\Sigma}_k), \quad (24)$$

where $\mathbf{x} \in \mathbb{R}^{10}$, $\boldsymbol{\nu}_k = [\nu_{k,1}, \dots, \nu_{k,10}]^\top$, and $\boldsymbol{\Sigma}_k = \xi_k^2 \mathbf{I}_{10}$ for $k \in \{1, 2, 3\}$, with \mathbf{I}_{10} being the 10×10 identity matrix. We set $\nu_{1,j} = -5$, $\nu_{2,j} = 6$, and $\nu_{3,j} = 3$ for all $j \in \{1, \dots, 10\}$. Moreover, we set $\xi_k = 8$ for all $k \in \{1, 2, 3\}$. The expected value of the target $\pi(\mathbf{x})$ is $E[X_j] = \frac{4}{3}$ for $j = 1, \dots, 10$, and the normalizing constant is $Z = 1$.

We use Gaussian proposal densities for all the compared methods. The initial means (adaptive parameters of the proposals) are selected randomly and independently in all techniques as $\mu_i^{(1)} \sim \mathcal{U}([-6 \times 6]^{10})$ for $i = 1, \dots, N$. We use the same isotropic covariance matrices for all the methods and proposal pdfs, $\mathbf{C}_i = \sigma^2 \mathbf{I}_{10}$, and we consider $\sigma \in \{1, 5, 20\}$. For every experiment, we run 200 independent simulations and compute the MSE in the estimation of $E[\mathbf{X}]$ (averaging the MSE of each component). We consider the same techniques as in the bi-dimensional example, testing $N \in \{100, 1000\}$ and different values of samples per iteration, $K \in \{2, 10, 20, 100\}$. We have tested different sets of parameters, always keeping the total number of samples fixed to $L = KNT = 2 \cdot 10^5$. Moreover, in this example we implement another variant of the SMC scheme [18], using a sequence of four tempered target densities, $\pi^{(1)}(\mathbf{x})$, $\pi^{(2)}(\mathbf{x})$, $\pi^{(3)}(\mathbf{x})$ and $\pi^{(4)}(\mathbf{x}) = \pi(\mathbf{x})$. These auxiliary targets have the same form as in Eq. (24), where the diagonal elements of each covariance matrix $\Sigma_k^{(s)}$, $s = 1, 2, 3, 4$ and $k = 1, 2, 3$, are respectively $\xi_k^{(1)} = 16$, $\xi_k^{(2)} = 12$, $\xi_k^{(3)} = 9$ and, finally, $\xi_k^{(4)} = 8$ (the true target). In addition, we also test this algorithm with the residual sampling (see for instance [31, 20]), instead of the standard multinomial resampling.

Table 4 shows that the proposed PMC schemes outperform the standard PMC in most of the cases. Indeed, a decrease of more than one order of magnitude in the MSE can often be attained by using DM-PMC with an appropriate value of K instead of the standard PMC. Finally, note that, although M-PMC behaves well for most of the parameters tested, overall the proposed methods yield the best performance in terms of MSE and robustness w.r.t. parameter choice.

In order to study the performance of the proposed schemes as the dimension of the state space increases, we change the dimension of the state space in (24). Namely, the target density is still a mixture of three isotropic Gaussians with the same structure for the mean vectors and covariance matrices as before, but now the dimension of \mathbf{x} is $D_x \in [1, 50]$. We have tested all the methods with $\sigma = 5$ and $N = 100$. Fig. 6 shows the evolution of the MSE in the estimation of the normalizing constant as a function of D_x . As expected, the performance of all the methods degrades as the dimension of the problem, D_x , becomes larger. Nonetheless, the performance of the proposed methods decays much more slowly than that of the standard PMC, thus allowing them to still provide a reasonably low MSE in higher dimensions. Note that, since the true normalizing constant of the target is $Z = 1$, when the methods behave poorly in high dimensions and the proposals do not discover the modes, the estimation is $\hat{Z} \approx 0$, and therefore the MSE tends to 1, which is the worst-case situation.

Algorithm	$N = 100$			$N = 1000$		
	$\sigma = 1$	$\sigma = 5$	$\sigma = 20$	$\sigma = 1$	$\sigma = 5$	$\sigma = 20$
Standard PMC [3]	12.43 (10.85-14.19)	8.11 (6.47-9.71)	1.24 (0.94-1.61)	12.68 (9.78-16.14)	5.94 (3.14-10.48)	0.53 (0.32-0.85)
GR-PMC ($K = 2$)	14.53 (13.29-16.07)	4.05 (2.52-6.24)	0.50 (0.43-0.58)	11.90 (7.86-17.65)	0.01 (0.01-0.02)	0.15 (0.12-0.20)
LR-PMC ($K = 2$)	11.55 (9.11-14.29)	12.77 (9.21-15.36)	78.31 (67.67-86.79)	2.52 (1.69-3.39)	0.82 (0.50-1.27)	29.44 (20.52-37.92)
GR-PMC ($K = 10$)	13.02 (11.69-14.48)	0.91 (0.48-1.58)	0.22 (0.20-0.24)	3.57 (1.82-6.37)	0.10 (0.00-0.27)	0.19 (0.14-0.25)
LR-PMC ($K = 10$)	8.15 (6.44-10.81)	0.21 (0.13-0.30)	1.85 (1.56-2.12)	4.34 (2.62-6.86)	0.01 (0.00-0.01)	1.61 (1.06-2.12)
GR-PMC ($K = 20$)	10.89 (9.82-11.92)	0.74 (0.35-1.32)	0.23 (0.20-0.26)	5.45 (2.49-9.43)	0.05 (0.02-0.09)	0.12 (0.08-0.16)
LR-PMC ($K = 20$)	6.92 (5.56-8.35)	0.16 (0.11-0.25)	0.77 (0.68-0.87)	4.59 (2.15-8.21)	0.04 (0.02-0.08)	0.55 (0.42-0.65)
GR-PMC ($K = 100$)	7.61 (6.57-8.60)	0.16 (0.10-0.29)	0.17 (0.15-0.18)	5.71 (3.28-9.78)	0.65 (0.15-1.46)	0.10 (0.07-0.14)
LR-PMC ($K = 100$)	7.05 (4.99-8.73)	0.41 (0.09-0.99)	0.28 (0.24-0.33)	5.48 (3.01-9.12)	0.17 (0.10-0.28)	0.19 (0.15-0.23)
M-PMC [14]	10.78 (9.53-19.78)	9.06 (4.40-12.72)	0.35 (0.20-0.64)	3.28 (2.77-4.88)	0.12 (0.07-0.50)	0.07 (0.05-0.12)
SMC [18]	4.99 (3.40-6.87)	0.92 (0.67-1.12)	0.45 (0.35-0.58)	11.45 (7.39-15.67)	1.75 (1.20-2.50)	0.38 (0.25-0.54)
SMC with tempering [18]	3.80 (2.76-4.90)	0.56 (0.48-0.65)	0.41 (0.29-0.50)	7.04 (4.75-9.82)	1.64 (1.12-2.03)	0.51 (0.41-0.67)
SMC with tempering and residual resampling [18]	2.79 (2.53-3.15)	0.54 (0.50-0.57)	0.26 (0.24-0.27)	7.29 (6.73-7.83)	1.24 (1.16-1.35)	0.43 (0.40-0.47)

Table 4: (**Ex. of Section 5.3**) MSE in the estimation of $E[\mathbf{X}]$, for $\sigma \in \{1, 5, 20\}$ and $K \in \{2, 10, 20, 100\}$, keeping the total number of evaluations of the target fixed to $L = 2 \cdot 10^5$. The dimension space of the target is $D_x = 10$. The best results for each value of σ are highlighted in bold-face.

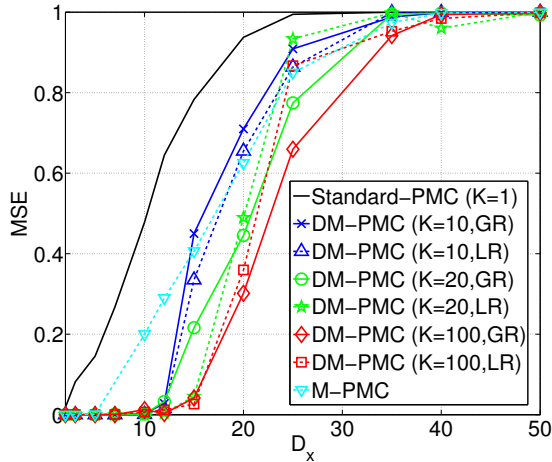


Figure 6: (**Ex. of Section 5.3**) MSE of the normalizing constant Z , using $N = 100$ proposals and a scale parameter $\sigma = 5$, as the dimension of the state space D_x increases.

5.4. Autoregressive filter with non-Gaussian noise

We consider the use of an autoregressive (AR) model contaminated by a non-Gaussian noise. This kind of filters is often used for modeling some financial time series (see for instance [34, Section 5] and [35]), where the noise is assumed to follow the so-called *generalized hyperbolic distribution* [36]. Namely, we consider the following observation model,

$$y_m = x_1 y_{m-1} + x_2 y_{m-2} + x_3 y_{m-3} + x_4 y_{m-4} + u_m, \quad (25)$$

where $m = 1, \dots, M$ is a time index, and u_m is a heavy-tailed driving noise:

$$u_m \sim p(u) \propto e^{\beta(u-\mu)} \frac{B_{\lambda-\frac{1}{2}}\left(\alpha\sqrt{\delta^2 + (u-\mu)^2}\right)}{\left(\sqrt{\delta^2 + (u-\mu)^2}\right)^{\frac{1}{2}-\lambda}},$$

where B_λ denotes the *modified Bessel function* [37]. The vector of unknowns, $\mathbf{x}^* = [x_1^*, x_2^*, x_3^*, x_4^*]^\top$, contains the coefficients of the AR model.

Given a set of observations $\mathbf{y} = [y_1, \dots, y_M]^\top$, the inference problem consists of obtaining statistical information about \mathbf{x}^* , by studying the corresponding posterior distribution $\tilde{\pi}(\mathbf{x}|\mathbf{y})$. More specifically, we have synthetically generated $M = 200$ observations, $\mathbf{y} = [y_1, \dots, y_M]^\top$, setting $\mathbf{x}^* = [0.5, 0.1, -0.8, 0.1]^\top$, $\lambda = 0.5$, $\alpha = 2$, $\beta = 1$, $\mu = -1$, and $\delta = 1$.⁶ Assuming improper uniform priors over the unknown coefficients, the objective is computing the expected value $\hat{\mathbf{x}} = \int_{\mathbb{R}^4} \mathbf{x} \tilde{\pi}(\mathbf{x}|\mathbf{y}) d\mathbf{x}$. Since we are using $M = 200$

⁶For the generation of i.i.d. samples of the generalized hyperbolic noise, we applied a fast and efficient MCMC technique (the FUS algorithm [38]), drawing samples from univariate distributions. After a few iterations, the resulting samples were virtually independent.

Algorithm	$N = 100$	$N = 1000$	$N = 5000$
Standard PMC [3]	13615.95 (13197.39-15021.45)	69.99 (62.23-72.24)	0.56 (0.56-0.68)
GR-PMC ($K = 5$)	1597.57 (1516.26-1727.82)	1.92 (1.72-2.22)	0.08 (0.07-0.10)
LR-PMC ($K = 5$)	31.04 (28.99-33.33)	0.36 (0.33-0.40)	0.20 (0.15-0.24)
GR-PMC ($K = 10$)	520.62 (472.44-558.27)	0.30 (0.26-0.40)	0.07 (0.05-0.10)
LR-PMC ($K = 10$)	14.99 (14.06-15.99)	0.29 (0.28-0.32)	0.21 (0.14-0.27)
GR-PMC ($K = 50$)	16.91 (15.43-20.42)	0.05 (0.04-0.08)	*
LR-PMC ($K = 50$)	1.89(1.61-2.12)	0.15 (0.14-0.21)	*
GR-PMC ($K = 100$)	2.23 (1.74-3.39)	0.10 (0.06-0.13)	*
LR-PMC ($K = 100$)	0.77 (0.62-0.90)	0.17 (0.09-0.18)	*
M-PMC [14]	182.10 (64.29-316.59)	0.07 (0.06-0.09)	0.05 (0.04-0.07)

Table 5: (**Ex. of Section 5.4**) MSE of $E[\mathbf{X}]$ for different values of K and N , keeping the total number of evaluations of the target fixed to $L = KNT = 2 \cdot 10^5$. The symbol * indicates combinations where the number of iterations $T < 1$, and therefore they cannot be performed.

observations (a large number for this example), we assume that the posterior pdf is quite sharp and concentrated around the true value, $\mathbf{x}^* = [0.5, 0.1, -0.8, 0.1]^\top$. Nevertheless, in practice we assume that the inference algorithms have no clue of which is that true value (i.e., we assume no a priori information). Therefore, \mathbf{x}^* is only used for evaluating the performance of the different methods in terms of MSE.

All the methods use Gaussian proposals, with the initial adaptive parameters of the individual proposals selected uniformly within the $[-6, 6]^4$ square, i.e., $\boldsymbol{\mu}_i^{(1)} \sim \mathcal{U}([-6, 6] \times [-6, 6] \times [-6, 6] \times [-6, 6])$, and the covariance matrices for all the Gaussians selected as $\mathbf{C}_i = \sigma^2 \mathbf{I}_4$, with $\sigma = 5$ for $i = 1, \dots, N$. As in the previous examples, we have tested different combinations of parameters, keeping the total number of evaluations of the target fixed to $L = NKT = 2 \cdot 10^5$. We have evaluated different values of $N \in \{100, 1000, 5000\}$ and $K \in \{5, 10, 50, 100\}$. We ran 500 independent simulations and computed the MSE in the estimation of $\hat{\mathbf{x}}$ w.r.t. the true value \mathbf{x}^* .

The results obtained by the different methods, in terms of MSE averaged over all the components of \mathbf{x} , are shown in Table 5. Note that some combinations of K and N would yield a number of iterations $T < 1$, since we set $T = L/(NK) = 2 \cdot 10^5/(NK)$. Therefore, those simulations cannot be performed and are indicated in the Table with the symbol *. Note that, for any choice of N , the alternative schemes proposed in the paper largely outperform the standard PMC. Furthermore, the advantage of using $K > 1$ can again be clearly seen for the three values of N tested. More specifically, the smallest the value of N the largest the value of K that should be used to attain the best results. Note also that M-PMC behaves particularly well in this scenario for high values of N , but its performance is very poor for $N = 100$ (unlike GR-PMC and LR-PMC, which can still provide a good performance for the right value of K).

5.5. Localization problem in a wireless sensor network

Let us consider a static target in a two-dimensional space. The goal consists on positioning the target within a wireless sensor network using only range measurements acquired by some sensors. This example appears in the signal processing literature for localization applications, e.g. in [39, 40, 41, 17]. In particular, let $\mathbf{X} = [X_1, X_2]^\top$ denote the random vector representing the position of the target in \mathbb{R}^2 plane. The

Algorithm	N= 100		N= 500	
	$\sigma = 1$	$\sigma = 2$	$\sigma = 1$	$\sigma = 2$
Standard PMC [3]	621.85 (542.98-685.76)	2424.35 (1916.39-2995.05)	167.52 (33.10-376.49)	756.46 (490.36-1077.76)
GR-PMC ($K = 20$)	7.51 (5.83-8.97)	28.02 (23.15-33.41)	0.87 (0.11-1.86)	9.30 (3.33-14.30)
LR-PMC ($K = 20$)	0.59 (0.50-0.68)	1.27 (0.89-1.65)	0.25 (0.10-0.54)	0.40 (0.35-0.45)
GR-PMC ($K = 50$)	1.82 (1.50-2.26)	7.00 (5.30-8.56)	0.52 (0.39-0.69)	1.72 (0.19-3.56)
LR-PMC ($K = 50$)	0.37 (0.32-0.44)	0.88 (0.70-1.04)	0.25 (0.13-0.32)	0.38 (0.30-0.47)
GR-PMC ($K = 100$)	0.74 (0.63-0.88)	1.66 (1.31-2.05)	0.32 (0.17-0.48)	0.80 (0.51-0.99)
LR-PMC ($K = 100$)	0.28 (0.25-0.33)	0.48 (0.39-0.58)	0.23 (0.14-0.33)	0.11 (0.11-0.11)
GR-PMC ($K = 200$)	0.43 (0.36-0.51)	0.57 (0.45-0.66)	0.36 (0.22-0.48)	0.23 (0.01-0.35)
LR-PMC ($K = 200$)	0.26 (0.23-0.29)	0.35 (0.28-0.42)	0.16 (0.12-0.20)	0.37 (0.37-0.37)
M-PMC [14]	7.75 (6.76-7.73)	32.77 (28.34-37.19)	1.07 (0.82-1.33)	1.66 (1.48-1.84)

Table 6: (Ex. of Section 5.5) MSE of the estimator of $E[\mathbf{X}]$ with different PMC algorithms.

measurements are obtained from 6 range sensors located at $\mathbf{h}_1 = [1, -8]^\top$, $\mathbf{h}_2 = [8, 10]^\top$, $\mathbf{h}_3 = [-15, -7]^\top$, $\mathbf{h}_4 = [-8, 1]^\top$, $\mathbf{h}_5 = [10, 0]^\top$ and $\mathbf{h}_6 = [0, 10]^\top$. The measurements are related to the target position through the following expression:

$$Y_{j,r} = -20 \log(\|\mathbf{x} - \mathbf{h}_j\|^2) + \Theta_j, \quad j = 1, \dots, 6, \quad r = 1, \dots, d_y, \quad (26)$$

where $\Theta_j \sim \mathcal{N}(\theta_j | \mathbf{0}, \omega_j^2 \mathbf{I})$, with $\omega_j = 5$ for all $j \in 1, \dots, 6$. Note that the total number of data is $6d_y$. We consider a wide Gaussian prior pdf with mean $[0, 0]^\top$ and covariance matrix $[\omega_0^2 \ 0; 0 \ \omega_0^2]^\top$ with $\omega_0 = 10$,

We simulate $6d_y = 360$ measurements from the model ($d_y = 60$ observations from each sensor), fixing $x_1 = 3.5$ and $x_2 = 3.5$. The goal consists in approximating the mean of the posterior distribution $\bar{\pi}(\mathbf{x}|\mathbf{y})$, through the improved PMC techniques proposed in this paper. In order to compare the different techniques, we computed the value of interest by using an extremely thin grid, yielding $E[\mathbf{X}] \approx [3.415, 3.539]^\top$.

We test the proposed methods and we compare them with the standard PMC [3] and the M-PMC [14]. In all cases, Gaussian proposals are used, with initial mean parameters selected uniformly within the $[1, 5] \times [1, 5]$ square, i.e., $\boldsymbol{\mu}_i^{(1)} \sim \mathcal{U}([-1, 5] \times [-4, 4])$ for $i = 1, \dots, N$. All the methods use the same isotropic covariance matrices for all the Gaussian proposals, $\mathbf{C}_i = \sigma^2 \mathbf{I}_2$ with $\sigma \in \{1, 2\}$. We have tried $N \in \{100, 500\}$ proposals. In the proposed methods, we test the values $K \in \{20, 50, 100, 200\}$. Note that again, we keep fixed the total number of evaluations to $L = KNT = 2 \cdot 10^5$.

Table 6 shows the MSE in estimation of the expected value of the posterior, with the different PMC methods. Again, the proposed methods largely beat the standard PMC for all the sets of parameters. The M-PMC algorithm is again competitive (especially with $N = 500$), but the proposed algorithms obtain better performance (in particular, the LR-PMC with a high K).

6. Conclusions

The population Monte Carlo (PMC) method is a well-known and widely used scheme for performing statistical inference in many signal processing problems. Three improved PMC algorithms are proposed in this paper. All of them are based on the deterministic

mixture (DM) approach, which provides estimators with a reduced variance (as proved in this paper) and increases the exploratory behavior of the resulting algorithms. Additionally, two of the methods draw multiple samples per mixand (both with local and global resampling strategies) to prevent the loss of diversity in the population of proposals. The proposed approaches are shown to substantially outperform the standard PMC on three numerical examples. The proposed improvements can be applied to other existing PMC implementations and other importance sampling techniques, to achieve similar benefits.

Acknowledgment

This work has been supported by the Spanish government’s projects AGES (S2010/BMD-2422), ALCIT (TEC2012-38800-C03-01), COMPREHENSION (TEC2012-38883-C02-01), DISSECT (TEC2012-38058-C03-01), and OTOSiS (TEC2013-41718-R); by the BBVA Foundation through project MG-FIAR (“I Convocatoria de Ayudas Fundación BBVA a Investigadores, Innovadores y Creadores Culturales”); by ERC grant 239784 and AoF grant 251170; by the National Science Foundation under Award CCF-0953316; and by the European Union’s 7th FP through the Marie Curie ITN MLPM2012 (Grant No. 316861).

References

References

- [1] C. P. Robert, G. Casella, *Monte Carlo Statistical Methods*, Springer, 2004.
- [2] Y. Iba, Population Monte Carlo algorithms, *Transactions of the Japanese Society for Artificial Intelligence* 16 (2001) 279–286.
- [3] O. Cappé, A. Guillin, J. M. Marin, C. P. Robert, Population Monte Carlo, *Journal of Computational and Graphical Statistics* 13 (4) (2004) 907–929.
- [4] G. Celeux, J. M. Marin, C. P. Robert, Iterated importance sampling in missing data problems, *Computational Statistics & Data Analysis* 50 (2006) 3386–3404.
- [5] M. C. A. M. Bink, M. P. Boer, C. J. F. Braak, J. Jansen, R. E. Voorrips, W. E. van de Weg, Bayesian analysis of complex traits in pedigreed plant populations, *Euphytica* 161 (2008) 85–96.
- [6] C. Bi, A Monte Carlo EM algorithm for de novo motif discovery in biomolecular sequences, *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 6 (2009) 370–386.
- [7] G. E. Barter, L. K. Purvis, N. P. Tecler, T. H. West, Analysis of detection systems for outdoor chemical or biological attacks, in: *Proc. IEEE Conf. on Technologies for Homeland Security*, 2009.
- [8] H. Bhaskar, L. Mihaylova, S. Maskell, Population based particle filtering, in: *IET Seminar on Target Tracking and Data Fusion: Algorithms and Applications*, 2008.
- [9] C. Andrieu, J. Thoms, A tutorial on adaptive MCMC, *Statistics and Computing* 18 (4) (2008) 343–373.
- [10] T. Li, M. Bolic, P. M. Djuric, Resampling methods for particle filtering: Classification, implementation, and strategies, *IEEE Signal Processing Magazine* 32 (3) (2015) 70–86.
- [11] R. Douc, A. Guillin, J. M. Marin, C. P. Robert, Convergence of adaptive mixtures of importance sampling schemes, *Annals of Statistics* 35 (2007) 420–448.
- [12] R. Douc, A. Guillin, J. M. Marin, C. P. Robert, Minimum variance importance sampling via population Monte Carlo, *ESAIM: Probability and Statistics* 11 (2007) 427–447.
- [13] A. Iacobucci, J.-M. Marin, C. Robert, On variance stabilisation in population Monte Carlo by double Rao-Blackwellisation, *Computational Statistics & Data Analysis* 54 (3) (2010) 698–710.
- [14] O. Cappé, R. Douc, A. Guillin, J. M. Marin, C. P. Robert, Adaptive importance sampling in general mixture classes, *Statistics and Computing* 18 (2008) 447–459.
- [15] E. Koblents, J. Míguez, A population Monte Carlo scheme with transformed weights and its application to stochastic kinetic models, *Statistics and Computing* (2013) 1–19.

- [16] J. M. Cornuet, J. M. Marin, A. Mira, C. P. Robert, Adaptive multiple importance sampling, *Scandinavian Journal of Statistics* 39 (4) (2012) 798–812.
- [17] L. Martino, V. Elvira, D. Luengo, J. Corander, An adaptive population importance sampler: Learning from the uncertainty, *IEEE Transactions on Signal Processing* 63 (16) (2015) 4422–4437.
- [18] P. D. Moral, A. Doucet, A. Jasra, Sequential Monte Carlo samplers, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 68 (3) (2006) 411–436.
- [19] J. S. Liu, *Monte Carlo Strategies in Scientific Computing*, Springer, 2004.
- [20] R. Douc, O. Cappé, E. Moulines, Comparison of resampling schemes for particle filtering, in: *Proc. 4th Int. Symp. on Image and Signal Processing and Analysis*, 2005, pp. 64–69.
- [21] P. Del Moral, A. Doucet, A. Jasra, On adaptive resampling strategies for sequential Monte Carlo methods, *Bernoulli* 18 (1) (2012) 252–278.
- [22] J. Geweke, Bayesian inference in econometric models using monte carlo integration, *Econometrica: Journal of the Econometric Society* (1989) 1317–1339.
- [23] D. W. Scott, *Multivariate density estimation: theory, practice, and visualization*, John Wiley & Sons, 2009.
- [24] M. Wand, M. Jones, *Kernel smoothing*, Chapman and Hall, 1994.
- [25] E. Veach, L. Guibas, Optimally combining sampling techniques for Monte Carlo rendering, In *SIGGRAPH 1995 Proceedings* (1995) 419–428.
- [26] A. Owen, Y. Zhou, Safe and effective importance sampling, *Journal of the American Statistical Association* 95 (449) (2000) 135–143.
- [27] G. H. Hardy, J. E. Littlewood, G. Pólya, *Inequalities*, Cambridge Univ. Press, 1952.
- [28] A. Kong, A note on importance sampling using standardized weights, University of Chicago, Dept. of Statistics, Tech. Rep 348.
- [29] V. Elvira, L. Martino, D. Luengo, M. F. Bugallo, Generalized multiple importance sampling, arXiv preprint arXiv:1511.03095.
- [30] V. Elvira, L. Martino, D. Luengo, M. F. Bugallo, Efficient multiple importance sampling estimators, *Signal Processing Letters, IEEE* 22 (10) (2015) 1757–1761.
- [31] O. Cappé, T. Rydeen, E. Moulines, *Inference in Hidden Markov Models*, Springer, 2005.
- [32] M. S. Arulampalam, S. Maskell, N. Gordon, T. Clapp, A tutorial on particle filters for online nonlinear/non-gaussian bayesian tracking, *Signal Processing, IEEE Transactions on* 50 (2) (2002) 174–188.
- [33] T. Li, S. Sun, T. P. Sattar, J. Corchado, Fight sample degeneracy and impoverishment in particle filters: A review of intelligent approaches, *Expert Systems with applications* 41 (8) (2014) 3944–3954.
- [34] R. H. Shumway, D. S. Stoffer, *Time series analysis and its applications*, Springer Science & Business Media, 2013.
- [35] J. D. Hamilton, R. Susmel, Autoregressive conditional heteroskedasticity and changes in regime, *Journal of Econometrics* 64 (1) (1994) 307–333.
- [36] E. Eberlein, Application of generalized hyperbolic lévy motions to finance, in: *Lévy processes*, Springer, 2001, pp. 319–336.
- [37] M. Abramowitz, I. A. Stegun, *Handbook of mathematical functions: with formulas, graphs, and mathematical tables*, no. 55, Dover Pub., 1972.
- [38] L. Martino, H. Yang, D. Luengo, J. Kannianen, J. Corander, A fast universal self-tuned sampler within Gibbs sampling, *Digital Signal Processing*, In press, 2015.
- [39] A. M. Ali, K. Yao, T. C. Collier, E. Taylor, D. Blumstein, L. Girod, An empirical study of collaborative acoustic source localization, *Proc. Information Processing in Sensor Networks (IPSN07)*, Boston.
- [40] A. T. Ihler, J. W. Fisher, R. L. Moses, A. S. Willsky, Nonparametric belief propagation for self-localization of sensor networks, *IEEE Transactions on Selected Areas in Communications* 23 (4) (2005) 809–819.
- [41] L. Martino, J. Míguez, Generalized rejection sampling schemes and applications in signal processing, *Signal Processing* 90 (11) (2010) 2981–2995.
- [42] J. Geweke, Bayesian inference in econometric models using Monte Carlo integration, *Econometrica* 24 (1989) 1317–1399.

Appendix A. Standard vs. deterministic mixture importance sampling

In this appendix we review the IS estimators, analyzing the properties (unbiasedness and variance) of the estimator with the DM weights. For the sake of clarity, we remove the temporal indexes.

Appendix A.1. Importance sampling estimators

Let us consider the estimator of Eq. (9) when we have a set of N proposal pdfs, $\{q_i(\mathbf{x})\}_{i=1}^N$. We draw exactly $K_i = 1$ sample from each proposal, i.e., $\mathbf{x}_i \sim q_i(\mathbf{x})$ for $i = 1, \dots, N$.⁷ If the normalizing constant Z is known, the IS estimator is then

$$\hat{I} = \frac{1}{NZ} \sum_{i=1}^N w_i f(\mathbf{x}_i). \quad (\text{A.1})$$

The difference between the standard and deterministic mixture (DM) IS estimators lies in the computation of the unnormalized weights. On the one hand, we recall the standard IS weights are given by

$$w_i = \frac{\pi(\mathbf{x}_i)}{q_i(\mathbf{x}_i)}, \quad (\text{A.2})$$

where $\pi(\mathbf{x}_i)$ is the target evaluated at the i -th sample (drawn from the i -th proposal). Substituting (A.2) into (A.1), we obtain the standard IS estimator,

$$\hat{I}_{IS} = \frac{1}{NZ} \sum_{i=1}^N \frac{f(\mathbf{x}_i)\pi(\mathbf{x}_i)}{q_i(\mathbf{x}_i)}. \quad (\text{A.3})$$

On the other hand, the weights in the DM approach are given by

$$w_i = \frac{\pi(\mathbf{x}_i)}{\sum_{j=1}^N q_j(\mathbf{x}_i)}. \quad (\text{A.4})$$

Substituting (A.4) into (A.1) we obtain the DM estimator

$$\hat{I}_{DM} = \frac{1}{NZ} \sum_{i=1}^N \frac{f(\mathbf{x}_i)\pi(\mathbf{x}_i)}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x}_i)}. \quad (\text{A.5})$$

Appendix A.2. Unbiasedness of the DM-IS estimator

It is well known that \hat{I}_{IS} in Eq. (A.3) is an unbiased estimator of the integral I define in Eq. (2) [1, 19]. In this section, we prove that the DM-IS estimator in Eq. (A.5) is

⁷From now on, we use $K_i = 1$, with $i = 1, \dots, N$, for the sake of clarity, but the analysis can be straightforwardly extended to any K_i .

also unbiased. Since $\mathbf{x}_i \sim q_i(\mathbf{x})$, we have

$$E[\hat{I}_{DM}] = \frac{1}{NZ} \sum_{i=1}^N E_{q_i} \left[\frac{f(\mathbf{X}_i)\pi(\mathbf{X}_i)}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{X}_i)} \right] \quad (\text{A.6})$$

$$= \frac{1}{NZ} \sum_{i=1}^N \int \frac{f(\mathbf{x}_i)\pi(\mathbf{x}_i)}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x}_i)} q_i(\mathbf{x}_i) d\mathbf{x}_i \quad (\text{A.7})$$

$$= \frac{1}{Z} \int \frac{f(\mathbf{x})\pi(\mathbf{x})}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x})} \left[\frac{1}{N} \sum_{i=1}^N q_i(\mathbf{x}) \right] d\mathbf{x} \quad (\text{A.8})$$

$$= \frac{1}{Z} \int f(\mathbf{x})\pi(\mathbf{x}) d\mathbf{x} = I. \quad (\text{A.9})$$

□

Appendix A.3. Variance of the DM-IS estimator

In this section, we prove that the DM-IS estimator in Eq. (A.5) always has a lower or equal variance than the standard IS estimator of Eq. (A.3). We also consider the standard mixture (SM) estimator \hat{I}_{SM} , where N samples are independently drawn from the mixture of proposals, i.e., $\mathbf{z}_i \sim \frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x})$, and

$$\hat{I}_{SM} = \frac{1}{NZ} \sum_{i=1}^N \frac{f(\mathbf{z}_i)\pi(\mathbf{z}_i)}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{z}_i)}. \quad (\text{A.10})$$

Note that obtaining an IS estimator with finite variance essentially amounts to having a proposal with heavier tails than the target. See [1, 42] for sufficient conditions that guarantee this finite variance.

Theorem 1. *For any target distribution, $\pi(\mathbf{x})$, any square integrable function w.r.t. $\pi(\mathbf{x})$, $f(\mathbf{x})$, and any set of proposal densities, $\{q_i(\mathbf{x})\}_{i=1}^N$, such that the variance of the corresponding estimators is finite, the variance of the DM estimator is always lower or equal than the variance of the corresponding standard IS and mixture (SM) estimators, i.e.,*

$$\text{Var}(\hat{I}_{DM}) \leq \text{Var}(\hat{I}_{SM}) \leq \text{Var}(\hat{I}_{IS}). \quad (\text{A.11})$$

Proof: The proof is given by Proposition 2 and Proposition 3. □

Proposition 2.

$$\text{Var}(\hat{I}_{SM}) \leq \text{Var}(\hat{I}_{IS}). \quad (\text{A.12})$$

Proof: The variance of the IS estimator is given by

$$\text{Var}(\hat{I}_{IS}) = \sum_{i=1}^N \frac{1}{N^2 Z^2} \int \frac{f^2(\mathbf{x})\pi^2(\mathbf{x})}{q_i(\mathbf{x})} d\mathbf{x} - \frac{I^2}{N}, \quad (\text{A.13})$$

where $I = \frac{1}{Z} \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$ is the true value of the integral that we want to estimate [29]. The variance of the SM estimator is given by

$$\begin{aligned} \text{Var}(\hat{I}_{SM}) &= \frac{1}{N^2} \sum_{i=1}^N \left(\frac{1}{Z^2} \int \frac{f^2(\mathbf{x})\pi^2(\mathbf{x})}{\psi(\mathbf{x})} d\mathbf{x} - I^2 \right) \\ &= \frac{1}{NZ^2} \int \frac{f^2(\mathbf{x})\pi^2(\mathbf{x})}{\psi(\mathbf{x})} d\mathbf{x} - \frac{I^2}{N}, \end{aligned} \quad (\text{A.14})$$

where $\psi(\mathbf{x}) = \frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x})$. Subtracting (A.14) and (A.13), we get

$$\text{Var}(\hat{I}_{SM}) - \text{Var}(\hat{I}_{IS}) = \frac{1}{N^2 Z^2} \int \left(\frac{N}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x})} - \sum_{i=1}^N \frac{1}{q_i(\mathbf{x})} \right) f^2(\mathbf{x})\pi^2(\mathbf{x}) d\mathbf{x}.$$

Hence, since $f^2(\mathbf{x})\pi^2(\mathbf{x}) \geq 0 \forall \mathbf{x}$, in order to prove the theorem it is sufficient to show that

$$\frac{1}{\frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x})} \leq \frac{1}{N} \sum_{i=1}^N \frac{1}{q_i(\mathbf{x})}. \quad (\text{A.15})$$

Now, let us note that the left hand side of (A.15) is the inverse of the arithmetic mean of $q_1(\mathbf{x}), \dots, q_N(\mathbf{x})$,

$$A_N = \frac{1}{N} \sum_{j=1}^N q_j(\mathbf{x}),$$

whereas the right hand side of (A.15) is the inverse of the harmonic mean of $q_1(\mathbf{x}), \dots, q_N(\mathbf{x})$,

$$\frac{1}{H_N} = \frac{1}{N} \sum_{i=1}^N \frac{1}{q_i(\mathbf{x})}.$$

Therefore, the inequality in (A.15) is equivalent to stating that $\frac{1}{A_N} \leq \frac{1}{H_N}$, or equivalently $A_N \geq H_N$, which is the well-known arithmetic mean–harmonic mean inequality for positive real numbers [27, 37]. Note that (A.15) can also be proved using Jensen's inequality in Eq. (16) with $\varphi(x) = \frac{1}{x}$, $\alpha_i = \frac{1}{N}$ and $z_i = q_i(\mathbf{x})$ for $i = 1, \dots, N$. \square

Proposition 3.

$$\text{Var}(\hat{I}_{DM}) \leq \text{Var}(\hat{I}_{SM}). \quad (\text{A.16})$$

Proof: The variance of \hat{I}_{DM} is computed

$$\begin{aligned}
\text{Var}(\hat{I}_{DM}) &= \frac{1}{N^2 Z^2} \sum_{i=1}^N \left(E_{q_i} \left[\frac{f^2(\mathbf{X}_i) \pi^2(\mathbf{X}_i)}{\psi^2(\mathbf{X}_i)} \right] - E_{q_i}^2 \left[\frac{f(\mathbf{X}_i) \pi(\mathbf{X}_i)}{\psi(\mathbf{X}_i)} \right] \right) \\
&= \frac{1}{N^2 Z^2} \sum_{i=1}^N \left(\int \frac{f^2(\mathbf{x}) \pi^2(\mathbf{x})}{\psi^2(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right) - \frac{1}{N^2 Z^2} \sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 \\
&= \frac{1}{N Z^2} \left(\int \frac{f^2(\mathbf{x}) \pi^2(\mathbf{x})}{\psi^2(\mathbf{x})} \left[\frac{1}{N} \sum_{i=1}^N q_i(\mathbf{x}) \right] d\mathbf{x} \right) - \frac{1}{N^2 Z^2} \sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 \\
&= \frac{1}{N Z^2} \int \frac{f^2(\mathbf{x}) \pi^2(\mathbf{x})}{\psi(\mathbf{x})} d\mathbf{x} - \frac{1}{N^2 Z^2} \sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 \tag{A.17}
\end{aligned}$$

Analyzing Eqs. (A.14) and (A.17), we see that proving $\text{Var}(\hat{I}_{DM}) \leq \text{Var}(\hat{I}_{SM})$ is equivalent to proving that

$$\begin{aligned}
\frac{1}{Z^2} \sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 &\geq N I^2 \\
\frac{1}{Z^2} \sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 &\geq N \left(\frac{1}{Z} \int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} \psi(\mathbf{x}) d\mathbf{x} \right)^2 \\
\sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 &\geq N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} \left(\frac{1}{N} \sum_{i=1}^N q_i(\mathbf{x}) \right) d\mathbf{x} \right)^2 \\
\sum_{i=1}^N \left(\int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 &\geq \frac{1}{N} \left(\sum_{i=1}^N \int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x} \right)^2 \tag{A.18}
\end{aligned}$$

By defining $a_i = a_i(\mathbf{x}) = \int \frac{f(\mathbf{x}) \pi(\mathbf{x})}{\psi(\mathbf{x})} q_i(\mathbf{x}) d\mathbf{x}$, (A.18) can be expressed more compactly as

$$N \sum_{i=1}^N a_i^2 \geq \left(\sum_{i=1}^N a_i \right)^2. \tag{A.19}$$

The inequality in Eq. (A.19) holds, since it corresponds to the definition of the Cauchy-Schwarz inequality [27],

$$\left(\sum_{i=1}^N a_i^2 \right) \left(\sum_{i=1}^N b_i^2 \right) \geq \left(\sum_{i=1}^N a_i b_i \right)^2, \tag{A.20}$$

with $b_i = 1$ for $i = 1, \dots, N$. Once more, (A.18) can also be proved by using Jensen's inequality in (16) with $\varphi(x) = x^2$, $\alpha_i = \frac{1}{N}$ and $z_i = a_i(\mathbf{x})$ for $i = 1, \dots, N$. \square

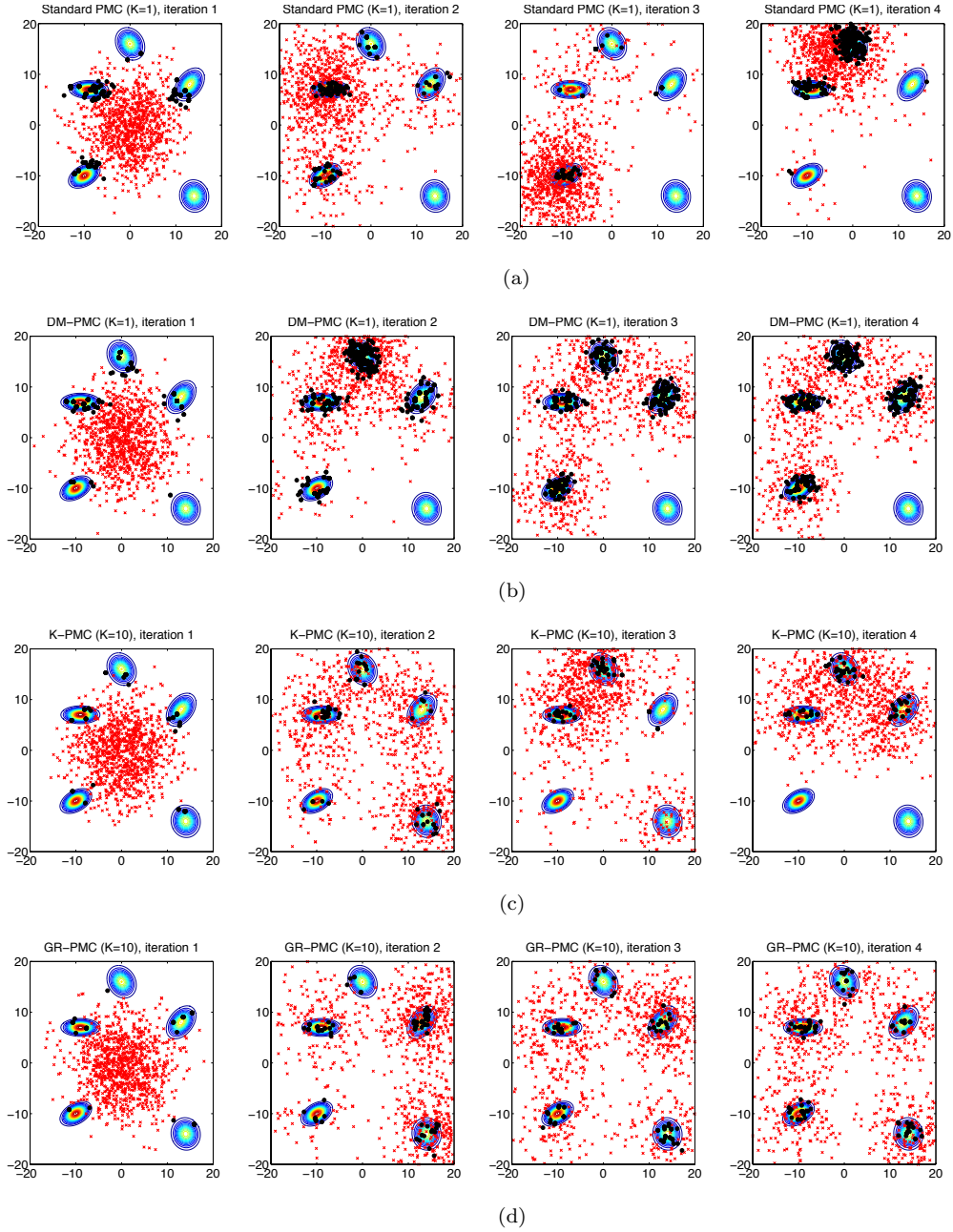


Figure A.7: **(Ex1-Section 5.2)** Evolution of the samples before (red crosses) and after resampling (black circles) for different schemes using $N = 100$ and $\sigma = 5$. The contour lines of the target density are also depicted. **(a)** Standard PMC. **(b)** DM-PMC ($K=1$). **(c)** K-PMC ($K=10$) with global resampling. **(d)** GR-PMC ($K=10$).