Joint Model Selection and Parameter Estimation by Population Monte Carlo Simulation

Mingyi Hong, Mónica F. Bugallo, and Petar M. Djurić

Abstract-In this paper, we study the problem of joint model selection and parameter estimation under the Bayesian framework. We propose to use the Population Monte Carlo (PMC) methodology in carrying out Bayesian computations. The PMC methodology has recently been proposed as an efficient sampling technique and an alternative to Markov Chain Monte Carlo (MCMC) sampling. Its flexibility in constructing transition kernels allows for joint sampling of parameter spaces that belong to different models. The proposed method is able to estimate the desired a posteriori distributions accurately. In comparison to the Reversible Jump MCMC (RJMCMC) algorithm, which is popular in solving the same problem, the PMC algorithm does not require burn-in period, it produces approximately uncorrelated samples, and it can be implemented in a parallel fashion. We demonstrate our approach on two examples: sinusoids in white Gaussian noise and direction of arrival (DOA) estimation in colored Gaussian noise, where in both cases the number of signals in the data is a priori unknown. Both simulations show the effectiveness of our proposed algorithm.

I. INTRODUCTION

Model selection is an important topic in signal processing. It has found application in various areas including array signal processing, communications, and speech signal processing and therefore has been studied extensively. A recent review provides some common approaches and criteria for model selection [1]. The model selection problem is often presented as a problem of joint model selection and parameter estimation. Many researchers have addressed it within the Bayesian framework. The main difficulty of this approach lies in solving multi-dimensional integrals. Some early works on model selection are based on the use of large sample theory and approximating the final posterior by Taylor expansion around the maximum likelihood (ML) estimates of the unknown parameters [2], [3], [4]. In [5], the authors developed an efficient iterative algorithm for carrying out the maximization needed for obtaining maximum a posterior (MAP) estimates. More recently, Reversible Jump Markov Chain Monte Carlo (RJMCMC) sampling [6] has been introduced for approximating joint posteriors and computing estimates of model order and parameters of interest [7], [8], [9]. Although computationally intensive, this algorithm was shown to have very good performance, especially when the sizes of available data are small. However, algorithms based on RJMCMC have several drawbacks. First, a burn-in period, whose samples are discarded, is required. Second, typical MCMC implementation may have poor mixing, i.e., the chain may converge to a

final distribution which depends on the starting point of the chain. Third, one may argue that it is hard to implement the algorithm in a parallel fashion because at each time step t, the algorithm produces only one sample, which depends on the sample produced in the previous time step.

PMC has recently been introduced in [10] and [11]. It is essentially an iterative sampling method which at each iteration employs importance sampling (IS) to produce a set of approximately uncorrelated samples from the target distribution. It also uses resampling [12] to prevent sample degeneration when needed. However, it is well known that for importance sampling, the importance function (IF) needs to be carefully chosen to ensure that the "region of importance" is reached quickly [13]. It is the ability of PMC to accommodate multiple importance functions (or rather, transition kernels), and to adaptively improve their sampling efficiency that makes it superior to pure importance sampling. Between different iterations, the algorithm can, based on certain criteria, change the structure of the transition kernels to ensure that the subsequent sampling procedure is carried out more efficiently.

In [11], a fixed number of pre-selected transition kernels have been used and each of them has been assigned different weights at different iterations. The efficiency of the algorithm has been demonstrated by an example with the posterior being a mixture Gaussian distribution. It has been shown that the produced samples by the algorithm accurately approximates the distribution. In [14], the authors have proposed an algorithm to adaptively choose transition kernels so that the asymptotic variance of the estimates decreases. In [15], it has been demonstrated that the PMC algorithm could progressively sample from distributions that had diminishing Kullback distance from the target distribution. Comparisons between MCMC and PMC have been made in [11] and [16]. In both cases, PMC has outperformed MCMC, mainly because of the slow mixing property of the MCMC. Improved performance of PMC in parameter estimation by the use of Rao-Blackwellization has been shown in [17].

In this paper, we propose to apply the PMC methodology to joint model selection and parameter estimation. We use a two-stage sampling procedure: we first sample the model order from a set of discrete transition kernels, and then we sample the parameters from a set of transition kernels that correspond to the sampled model. This two-stage sampling procedure allows us to sample from parameter spaces with different dimensions. When the samples are properly weighted, the samples and the weights produce approximations of the desired posteriors.

The paper is organized as follows. In Section II we pro-

The work has been supported by the National Science Foundation under Award CCF-0515246 and the Office of Naval Research under Award N00014-09-1-1154.

vide the formulation of the problem, and in Section III, we describe the steps of the algorithm. There we also give a proof of convergence of the algorithm. In Section IV, we demonstrate the performance of the proposed algorithm with two experiments: (a) detection of sinusoids in white Gaussian noise and estimation of their frequencies and (b) detection of number of sources whose signals impinge on an array of sensors. In the latter case the signals are corrupted by additive colored Gaussian noise and the objective, too, is to estimate the Direction of Arrival (DOA) of the signals. In both examples we compare our algorithm with the RJMCMC algorithm. Our main objective is to demonstrate that PMC is a valuable alternative to RJMCMC.

II. FORMULATION OF THE PROBLEM

In this section, we formulate the problem of joint model selection and parameter estimation in a Bayesian framework. Assume we have an observation vector \mathbf{y} which contains d_y data samples. We also have K competing models $\mathcal{M}_0, ..., \mathcal{M}_{K-1}$, and one of them generates the observations. Associated with each model, there is a vector of parameters $\boldsymbol{\theta}_k \in S_k$, where S_k denotes the parameter space of \mathcal{M}_k . The objective is to identify the true model as well as to estimate the parameters $\boldsymbol{\theta}_k$ associated with the model.

We can view the model order as a realization of a discrete random variable, so that the total parameter space Θ could be expressed as follows: $\Theta = \bigcup_{k=0}^{K-1} \{k\} \times S_k$, which is a union of disjoint subspaces. Note that each S_k may have different dimension and may include different parameters. The objective of Bayesian inference is to obtain the posterior $p(k, \theta | \mathbf{y})$, which can then be used (if needed) to compute point estimates.

In the Bayesian context, one typically employs the MAP model selection rule, which can be expressed as

$$k^{*} = \arg_{k} \max \left\{ p(k|\mathbf{y}) \right\}$$
$$= \arg_{k} \max \left\{ \int_{\boldsymbol{\theta} \in \mathcal{S}_{k}} p(k, \boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta} \right\}$$
$$= \arg_{k} \max \left\{ \int_{\boldsymbol{\theta} \in \cup_{\widetilde{k}=0}^{K-1} \mathcal{S}_{\widetilde{k}}} \Pi_{k}(\widetilde{k}) p(\widetilde{k}, \boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta} \right\}$$
(1)

where $\Pi_k(\tilde{k})$ is an indicator function that takes the value 1 when $\tilde{k} = k$ and is 0 otherwise. The difficulty of using (1) for model selection is that the posterior is usually highly nonlinear in θ , and the integration does not have a close form expression. In that case, one can resort to a Monte Carlo technique to approximate the integral $\int_{\Theta} f(\theta) p(\theta) d\theta$ by first drawing a sample $\theta^{(i)}$ of size N directly from the distribution $p(\theta)$, and then performing Monte Carlo integration by

$$\int_{\Theta} f(\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \simeq \frac{1}{N} \sum_{i=1}^{N} f(\boldsymbol{\theta}^{(i)}).$$
(2)

From the strong law of large numbers, the above approximation converges to the true value of the integration with probability one. Thus, if one can draw the samples $\{(k^{(i)}, \boldsymbol{\theta}^{(i)})\}_{i=1}^{N}$ directly from the posterior $p(k, \boldsymbol{\theta}|\mathbf{y})$, then (1) becomes

$$k^* \simeq arg_k \max\left\{\frac{1}{N}\sum_{i=1}^N \Pi_k(k^{(i)})\right\}.$$
 (3)

As stated above, $\Pi_k(k^{(i)})$ is an indicator function, and $\sum_{i=1}^N \Pi_k(k^{(i)})$ essentially calculates the total number of drawn samples $\{(k^{(i)}, \boldsymbol{\theta}^{(i)})\}_{i=1}^N$ with $k^{(i)} = k$. Thus, one can calculate k^* by first drawing samples $\{(k^{(i)}, \boldsymbol{\theta}^{(i)})\}_{i=1}^N$ from $p(k, \boldsymbol{\theta}|\mathbf{y})$, then selecting the model order k that is most frequently sampled.

It is well known that one can directly generate samples from the posterior distribution only in a very few cases. When such generation is impossible, one may resort to the use of importance sampling. In the following, we first briefly review this technique and then apply it to solve the model selection problem.

Importance sampling has been used mainly for numerical integration [13] and has many applications in signal processing and communications [18], [19]. For example, in signal processing, in order to obtain the minimum mean square (MMSE) estimate of a parameter x from the observation y, we have to solve the following integration:

$$\hat{x} = \int x p(x|y) dx$$

where p(x|y) is the posterior distribution of x, and \hat{x} is the MMSE estimate of x. The above integration is usually hard to solve directly. It may also be hard to draw samples $x^{(i)}$ directly from the distribution p(x|y) in order to perform the classical Monte Carlo numerical integration. Alternatively, we can draw samples $x^{(i)}$ from an *importance function* (IF) q(x), and perform numerical integration as follows:

$$\hat{x} = \int xp(x|y)dx$$
$$= \int x \frac{p(x|y)}{q(x)}q(x)dx$$
$$\simeq \frac{1}{N} \sum_{i=1}^{N} \frac{x^{(i)}p(x^{(i)}|y)}{q(x^{(i)})}$$

It has been shown that when the sample size N is large, the estimate \hat{x} obtained by the above importance sampling algorithm converges to the mean of the posterior. If we define $\tilde{w}^{(i)} = \frac{p(x^{(i)}|y)}{q(x^{(i)})} \frac{1}{N}$, the above estimate can be interpreted as the weighted mean of the drawn samples $x^{(i)}$. If p(x|y) is only known up to its proportionality constants (in other words, p(x|y) is unscaled), then we calculate $w^{(i)}$ by

$$\bar{w}^{(i)} = \frac{p(x^{(i)}|y)}{q(x^{(i)})}$$

and normalize the weights according to

$$w^{(i)} = \frac{\bar{w}^{(i)}}{\sum_{j=1}^{N} \bar{w}^{(j)}}$$

One interesting interpretation of the sample weight pairs $\{x^{(i)}, w^{(i)}\}$ is that they form a discrete random measure $\chi = \{x^{(i)}, w^{(i)}\}$, where the samples $x^{(i)}$ constitute the support of this measure, and $w^{(i)}$ are weights associated to these samples. When the sample size is large, the measure χ "approximates" the posterior distribution of x.

In light of the above importance sampling algorithm, we can draw samples $\{(k^{(i)}, \theta^{(i)})\}_{i=1}^N$ (where, $\theta^{(i)} \in S_{k^{(i)}}$) from an IF $g(k, \theta)$, and assign each sample a proper weight $w^{(i)}$, so that the random measure $\{(k^{(i)}, \theta^{(i)}), w^{(i)}\}_{i=1}^N$ approximates the posterior distribution. We can then approximate the integration in (1) according to

$$k^{*} = \arg_{k} \max\left\{ \int_{\boldsymbol{\theta} \in \bigcup_{\overline{k}=0}^{K-1} S_{\widetilde{k}}} \Pi_{k}(\widetilde{k}) p(\widetilde{k}, \boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \right\}$$
$$= \arg_{k} \max\left\{ \int_{\boldsymbol{\theta} \in \bigcup_{\overline{k}=0}^{K-1} S_{\widetilde{k}}} \Pi_{k}(\widetilde{k}) \frac{p(\widetilde{k}, \boldsymbol{\theta} | \mathbf{y})}{g(\widetilde{k}, \boldsymbol{\theta})} g(\widetilde{k}, \boldsymbol{\theta}) d\boldsymbol{\theta} \right\}$$
$$\simeq \arg_{k} \max\left\{ \frac{1}{N} \sum_{i=1}^{N} \Pi_{k}(k^{(i)}) \frac{p(k^{(i)}, \boldsymbol{\theta}^{(i)} | \mathbf{y})}{g(k^{(i)}, \boldsymbol{\theta}^{(i)})} \right\}$$
$$= \arg_{k} \max\left\{ \frac{1}{N} \sum_{i=1}^{N} \Pi_{k}(k^{(i)}) w^{(i)} \right\}$$
(4)

where

$$w^{(i)} \triangleq \frac{p(k^{(i)}, \boldsymbol{\theta}^{(i)} | \mathbf{y})}{g(k^{(i)}, \boldsymbol{\theta}^{(i)})}.$$
(5)

The above expression indicates that one can approximate the a posteriori probability of the model with index k, $p(k|\mathbf{y})$ by summing up all the weights of the samples $\{(k^{(i)}, \boldsymbol{\theta}^{(i)})\}_{i=1}^{N}$, where $k^{(i)} = k$. With the MAP criterion, we choose the model whose sum of weights is maximum. Moreover, one can use the weights to estimate the unknown parameters of each model k by

$$\widehat{\boldsymbol{\theta}}^{k} = \frac{\sum_{i=1}^{N} w^{(i)} \ \boldsymbol{\theta}^{(i)} \ \Pi_{k}(k^{(i)})}{\sum_{i=1}^{N} w^{(i)} \ \Pi_{k}(k^{(i)})}.$$
(6)

The above importance sampling procedure may suffer from poor choices of IF which can, for example, make most of the sample weights negligible and cause the subsequent estimation inaccurate. This situation may arise when the IF generates samples that concentrate on regions with low probability mass under the target posterior. The choice of the IF is both art and science, and many criteria for selecting good IFs have been developed [13]. One can, for example, choose an IF with heavy tails that dominate the tails of the target distribution, or work with an IF that mimics the behavior of the target distribution. However, both of these strategies require certain information about the target distribution, which is usually not available. In the following section, we propose to use a PMC algorithm which obtains samples from an IF that is adaptively modified so that the "quality" of the samples improves with iterations and one can evaluate the integral in (1) with greater accuracy.

III. PROPOSED ALGORITHM

PMC has the ability to progressively learn about the target distribution and to adaptively modify the IF based on the gathered information so that the sampling procedure becomes more efficient. We first introduce various elements of the PMC algorithm in Section III-A, and then present the algorithm used to solve our model selection problem in Section III-B. We present the convergence result for the proposed algorithm in Section III-C.

A. Generic PMC Methodology

One way to overcome the main difficulty encountered by IS - the poor choice of IFs - is to introduce into the sampling procedure multiple IFs with different properties, and iteratively and adaptively select them according to their performance. This iterative procedure is a learning process, during which we gain knowledge about the target distribution. For example, in the initial stages of the sampling, it is preferred to have IFs with heavy tails so that the parameter space could be explored fully, while in the later stages, the IFs with good local exploring property may be preferred to stabilize the samples. This process could be implemented in the following way. Before the start of the algorithm, D IFs g_1, \dots, g_D are selected. Define the vector of random variables $\mathbf{x} \triangleq [k \ \boldsymbol{\theta}]$ and the vector of samples $\mathbf{x}^{(i,t)} \triangleq [k^{(i,t)} \boldsymbol{\theta}^{(i,t)}]$, where t stands for iteration. Then the overall IF $q^{(t)}(\mathbf{x})$ can be constructed at each iteration t as follows

$$g^{(t)}(\mathbf{x}) = \sum_{d=1}^{D} \alpha_d^{(t)} g_d(\mathbf{x})$$
$$\sum_{d=1}^{D} \alpha_d^{(t)} = 1.$$
(7)

It is clear that $g^{(t)}(\mathbf{x})$ is a mixture of D functions and that $\alpha_d^{(t)}$ could be interpreted as frequency of using the IF $g_d(\mathbf{x})$ for sampling. The following procedure is used to sample from $g^{(t)}(\mathbf{x})$. Let $d^{(i,t)} \in \{1, \dots, D\}$ denote the index of the IF that is used for sampling $\mathbf{x}^{(i,t)}$. Then $d^{(i,t)}$ is determined by drawing it from a multinomial distribution $\mathcal{M}(\alpha_1^{(t)}, \dots, \alpha_D^{(t)})$ followed by generating $\mathbf{x}^{(i,t)}$ from $g_{d^{(i,t)}}(\mathbf{x})$.

The significance of the above construction of $g^{(t)}(\mathbf{x})$ is that, one can modify $g^{(t)}(\mathbf{x})$ at each iteration by changing the weights $\alpha_d^{(t)}$, according to the performance of each g_d in the past. One criterion to evaluate the performance of the g_d is the sum of the weights of the samples they generate [11]. This criterion favors the IFs that focus on exploring the so called "region of importance" of the target distribution, and thus generate samples with large total weights. According to this criterion, $\alpha_d^{(t)}$ is updated by the following equation:

$$\alpha_d^{(t)} = \frac{\sum_{i=1}^N w^{(i,t-1)} \Pi_{g_d}(\mathbf{x}^{(i,t-1)})}{\sum_{i=1}^N w^{(i,t-1)}}$$
(8)

where $\Pi_{g_d}(\mathbf{x}^{(i,t-1)})$ is an indicator function that takes a value 1 if the sample $\mathbf{x}^{(i,t-1)}$ is generated by the IF g_d , and 0 otherwise.

These time-varying IFs $g^{(t)}(\mathbf{x})$ generate N samples at each iteration, and with the weights $\{w^{(i,t)}\}_{i=1}^N$, they form a random measure $\mu^{(t)} \triangleq \{\mathbf{x}^{(i,t)}, w^{(i,t)}\}_{i=1}^N$ that approximates the distribution of the random variable $\mathbf{x}^{(t)}$. The objective of the algorithm is to ensure that this distribution converges to the target distribution when t is large.

For monitoring the sampling efficiency, we can use the entropy relative to uniformity [20], i.e.,

$$H^{(t)} = -\sum_{i=1}^{N} w^{(i,t)} \frac{\log w^{(i,t)}}{\log N}$$
(9)

where $H^{(t)}$ is a measure of uniformity of the weights $w^{(i,t)}$ at time t. If the IF converges to the target distribution, then the weight of each sample $\mathbf{x}^{(i,t)}$ converges to $\frac{1}{N}$, and $H^{(t)}$ converges to 1, and every sample can be seen as drawn approximately from the target distribution.

For the presented iterative IS algorithm, it is beneficial to relate the samples in a current stage with the samples from the previous stage. PMC introduces dependence of the current samples $\{\mathbf{x}^{(i,t)}\}_{i=1}^N$ on the samples in the previous iteration $\{\mathbf{x}^{(i,t-1)}\}_{i=1}^N$ by replacing the IFs $g_d(\mathbf{x})$ (which do not depend on the past) with transition kernels. This idea is closely related to smoothly approximating the posterior distribution using the kernel technique [20], except that PMC uses multiple kernels at each stage. Recall that $\mathbf{x}^{(t)}, t \in \{1, 2, ..., T\}$ is the set of random variables approximated by the iterative importance sampling algorithm at each iteration t. A transition kernel $Q(\mathbf{x}^{(i,t-1)}, \mathbf{x})$ is defined as follows:

$$Q(\mathbf{x}^{(i,t-1)},\mathbf{x}) = \frac{\partial P(\mathbf{x}^{(t)} \le \mathbf{x} | \mathbf{x}^{(t-1)} = \mathbf{x}^{(i,t-1)})}{\partial \mathbf{x}}.$$
 (10)

The transition kernel is a generalization of a transition matrix of discrete state Markov chains, where $Q(\mathbf{x}^{(i,t-1)}, \mathbf{x})$ is the probability density of the state of the chain at iteration t conditional on the chain being in state $\mathbf{x}^{(i,t-1)}$ at iteration t-1. We can replace our previous IF with a transition kernel $Q(\mathbf{x}^{(i,t-1)},\mathbf{x})$, where $Q(\mathbf{x}^{(i,t-1)},\mathbf{x})$ can take the form, for example, of a Gaussian kernel (in the case when \mathbf{x} is one-dimensional): $Q(x^{(i,t-1)},x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-x^{(i,t-1)})^2}{2\sigma^2}\}$. It is clear that the new samples $x^{(i,t)}$ drawn from this kernel will be around $x^{(i,t-1)}$, and the shape of the kernel is determined by σ^2 . In this case, the kernel $Q(x^{(i,t-1)},x)$ is a family of normal distributions with $x^{(i,t-1)}$ as location parameter.

PMC uses a mixture of multiple kernels $\{Q_d(., \mathbf{x})\}_{d=1}^D$ in each iteration to improve the sampling efficiency, and the sampling is performed as in (7), i.e.,

$$\mathbf{x}^{(i,t)} \sim \sum_{d=1}^{D} \alpha_d^{(t)} Q_d(\mathbf{x}^{(i,t-1)}, \mathbf{x})$$
(11)

$$\sum_{d=1}^{D} \alpha_d^{(t)} = 1$$
 (12)

where $\alpha_d^{(t)}$ is determined as in (8).

We summarize the procedure by explaining how we generate one particular sample at iteration t. First, we select the index d of the transition kernel from the multinomial distribution $\mathcal{M}(\alpha_1^{(t)}, \cdots, \alpha_D^{(t)})$; then we generate a sample $\mathbf{x}_i^{(t)}$ by the kernel with index d. Because the sample $\mathbf{x}_i^{(t)}$ is actually generated by the kernel $Q_d(\mathbf{x}^{(i,t-1)}, \mathbf{x}^{(i,t)})$, the weight of this sample can be obtained by

$$w^{(i,t)} \propto \frac{p(\mathbf{x}^{(i,t)}|\mathbf{y})}{Q_d(\mathbf{x}^{(i,t-1)}, \mathbf{x}^{(i,t)})}$$
(13)

which follows from (5). We use the proportionality operator \propto here to include the case when the posterior distribution is not scaled. Equation (13) represents the weight in the original PMC algorithm, and it has been used in many applications of the PMC, for example, [16], [17], [21].

We summarize a generic PMC implementation as follows. For each iteration t,

- 1. Generate samples from the mixture of transition kernels.
- 2. Compute the weights of the samples $w^{(i,t)}$.
- 3. Perform estimation based on the weights.
- 4. Resample.
- 5. Compute the weights $\alpha_d^{(t+1)}$ for each kernel $Q_d(\cdot, \cdot)$.

B. PMC for Model Selection

samples after resampling.

In a parameter estimation problem only, as demonstrated in [11], [17], the above steps are sufficient to produce approximation of the target distribution. Under model uncertainty, a direct extension of the above method would be to run K parallel PMC algorithms, one for each model, and compare their performance to determine the model. Of course, this naive approach is very computationally expensive when the number of competing models K is large, because the computational load is "equally" distributed upon the different models. Similar observation, which serves as inspiration for using RJMCMC instead of multiple MCMC under model uncertainty, has been made in [8].

Observe that there is a natural hierarchy in our full parameter space Θ : the space for the parameter θ is determined only by the choice of k. After the determination of k, it is sufficient to sample only from the parameter space S_k to produce samples that approximate the distribution $p(\theta|k, \mathbf{y})$. It is thus natural to decompose our sampling kernels into two components: one for sampling of the model order k from the index set $\{k\}_{k=0}^{K-1}$, and the other for sampling from the parameter space S_k . Let $Q(\cdot, \cdot)$ denote the kernel we use for sampling the model order, and let $Q_{\theta_k}(\cdot, \cdot)$ denote the kernel we use for sampling θ^k . We propose the following two-stage sampling scheme:

- At iteration t,
- 1) draw a model order sample $k^{(i,t)}$ from $Q(k^{(i,t-1)},k)$ and
- 2) draw a parameter sample $\theta^{(i,t)}$ from $Q_{\theta_{k^{(i,t)}}}(\widehat{\theta}_{k^{(i,t)}}^{(t-1)}, \theta)$ where $\widehat{\theta}_{k^{(i,t)}}^{(t-1)}$ is the estimate of the parameter of the specific model at iteration t-1.

It is worth mentioning that the model order is a discrete random variable, so that the kernel $Q(\cdot, \cdot)$ is represented by a $K \times K$ matrix with transition probabilities. For example, if $Q(\cdot, \cdot) = \mathbf{I}_{K \times K}$, which is a K by K identity matrix, we will always have $k^{(i,t-1)} = k^{(i,t)}$.

In the following we use D predetermined kernels $Q_d(.,.)$ for generating model order, where $d = 1, 2, \dots, D$, and for each model k, we have a single kernel $Q_{\theta_k}(\cdot, \cdot)$ for generating the parameters that belong to the parameter space S_k . One reason of using multiple model kernels is to improve the sampling efficiency, as stated in Section III-A. An extension can be easily made to generate parameters using multiple kernels under each model. However, the rationale is the same, and we keep our choice of single kernel for the parameters in order to maintain the presentation clear.

Based on the above decomposition of transition kernels, the sampling of $(k^{(i,t)}, \theta_{k^{(i,t)}}^{(i,t)})$ can be expressed by

$$(k^{(i,t)},\boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)}) \sim \sum_{d=1}^{D} \alpha_d^{(t)} Q_d(k^{(i,t-1)},k) Q_{\theta_k}(\widehat{\boldsymbol{\theta}}_k,\boldsymbol{\theta})$$

and the weight for each sample can be expressed according to (13) as follows:

$$w^{(i,t)} \propto \frac{p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} | \mathbf{y})}{Q_{d_{i}^{(t)}}(k^{(i,t-1)}, k^{(i,t)}) Q_{\boldsymbol{\theta}_{k^{(i,t)}}}(\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})}$$
(14)

where $\sum_{i=1}^{N} w^{(i,t)} = 1$. The estimate of the marginal posterior for the model order can be obtained by summing up the normalized weights for each model

$$P(k|\mathbf{y}) \simeq \sum_{i=1}^{N} w^{(i,t)} \Pi_k(k^{(i,t)}).$$

The resampling procedure equalizes the weight of the samples, and is commonly carried out by duplicating samples proportional to their weights [24]. Namely, for $i = 1, \dots, N$, we first sample the index $j_i \sim \mathcal{M}(w^{(1,t)}, \dots, w^{(N,t)})$ and then we let $\mathbf{x}^{*(i,t)} = \mathbf{x}^{(j_i,t)}$. After resampling, the models with large total weights (thus large estimated marginal posterior) occupy greater portion of the sample.

Now, we can provide another justification for using multiple kernels $Q_d(\cdot, \cdot)$ for sampling the model order. Let U denote a matrix with entries all ones. By choosing either $Q_d(\cdot, \cdot) =$ $\frac{1}{K}$ U or $Q_d(\cdot, \cdot) =$ I, and assigning $\alpha_d^{(1)} = \frac{1}{D}$ for all d in the initial stage of the algorithm, there will always be equal number of samples for each model at every iteration, which corresponds to the naive approach mentioned in the beginning of this section. Consequently, it is desirable to find a transition kernel that distributes the right amount of computational time to each competing model, preferably according to their true posterior, i.e., the models with high posteriors should get more "attention". Since the true posterior is not known, one can use multiple transition kernels and let the algorithm choose the most efficient distribution of computational time. In Section IV, we demonstrate that, indeed, by using D different transition kernels, we get better performance than when we use only one kernel.

We also provide some heuristic guidelines for designing the set of predetermined kernels $Q_d(\cdot, \cdot)$ for model order. First, we would like to build our current computation based on the distribution obtained from the previous iteration. As a result, we would expect that the majority of the model order samples $\{k^{(i,t)}\}_{i=1}^N$ generated at iteration t correspond to the model orders with the largest total weights in iteration t - 1. Second, we would also like to improve upon the previous distribution by exploring the parameter spaces more thoroughly. One way to achieve this is to allow portions of the majority samples in iteration t-1 to "move" to the other model orders, and this behavior is determined by the off-diagonal At t = 0For $i = 1, \dots, N$ Draw $k^{(i,0)}$ from the prior p(k)Draw $\theta^{(i,0)}$ from the prior $p(\theta_{k^{(i,0)}})$ Compute the normalized weights $w^{(i,0)} \propto p(\mathbf{y}|k^{(i,0)}, \theta_{k^{(i,0)}}^{(i,0)})$ Obtain estimates $\hat{k}^{(0)}, \hat{\theta}^{(0)}$ using $\{w^{(i,0)}\}_{i=1}^{N}$ and the samples Assign $\alpha_d^{(1)} = 1/D, d = 1, \dots, D$ Resample and get $k^{*(i,0)}, \theta^{*(i,0)}$ For $t = 1, \dots, T$ For $i = 1, \dots, N$ Determine $d^{(i,t)}$ by sampling from $\mathcal{M}(\alpha_1^{(t)}, \dots, \alpha_D^{(t)})$ Draw $k^{(i,t)}$ from $Q_{d^{(i,t)}}(k^{*(i,t-1)}, k)$ Draw $\theta_{k^{(i,t)}}^{(i,t)}$ from $Q_{\theta_{k^{(i,t)}}}(\hat{\theta}_{k^{(i,t)}}^{(t-1)}, \theta)$ Compute the normalized weights $w^{(i,t)} \propto \frac{p(\mathbf{y}|k^{(i,t)}, \theta_{k^{(i,t)}})}{Q_{d^{(i,t)}}(k^{*(i,t-1)}, k^{(i,t)})Q_{\theta_{k^{(i,t)}}}(\hat{\theta}_{k^{(i,t)}}^{(t-1)}, \theta_{k^{(i,t)}}^{(i,t)})}$ Obtain estimates $\hat{k}^{(t)}, \hat{\theta}^{(t)}$ Assign $\alpha_d^{(t+1)} = \sum_{i=1}^N w^{(i,t)} \Pi_d(d^{(i,t)})$ Resample to obtain $k^{*(i,t)}$ and $\theta^{*(i,t)}$.



entries of $Q_d(\cdot, \cdot)$. Specifically, if $Q_d(2,3) = 0.5$, then 50 percent of the samples that have model order $k^{(i,t-1)} = 2$ at iteration t-1 will be moved to $k^{(i,t)} = 3$ in the next iteration. Now it is clear that the above requirements are somewhat conflicting with each other: the first requirement translates to the strategy of selecting $Q_d(\cdot, \cdot)$ that have large diagonal entries, while the second requirement dictates how to choose $Q_d(\cdot, \cdot)$ with off-diagonal entries not too small. As a result, we suggest to design multiple kernels that achieve a trade-off between these requirements: some kernels may have negligible off-diagonal entries, while some may have relatively large off-diagonal entries. Our choice of transition kernels in the following simulation section is an example of the above suggested design. We will also show in the simulations that strategies that do not obey the above suggestions may lead to poor performance of the PMC algorithm.

In Table I, we summarize the proposed algorithm and in Fig. 1, we present a graphical illustration of it.

C. Convergence

A natural question is whether the above stated algorithm can approximate the target distribution, or specifically, if for each iteration t, when the number of samples goes to infinity, we have

$$\lim_{N \to \infty} \sum_{i=1}^{N} w^{(i,t)} f(\mathbf{x}^{(i,t)}) = \int f(\mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$
(15)

for all function $f(\cdot)$ that satisfy certain regularity conditions. If (15) is true, then when the function $f(\cdot)$ is $\Pi_k(k^{(i,t)})$, the



Fig. 1. Graphical illustration of the proposed algorithm.

above equation becomes

$$\lim_{N \to \infty} \sum_{i=1}^{N} w^{(i,t)} \Pi_k(k^{(i,t)})$$
$$= \sum_{j=1}^{K} \Pi_k(j) \int p(j, \boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}$$
(16)
$$= p(k | \mathbf{y}).$$
(17)

We prove the following theorem regarding the convergence of the PMC algorithm for model selection.

Theorem 1: For the PMC algorithm detailed in Table I, we have the following convergence result:

$$\lim_{N \to \infty} \sum_{i=1}^{N} w^{(i,t)} \Pi_k(k^{(i,t)}) = p(k|\mathbf{y}).$$
(18)

Proof: The proof is given in the Appendix.

D. Discussion

It is easy to extend the above algorithm to support multiple kernels for exploring the parameter space $S_{k^{(i)}}$ instead of using a single kernel. Alternatively, we can adaptively change the single kernel $Q_{\theta_{k^{(i)}}}(\cdot, \cdot)$ to achieve improved sampling efficiency. Assume that $Q_{\theta_{k^{(i,t)}}}(\theta^{(t)}, \theta^{(t+1)})$ is a Gaussian kernel $\mathcal{N}(\theta^{(t)}, \sigma_{k_{i,t}}^2 \mathbf{I})$. We propose to use a time-varying kernel $Q_{\theta_{k^{(i,t)}}}(\cdot, \cdot) = \mathcal{N}(\theta^{(t)}, \sigma_{k^{(i,t)}}^2 \mathbf{I})$ as follows. If $k^{(i,t)}$ is the MAP estimate of the model order at iteration t, let

$$\sigma_{k^{(i,t+1)}}^2 = \frac{t}{t+c} \sigma_{k^{(i,t)}}^2$$
(19)

and else, let

$$\sigma_{k^{(i,t+1)}}^2 = \frac{t}{t+c} \sigma_{k^{(i,t)}}^2 + \frac{c}{t+c} \sigma_{k^{(i,1)}}^2$$
(20)

where c is a constant. It is clear that if t is large enough and the estimates for $k^{(i,t)}$ stabilize, then $\sigma_{k^{(i,t)}}^2$ stabilizes for all $k^{(i,t)}$. By setting $\sigma_{k^{(i,1)}}^2$ relatively large, the above scheme allows us to have kernels with heavy tails at the initial stages of the algorithm so that we can explore the whole parameter space, and have lighter tails to focus on local sampling in the later stages of the algorithm, Once the order estimates stabilize, $\sigma_{k^{(i,t)}}^2$ shrinks and allows the kernel to explore the local parameter space instead. We will demonstrate the improvement of the sampling efficiency for this choice of $Q_{\theta_{k^{(i,t)}}}(\cdot, \cdot)$.

It is also worth mentioning that our proposed algorithm is capable of sampling nuisance parameters if their conditional distributions are known. Specifically, denote with z the vector of nuisance parameters. Then we can sample z at iteration t according to $\mathbf{z}^{(i,t)} \sim p(\mathbf{z}^{(i,t)}|k^{*(i,t)}, \boldsymbol{\theta}^{*(i,t)}, \mathbf{y})$, where $\boldsymbol{\theta}^{*(i,t)}$ and $k^{*(i,t)}$ are the resampled parameters of interest and the model order, respectively.

IV. SIMULATIONS

In this section, we present extensive Monte Carlo simulation results of the proposed algorithm. Two simulation experiments are considered. The first simulation experiment is the joint detection and estimation of sinusoids in white Gaussian noise [5]. The second simulation experiment is the joint detection of number of sources and estimation of direction-of-arrival (DOA) of signals emitted by the sources [25]. The signals impinge on an array of sensors and are corrupted by colored Gaussian noise. In subsections IV-A and IV-B, we briefly present the mathematical model for these problems. For their detailed description, we refer the readers to [8] and [9]. In subsections IV-C and IV-D, we present simulation setups and results for each problem, respectively.

A. Detection and Estimation of Sinusoids in White Noise

We have an observation vector \mathbf{y} with d_y data samples. The observations are generated by one of the following K models:

$$\mathcal{M}_0: y[n] = \epsilon[n]$$

$$\mathcal{M}_k: y[n] = \sum_{j=1}^k a_{c,j} \cos(2\pi f_j n) + a_{s,j} \sin(2\pi f_j n) + \epsilon[n]$$

where $n = 0, 1, \dots, d_y - 1, k = 1, \dots, K - 1, \epsilon[n]$ are iid noise samples, and $\epsilon[n] \sim \mathcal{N}(0, \sigma^2)$. In matrix form, the observation vector can be expressed as follows:

$$\mathbf{y} = \mathbf{H}(\mathbf{f}_k)\mathbf{a}_k + \boldsymbol{\epsilon} \tag{21}$$

where \mathbf{y} is the $d_y \times 1$ observation vector, $\boldsymbol{\epsilon}$ is a $d_y \times 1$ Gaussian noise vector, i.e., $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, \mathbf{f}_k is a $k \times 1$ frequency vector defined as $\mathbf{f}_k \triangleq [f_1 \ f_2 \ \cdots \ f_k]^\top$, \mathbf{a}_k is a $2k \times 1$ amplitude vector, and $\mathbf{H}(\mathbf{f}_k)$ is a $d_y \times 2k$ matrix whose elements can be expressed according to

$$\mathbf{H}(\mathbf{f}_k)_{n,2j-1} = \cos(2\pi f_j n)$$
$$\mathbf{H}(\mathbf{f}_k)_{n,2j} = \sin(2\pi f_j n).$$

The unknown parameter vector is $\boldsymbol{\theta}_k \triangleq [\sigma^2 \mathbf{a}_k^\top \mathbf{f}_k^\top]^\top$. Our objective is to jointly determine which model generated the observations \mathbf{y} and estimate the parameters \mathbf{f}_k . Note that σ^2 and \mathbf{a}_k are treated as nuisance parameters. Therefore, we use the Bayesian methodology: we integrate out the nuisance parameters and determine the model order and frequency from the joint posterior distribution $p(k, \mathbf{f}_k | \mathbf{y})$.

We first assign prior distributions to the parameters. We assume the models have equal prior probability, thus,

$$p(k) = \frac{1}{K}.$$
(22)

Note, that in the RJMCMC setting, the prior for model order is usually set to be a truncated Poisson distribution $p(k) \propto \Lambda^k exp(-\Lambda)/k! \times \Pi\{0 \le k \le K-1\}$ with hyperparameter Λ in order to facilitate the 'Birth' and 'Death' moves. We use the Jeffreys' prior for the noise variance [26],

$$p(\sigma^2) \propto \frac{1}{\sigma^2}.$$
 (23)

Then we can write the joint posterior distribution $p(k, \theta_k | \mathbf{y})$ as follows.

$$p(k, \boldsymbol{\theta}_{k} | \mathbf{y}) \propto p(\mathbf{y}, k, \boldsymbol{\theta}_{k})$$

= $p(\mathbf{y} | k, \boldsymbol{\theta}_{k}) p(\boldsymbol{\theta}_{k}, k)$
 $\propto p(\mathbf{y} | k, \mathbf{a}_{k}, \mathbf{f}_{k}, \sigma^{2}) p(k, \mathbf{a}_{k}, \mathbf{f}_{k} | \sigma^{2}) p(\sigma^{2}) (24)$

Assuming, that \mathbf{f}_k has a uniform prior on $[0, \frac{1}{2}]^k$ and \mathbf{a}_k has a zero mean normal prior, we have

$$p(k, \mathbf{a}_k, \mathbf{f}_k | \sigma^2) \propto \frac{1}{|2\pi\sigma^2 \boldsymbol{\Sigma}_k|^{1/2}} \exp\left\{-\frac{\mathbf{a}_k^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{a}_k}{2\sigma^2}\right\} 2^k \quad (25)$$

where $\Sigma_k^{-1} = \delta^{-2} \mathbf{H}^{\top}(\mathbf{f}_k) \mathbf{H}(\mathbf{f}_k)$, with δ^2 being a hyperparameter that can be integrated out numerically if we choose for it an inverse gamma prior of the form $\mathcal{IG}(\alpha_{\delta}, \beta_{\delta})$. (See section V-C of [8] for detailed discussion.)

After integrating out \mathbf{a}_k and σ^2 in (24), we obtain

$$p(k, \mathbf{f}_k | \mathbf{y}) \propto (\mathbf{y}^\top \mathbf{P}_k^\perp \mathbf{y})^{-d_y/2} \frac{2^k}{(\delta^2 + 1)^k}$$
(26)

where

$$\mathbf{P}_k^{\perp} = \mathbf{I} - \mathbf{H}(\mathbf{f}_k) \mathbf{M}_k^{-1} \mathbf{H}^{ op}(\mathbf{f}_k)$$

and

$$\mathbf{M}_k = \mathbf{H}^\top(\mathbf{f}_k)\mathbf{H}(\mathbf{f}_k) + \boldsymbol{\Sigma}_k^{-1}.$$

In the simulations in subsection IV-C, we show that our PMC algorithm is capable of approximating the marginalized distribution (26).

B. Detection of number of sources and estimation of DOA

We have an $M \times d_y$ complex observation matrix Y. Each column of Y, $\mathbf{y}[n]$, is an $M \times 1$ vector representing the data received by a linear array of M sensors that can be expressed as

$$\mathbf{y}[n] = \mathbf{H}(\boldsymbol{\phi}_k)\mathbf{a}[n] + \boldsymbol{\epsilon}[n]$$
(27)

where $\epsilon[n]$ is an $M \times 1$ zero mean Gaussian noise with covariance matrix Σ , $\mathbf{a}[n]$ is a $k \times 1$ amplitude vector, and

 $\mathbf{H}(\boldsymbol{\phi}_k)$ is an $M\times k$ matrix whose elements can be expressed by

$$\mathbf{H}(\boldsymbol{\phi}_k)_{m,l} = \exp\{j(m-1)\boldsymbol{\phi}_l\}$$
(28)

where $j = \sqrt{-1}$, $m = 0, 1, \dots, M-1$ and $l = 1, \dots, k$. The vector ϕ_k is defined by $\phi_k \triangleq [\phi_1, ..., \phi_k]^\top$, where ϕ_l is given by $\phi_l = \omega_0 \lambda_0 \sin(\varphi_l)$ with φ_l being the angle between the l^{th} incident signal and the sensor array, ω_0 the carrier frequency of the received signal, λ_0 the distance between the sensors, and v the propagation speed of the signal. In summary, each element $y_m[n]$ of $\mathbf{y}[n]$ could be expressed as

$$y_m[n] = \sum_{l=1}^k a_l[n] \exp\{j(m-1)\phi_l\} + \epsilon_m[n].$$
 (29)

We need to determine which one of the following K models generates y:

$$\mathcal{M}_0: \mathbf{y}[n] = \boldsymbol{\epsilon}[n]$$

$$\mathcal{M}_k: \mathbf{y}[n] = \mathbf{H}(\boldsymbol{\phi}_k)\mathbf{a}[n] + \boldsymbol{\epsilon}[n] \quad k = 1, \cdots, K - 1.$$

Besides determining the model order, we need to estimate ϕ_k . Again we integrate out the nuisance parameters $\mathbf{A} \triangleq [\mathbf{a}[1] \mathbf{a}[2] \dots \mathbf{a}[d_y]]$ and $\boldsymbol{\Sigma}$, and estimate the model order k and ϕ_k using the marginalized posterior distribution $p(k, \phi_k | \mathbf{y})$.

Define the projector on the signal subspace by [5]

$$\mathbf{P}_{k} = \mathbf{H}(\phi_{k}) \left(\mathbf{H}(\phi_{k})^{\mathsf{H}} \mathbf{H}(\phi_{k}) \right)^{-1} \mathbf{H}(\phi_{k})^{\mathsf{H}}$$
$$= \mathbf{U}_{s}(\phi_{k}) \mathbf{U}_{s}^{\mathsf{H}}(\phi_{k})$$
(30)

and let the projector on the noise subspace be

$$\mathbf{P}_{k}^{\perp} = \mathbf{I} - \mathbf{P}_{k} = \mathbf{U}_{\epsilon}(\boldsymbol{\phi}_{k}) \mathbf{U}_{\epsilon}^{\mathsf{H}}(\boldsymbol{\phi}_{k})$$
(31)

where $\mathbf{U}_s(\phi_k)$ and $\mathbf{U}_\epsilon(\phi_k)$ are orthogonal matrices whose columns span the signal subspace and the noise subspace, respectively. Define

$$\mathbf{z}_{s}[n] = \mathbf{U}_{s}^{\mathsf{H}}(\boldsymbol{\phi}_{k})\mathbf{y}[n]$$
(32)

$$\mathbf{z}_{\epsilon}[i] = \mathbf{U}_{\epsilon}^{\mathsf{H}}(\boldsymbol{\phi}_{k})\mathbf{y}[n].$$
(33)

Let $\mathbf{Z}_s \triangleq [\mathbf{z}_s[1], \cdots, \mathbf{z}_s[d_y]]$, and define \mathbf{Z}_{ϵ} similarly. We can show that

$$p(\mathbf{Z}_{s}, \mathbf{Z}_{\epsilon} | \tilde{\mathbf{A}}, \boldsymbol{\phi}, k, \mathbf{W}^{-1}) \simeq \pi^{-d_{y}k} | \mathbf{C}^{-1} |^{d_{y}}$$

$$\times \exp\left\{-\sum_{n=1}^{d_{y}} (\mathbf{z}_{s}[n] - \tilde{\mathbf{a}}[n])^{\mathsf{H}} \mathbf{C}^{-1} (\mathbf{z}_{s}[n] - \tilde{\mathbf{a}}[n])\right\}$$

$$\pi^{-d_{y}(M-k)} | \mathbf{W}^{-1} |^{d_{y}} \exp\left\{-\sum_{n=1}^{N} \mathbf{z}_{\epsilon}^{\mathsf{H}}[n] \mathbf{W}^{-1} \mathbf{z}_{\epsilon}[n]\right\}$$
(34)

where $\mathbf{A} = [\mathbf{\tilde{a}}[1] \ \mathbf{\tilde{a}}[2] \ \cdots \ \mathbf{\tilde{a}}[d_y]], \ \mathbf{\tilde{a}}[n] \triangleq \mathbf{U}_s^{\mathsf{H}}\mathbf{H}\mathbf{a}[n], \ \mathbf{C} \triangleq \mathbf{U}_s^{\mathsf{H}}\mathbf{\Sigma}\mathbf{U}_s$, and $\mathbf{W} \triangleq \mathbf{U}_{\epsilon}^{\mathsf{H}}\mathbf{\Sigma}\mathbf{U}_{\epsilon}$. We will then assign the priors to the different sets of parameters and integrate out the nuisance

parameters. Let [25]

$$p(k) = \frac{1}{K} \tag{35}$$

$$p(\phi_k|k) = U[0, 2\pi]^k$$
(36)
$$\mathbf{W}^{-1}|\phi_k, k) \propto |\mathbf{W}^{-1}|^{-(M-k)}$$
(37)

$$p\left(\tilde{\mathbf{A}}|\boldsymbol{\phi}_{k},k,\mathbf{W}^{-1}\right) = \prod_{n=1}^{d_{y}} \mathcal{N}(0,\rho^{2}\mathbf{I}_{k})$$
(38)

where ρ^2 is a hyperparameter that can be determined numerically (see Section V of [9] for discussion). Finally, after integrating out \mathbf{W}^{-1} and $\tilde{\mathbf{A}}$, we get

$$p(k, \phi | \mathbf{Z}_{\epsilon}) \propto \frac{\pi^{1/2(M-k)(M-k-1)} \prod_{n=1}^{M-k} \Gamma(N_y - n + 1)}{(2\pi)^k (\rho^2)^{kN_y}} \times |\widehat{\mathbf{R}}|^{-d_y}$$
(39)

where $\widehat{\mathbf{R}} \triangleq \sum_{i=1}^{d_y} \mathbf{z}_{\epsilon}[n] \mathbf{z}_{\epsilon}[n]^{\mathsf{H}}$, and $\Gamma(d_y - n + 1)$ is the gamma function with argument $d_y - n + 1$. In the simulation in section IV-D, we show that our PMC algorithm can approximate (39).

C. Simulation Results for Sinusoids in White Noise

We used the following setup for the experiment. We set $d_y = 64$, k = 2, and $f_1 = 0.2$, $f_2 = 0.2 + 1/(l_e \times d_y)$, where $l_e = 1, 2, 3, 4$, $-\arctan(a_{s_1}/a_{c_1}) = 0$ and $-\arctan(a_{s_2}/a_{c_2}) = \pi/4$. We tested the detection performance for SNR=3dB and SNR=10dB. The SNR was defined as $10 \log_{10} \left\{ \frac{a_{s_1}^2 + a_{c_1}^2}{2\sigma^2} \right\}$, and both sinusoids had the same SNR.

In our simulations we assumed K = 5. For the PMC algorithm we generated 3,000 samples in each iteration, and we ran the algorithm for a total of 10 iterations. As a prior for the hyperparameter δ^2 , we used $\mathcal{IG}(2, 10)$. We employed three 5×5 matrices as transition kernels for model order: $Q_1(i,i) = 0.4, Q_1(i,j) = 0.15, \text{ for } i \neq j; Q_2(i,i) = 0.80,$ $Q_2(i,j) = 0.05, \text{ for } i \neq j; Q_3(i,i) = 0.92, Q_3(i,j) =$ 0.02, for $i \neq j$. As discussed in the end of section III-B, by using these transition matrices, at each iteration t of the algorithm, a majority of the model order samples $\{k_{i,j}\}_{j=1}^N$ represented the models with large total weights in the previous iteration t - 1, while we still allocated a small portion of the samples to represent those models that have small, even negligible, total weights. The rationale for the above choice of transition kernel is as follows. Even if after resampling at iteration t-1 all the samples represented model k, the other models would not become extinct at t. At iteration t = 0, instead of sampling uniformly on $[0, 0.5]^k$, which is the prior for frequency, we chose to sample from an IF $g^{(0)}(\mathbf{f}_k)$ to make sure that the samples reach the region of interest quickly. The function $g^{(0)}(\mathbf{f}_k)$ is defined by $g^{(0)}(\mathbf{f}_k) =$ $\mathcal{N}(\widehat{\mathbf{f}}_k, \mathbf{C}_k^{(0)})$, where $\widehat{\mathbf{f}}_k$ is a $k \times 1$ vector whose values are the frequencies of the largest k peaks of the periodogram of the data, and $\mathbf{C}_k^{(0)}$ is a $k \times k$ diagonal matrix whose diagonal elements, say $\mathbf{C}_k^{(0)}[i,i]$ represent a quarter of the width of the peak of the periodogram located at frequency f_i . We chose the following form of time-varying transition kernel $Q_{\theta_{k^{(i,t)}}}(\cdot, \theta) = \mathcal{N}\left(\hat{\theta}_{k^{(i,t-1)}}, \mathbf{C}_{k}^{(t)}\right)$, and let $\mathbf{C}_{k}^{(1)} = \mathbf{C}_{k}^{(0)}$.

When t > 1, if $k = \arg \max \sum_{i=1}^{N} w^{(i,t-1)} \Pi_k(k^{(i,t-1)})$, we have $\mathbf{C}_k^{(t)} = \frac{t-1}{t} \mathbf{C}_k^{(t-1)}$; if not, then $\mathbf{C}_k^{(t)} = \frac{t-1}{t} \mathbf{C}_k^{(t-1)} + \frac{1}{t} \mathbf{C}_k^{(1)}$. Of course, other choices of the transition kernels are possible, but we found that our choice resulted in good performance.

The averaged estimates of marginal distribution of model order k = 1, 2, 3 when SNR = 3dB and $l_e = 3$ is shown in Fig. 2, from which we observe that the estimates stabilized fairly quickly. Note, that the estimates of other models were all zeros and we did not show them in this figure. We also show in Fig. 3 the averaged estimates of the frequencies along with the true values of the frequencies. It is also clear that the estimates converge fast.



Fig. 2. Estimated p(k|y) vs number of iterations.



Fig. 3. Estimated f_1, f_2 vs number of iterations.

We also compared the sampling efficiencies of different choices of kernels. We first proposed a set of alternative kernels for model order that *does not* satisfy the requirements we suggested in III-B. The alternative kernels were defined as follows: $\bar{Q}_1(i, 4) = 0.4$, $\bar{Q}_1(i, j) = 0.15$, for $j \neq 4$; $\bar{Q}_2(i, 4) = 0.80$, $\bar{Q}_2(i, j) = 0.05$, for $j \neq 4$; $\bar{Q}_3(i, 4) =$ 0.92, $\bar{Q}_3(i, j) = 0.02$, for $j \neq 4$. It is clear that none of these kernels obey our first requirement, i.e., dominant diagonal entries. In Fig. 4, we calculated the averaged (over 100 realizations) entropy with respect to uniformity $H^{(t)}$ defined in (9) of three PMC implementations: 1) using $Q_d(\cdot, \cdot)$ for model order k and time-varying kernel for parameters θ ; 2) using $Q_d(\cdot, \cdot)$ for model order k and time-invariant kernel for parameters θ ; 3) using the alternative kernel $\bar{Q}_d(\cdot, \cdot)$ for model order k and time-varying kernel for parameters θ . We can see that by employing the time-varying kernel, the sampling efficiency increased steadily with each iteration. We can also see that the implementation with the alternative kernels $\bar{Q}_d(\cdot, \cdot)$ has lower efficiency than the use of the kernels $Q_d(\cdot, \cdot)$.



Fig. 4. Efficiency vs Iteration.

We implemented the RJMCMC algorithm in the following way. The algorithm was run for 30,000 iterations with 5,000 samples as burn-in period. We also used two different proposals for updating the frequency, one for local exploration and one for exploration of the "region of importance", see Section IV-A of [8]. Note that we chose the total number of samples generated by RJMCMC (30,000) to be large enough so that the estimated posterior can be stabilized. Also note, that this number was the same as the total number of samples used by the PMC (3,000 × 10) [16].

We compared these two algorithms under scenarios where the SNR = 3dB, 10dB and the spacing parameter $l_e = 1, 2, 4$. For each of the different scenarios, both algorithms were run for 100 realizations, and the comparison of the detection performance is shown in Table II. The entries in the table are the number of times a particular model k was chosen out of 100 realizations. It can be seen that the performance of the two algorithms was comparable. Note that we also present the performance of the PMC algorithm that employed only one transition kernel, and the algorithm with alternative kernels $\bar{Q}_d(\cdot, \cdot)$ defined above. We observed that for these two choices of kernels, the performance degraded when the sinusoids were closely spaced and the SNR was low. The last result supports our claim that using multiple kernels for model order is beneficial, it also supports our heuristics of intelligently choosing these predetermined kernels (as discussed in the end of section III-B).

We also observe that when the total number of samples was relatively small, the estimation performance of the RJMCMC deteriorated. For example, when the total number of samples was set to 10,000 and 2,000 samples were used for the burnin period, and we used 2000 samples per iteration with 5 iterations for PMC, the PMC outperformed the RJMCMC for low SNRs. Fig. 5 shows the estimation performance of the two algorithms under the above settings and when $l_e = 3$. It is clear that when the SNR is below 6dB, the PMC performs

TABLE II COMPARISON OF DETECTION PERFORMANCE.

Algorithm	SNR	l	$k \leq 1$	k = 2	$k \ge 3$
		1	0	100	0
PMC	3 dB	2	1	99	0
(multiple kernel)		4	77	23	0
		1	2	98	0
PMC	3 dB	2	9	90	1
(single kernel)		4	89	9	2
		1	0	98	2
PMC	3 dB	2	1	92	7
(alternative \bar{Q}_d)		4	89	11	0
		1	0	100	0
PMC	10 dB	2	0	100	0
(multiple kernel)		4	0	100	0
		1	1	99	0
PMC	10 dB	2	2	98	0
(single kernel)		4	5	95	0
		1	0	99	1
PMC	10 dB	2	0	96	4
(alternative \bar{Q}_d)		4	0	96	4
		1	0	99	1
	3dB	2	1	95	4
RJMCMC		4	81	19	0
		1	0	100	0
	10dB	2	0	98	2
RJMCMC		4	1	99	0



Fig. 5. MSE vs SNR.

better than the RJMCMC. In this figure we also plotted the CRLB [27].

D. Simulation Results for DOA

We used the following setup for the model: $d_y = 30, k = 2$, $\phi_1 = 20^\circ$, $\phi_1 = 45^\circ$, the number of sensors was M = 5and the amplitudes of the signals were fixed at $a_1 = 10$, $a_2 = 10$. In order to generate a spatially colored noise, we used a second order AR process with poles $0.9 \exp\{-j1.05\pi\}$

 TABLE III

 COMPARISON OF PERFORMANCE OF DETECTION

Algorithm	SNR	$k \leq 1$	k = 2	$k \ge 3$
	-1 dB	44	56	0
	0 dB	45	55	0
PMC	1 dB	37	63	0
4,000 samples	2 dB	35	65	0
	3 dB	14	82	4
	-1 dB	66	34	0
	0 dB	60	40	0
RJMCMC	1 dB	48	52	0
4,000 samples	2 dB	37	63	0
	3 dB	27	83	0
	-1 dB	48	52	0
	0 dB	49	51	0
RJMCMC	1 dB	40	60	0
15,000 samples	2 dB	30	69	1
	3 dB	15	85	0

and $0.9 \exp \{-j0.90\pi\}$ whose driving noise was a circularly complex white Gaussian process with identical σ^2 . The SNR was defined by $10 \log 10(\frac{a_1^2}{2\sigma^2})$. The hyperparameter ρ^2 was determined according to the criteria developed in Section V of [9] for different SNRs. As mentioned in [9], this setup was very difficult because the two sources were within a beamwidth of the receiver array, and the data size was very small.

For the PMC algorithm we used 1,000 samples in each iteration, and we ran the algorithm for a total of four iterations. We used three 5×5 matrices as transition kernels for the model order, which were identical to those used in the previous experiment. At iteration t = 0, we sampled the model order and the DOAs uniformly. Note that in this situation, we did not have a natural candidate for the initial IF, as in the previous example, so we sampled the parameters from their priors (which were $U[0, 2\pi]^k$). In the subsequent iterations, we chose to use the exact time-varying transition kernel as in the previous experiment, except that $\mathbf{C}_k^{(1)} = (0.1\pi)^2 \times \mathbf{I}_{k \times k}$. Note, that $\mathbf{C}_k^{(1)}$ was preselected, as suggested in [11], when no obvious candidate was available.

We ran the RJMCMC algorithm for two settings: 4,000 iterations with 1,000 samples used for the burn-in period, and 15,000 iterations with 5,000 samples for burn-in. In the first setting, the total number of samples generated by the RJMCMC (4,000) was the same as the total number of samples used by the PMC (1,000 \times 4). Notice that in this setting, our choice of total number of samples is significantly smaller than that used in [9], which corresponds to our second setting. The purpose here was to demonstrate the instability of the RJMCMC algorithm when the sample size was small. In both settings, we used two proposals for updating the DOA, one for local exploration and one for exploring the whole parameter space. See Section IV.A of [9] for details.

Both algorithms were run for 100 realizations, and the detection performance of the algorithms is compared in Table



Fig. 6. Estimated p(k|y) vs number of iterations obtained by the RJMCMC algorithm.



Fig. 7. Estimated p(k|y) vs number of iterations obtained by the PMC algorithm.

III. Again, each entry of the table represents the number of times a particular model is selected in 100 realizations. It can be seen that when the total number of samples was small and the SNR was low, the performance of the RJMCMC algorithm was worse than that of the PMC algorithm. In Fig. 6 we see a typical realization of the estimates of p(k = 1|y) and p(k = 2|y) produced by the RJMCMC algorithm with 4,000 samples and for SNR = -1dB. It is clear that the estimated marginalized posterior was not stabilized within the window of 4,000 samples, and the estimates based on these samples are not satisfactory. In Fig. 7 we see the averaged estimates of p(k=1|y) and p(k=2|y) for the PMC algorithm. Since the estimates of the other models were all zeros, we did not show them in the figure. It is clear that although there are only four iterations of the algorithm, the estimates improved with each iteration and they stabilized quickly.

E. Discussion

In the previous simulations we demonstrated the performance of the PMC by comparing it with that of the RJMCMC. We showed that the PMC had comparable performance with RJMCMC when the simulated sample size was large, and that the PMC outperformed the RJMCMC when the sample size was small. In both [8] and [9], the burn-in periods were determined in a heuristic fashion. We speculate that the choice of the updating proposals and the starting point of the algorithm heavily influenced the length of the burn-in period.

Another advantage of the PMC over the RJMCMC, as mentioned in Section I, is the potential for its parallel implementation. Although a thorough investigation of this topic is beyond the scope of this paper, we could easily identify several similarities between PMC and particle filtering, whose parallel implementation has been studied (see [28] for a recent development of this topic). We argue that RJMCMC is not suitable for parallel implementation mainly because each iteration of the algorithm produces a single sample and this sample is dependent on the sample produced in the previous iteration.

V. CONCLUSION

We have proposed a general algorithm to carry out the Bayesian computation required for selecting the MAP model order. We have studied the convergence result of the proposed algorithm and demonstrated its performance on two typical signal processing problems. Indeed the proposed algorithm is flexible enough to approximate the posterior distribution for both problems, and it is computationally more efficient than the popular RJMCMC algorithm.

One future direction of research is to investigate theoretically the convergence rate of the proposed algorithm, which could shed new light to its behavior. It would also facilitate researchers in proposing new structures for the transition kernels for improving its rate of convergence. We also believe that a thorough study of the parallel implementation of the PMC algorithm would be of great practical interest.

APPENDIX I **PROOF OF THEOREM 1**

The proof is an adaptation of the proofs for PMC algorithm in [15]. First note that we have the following equations for importance sampling.

$$E(f(x)) = \frac{\int f(x)w(x)q(x)dx}{\int w(x)q(x)dx}$$

$$\approx \frac{\frac{1}{N}\sum_{i=1}^{N}\tilde{w}^{(i)}f(x^{(i)})}{\frac{1}{N}\sum_{i=1}^{N}\tilde{w}^{(i)}}$$
(40)

$$= \sum_{i=1}^{N} w^{(i)}f(x^{(i)})$$

where q(x) is the sampling distribution and $\tilde{w}^{(i)}$ is unnormalized weight. In the following, we will set out to prove that

$$\lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} \tilde{w}^{(i,t)} f(\mathbf{x}^{(i,t)}) = \int f(\mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$
(41)

which is the numerator of (40), and the convergence of the normalization constant will be automatically established when we plug in f(x) = 1 ($\lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} \tilde{w}^{(i,t)} = 1$) [15], and (15) would follow.

We then state a Lemma. Let \rightarrow_p denote convergence in probability.

Lemma 1: Let (Ω, A) be a measurable space. Assume that the following conditions hold:

- A sequence $\{\xi^{(i)}\}_{i=1}^{N}$ is independent given G_N , where G_N is a σ algebra in A. $\left\{\sum_{i=1}^{N} E[\xi^{(i)}|G_N]\right\}$ is bounded in probability, e.g., i)
- ii)

$$\lim_{a \to \infty} \sup_{N>1} p\left(\sum_{i=1}^{N} E[\xi^{(i)}|G_N] \ge a\right) = 0 \qquad (42)$$

iii) $\forall b > 0, \sum_{i=1}^{N} E\left[\xi^{(i)} \Pi\{|\xi^{(i)}| > b\}|G_N\right] \to_p 0.$ Then.

$$\sum_{i=1}^{N} (\xi^{(i)} - E(\xi^{(i)}|G_N)) \to_p 0.$$
(43)

Proof: The proof of the lemma involves the use of Doob Maximal Inequality. We refer the interested readers to [29] for the details of the proof.

We verify the conditions on Lemma 1 to prove the convergence of the PMC algorithm.

We first state our assumptions. We assume that the function f(x) is absolutely integrable under p(x|y), which is to say

$$\int |f(x)|p(x|y)dx < \infty.$$
(44)

Notice, that under this assumption, f(x) need not be bounded and that if for a set $M = \{x: |f(x)| = \infty\}$, we must have for all $x \in M$

$$|f(x)|p(M) = C < \infty$$
, and $|f(x)|p(x|y) = C < \infty$ (45)

where C is a constant. Otherwise the absolute integrability of f is contradicted. We also assume that the kernels $Q_d(\cdot, \cdot)$ and $Q_{\theta_k}(\cdot, \cdot)$ do not take value of 0 (or that $\frac{1}{Q_d(\cdot, \cdot)}$ and $\frac{1}{Q_{\theta_k}(\cdot, \cdot)}$ are both finite).

Let

$$\xi^{(i)} = \frac{1}{N} \tilde{w}^{(i,t)} f(\mathbf{x}^{(i,t)}) = \frac{1}{N} \frac{p(\mathbf{x}^{(i,t)} | \mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_{d^{(i,t)}}(k^{*(i,t-1)}, k^{(i,t)}) Q_{\theta_{k^{(i,t)}}}(\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})}$$
(46)

where $\mathbf{x}^{(i,t)}$ is the vector $[k^{(i,t)} \quad \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)^{\top}}]$. We let $G_N =$ $\sigma\{\{\mathbf{x}^{*(i,t-1)}\}_{i=1}^{N}, \{\alpha_{d}^{(t)}\}_{d=1}^{D}\}, \text{ which is, essentially, a } \sigma \text{ alge$ bra generated by the sequence $\{\mathbf{x}^{*(i,t-1)}\}_{i=1}^N$ and $\{\alpha_d^{(t)}\}_{d=1}^D$.

Notice that in the proposed algorithm the step t = 0, is just conventional importance sampling, so the proof for this case comes from the law of large numbers of importance sampling. Based on this observation, we check all the three conditions in Lemma 1, for $t \ge 1$. We use induction and that the convergence in (15) is established for t - 1.

1) At iteration t, the $k^{(i,t)}$'s are iid and drawn from $Q_{d^{(i,t)}}(k^{*(i,t-1)}, k^{(i,t)})$, and the $d^{(i,t)}$'s are iid and generated from $\mathcal{M}(\alpha_1^{(t)}, ..., \alpha_D^{(t)})$. Thus, we have that $\{k^{(i,t)}\}_{i=1}^{N}$ are independent from each other conditionally on $k^{*(i,t-1)}$ and $\alpha_d^{(t)}$. We also have that each $\boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)}$ is drawn from $Q_{\boldsymbol{\theta}_{k^{(i,t)}}}(\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})$. Then, the

 $\hat{\theta}_{k^{(i,t)}}^{(t-1)}$ is the estimate of θ for the model order $k^{(i,t)}$ and is obtained from $\{\mathbf{x}^{(i,t-1)}\}_{i=1}^{N}$ at iteration t-1). 2) Since

$$E\left[\xi^{(i)}|G_{N}\right] = E\left[\frac{1}{N}\tilde{w}^{(i,t)}f(\mathbf{x}^{(i,t)})|G_{N}\right] = \frac{1}{N}\int \frac{p(\mathbf{x}^{(i,t)}|\mathbf{y})f(\mathbf{x}^{(i,t)})}{Q_{d^{(i,t)}}(k^{*(i,t-1)},k^{(i,t)})Q_{\theta_{k^{(i,t)}}}(\widehat{\theta}_{k^{(i,t)}}^{(t-1)},\theta_{k^{(i,t)}}^{(i,t)})} > Q_{d^{(i,t)}}(k^{*(i,t-1)},k^{(i,t)})Q_{\theta_{k^{(i,t)}}}(\widehat{\theta}_{k^{(i,t)}}^{(t-1)},\theta_{k^{(i,t)}}^{(i,t)})d\mathbf{x}^{(i,t)}} = \frac{1}{N}\int p(\mathbf{x}^{(i,t)}|\mathbf{y})f(\mathbf{x}^{(i,t)})d\mathbf{x}^{(i,t)}.$$
(47)

The integral in the above result is bounded for all i and we have

$$\sum_{i=1}^{N} E[\xi^{(i)}|G_N] = \frac{1}{N} \sum_{i=1}^{N} \int p(\mathbf{x}^{(i,t)}|\mathbf{y}) f(\mathbf{x}^{(i,t)}) d\mathbf{x}^{(i,t)} < \infty.$$

The last inequality comes from the assumption that $f(\cdot)$ is absolutely integrable w.r.t. $p(\mathbf{x}|\mathbf{y})$.

3) From the condition (iii) of the lemma, we have the following:

$$\sum_{i=1}^{N} E[\xi^{(i)} \Pi\{|\xi^{(i)}| > \delta\}|G_N]$$

independence of $\theta_{k^{(i,t)}}^{(i,t)}$ conditional on $\{\mathbf{x}^{*(i,t-1)}\}_{i=1}^N$ of $\mathbf{x}^{(i,t)}$, d, $k^{*(i,t-1)}$ and $\hat{\theta}_{k^{(i,t)}}^{(t-1)}$, respectively. Then the above comes from the independence of $k^{(i,t)}$ (notice, that equation becomes

$$\frac{1}{N}\sum_{i=1}^{N}\sum_{d=1}^{N}\alpha_{d}^{(t)}\int F(\cdot)\Pi\{|M(\cdot)|>\delta\}d\mathbf{x}^{(i,t)}.$$
(49)

Notice that $\int F(\cdot)\Pi\{|M(\cdot)| > \delta\}d\mathbf{x}^{(i,t)}$ is a function of $k^{*(i,t-1)}$, d, and $\{\widehat{\theta}_{k}^{(t-1)}\}_{k=1}^{K}$. Evidently, we have the follow- \times ing inequality:

$$\int F(\cdot) \Pi\{|M(\cdot)| > \delta\} d\mathbf{x}^{(i,t)}$$

$$\leq \int |F(\cdot)| \Pi\{|M(\cdot)| > \delta\} d\mathbf{x}^{(i,t)}$$

$$\leq \int |F(\cdot)| d\mathbf{x}^{(i,t)} = \int |p(\mathbf{x}^{(i,t)}|\mathbf{y}) f(\mathbf{x}^{(i,t)})| d\mathbf{x}^{(i,t)} (50)$$

$$< \infty.$$

Consequently, we have that $\int F(\cdot) \Pi\{|M(\cdot)| > \delta\} d\mathbf{x}^{(i,t)}$ is a function of $k^{*(i-1,t)}$, and it takes a finite value (it is also, as mentioned above, a function of d and $\{\widehat{\theta}_k^{(t-1)}\}_{k=1}^K$). Evidently, the above function is integrable under $p(k|\mathbf{y})$. Let us denote the integral $\int F(\cdot) \Pi\{|M(\cdot)| > \delta\} d\mathbf{x}^{(i,t)}$ by $h(k^{*(i-1,t)})$. Recall that $k^{*(i,t-1)}$ is a resampled particle, so the weights are normalized. By induction on equation (15), we have

$$\begin{split} &= \frac{1}{N} \sum_{i=1}^{N} E(\tilde{w}^{(i,t)} f(\mathbf{x}^{(i,t)}) \Pi\{|\frac{1}{N} \tilde{w}^{(i,t)} f(\mathbf{x}^{(i,t)})| > \delta\}|G_{N}) \\ &= \frac{1}{N} \sum_{i=1}^{N} E\Big(\frac{p(\mathbf{x}^{(i,t)} | \mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_{d^{(i,t)}} (k^{*(i,t-1)}, k^{(i,t)}) Q_{\theta_{k^{(i,t)}}} (\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})} \times \\ &= \lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} \int F(\cdot) \Pi\{|M(\cdot)| > \delta\} d\mathbf{x}^{(i,t)} \\ &= \lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} h(\tilde{k}^{(i-1,t)}) \\ \Pi\left\{\frac{1}{N} \left| \frac{p(\mathbf{x}^{(i,t)} | \mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_{d^{(i,t)}} (k^{*(i,t-1)}, k^{(i,t)}) Q_{\theta_{k^{(i,t)}}} (\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})} \right| > \delta\right\} |G_{N}\Big)_{:p} \lim_{N \to \infty} \int h(k) p(k|\mathbf{y}) dk \\ & \xrightarrow{K} \end{split}$$

Notice, that the expectation is taken on the random variable $d^{(i,t)}$ and $\mathbf{x}^{(i,t)}$, respectively. Since

$$E_{x,d}(f(x,d)) = E_d(E_{x|d}(f(x,d)))$$
(48)

the above long equation becomes, as in equation (47),

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{d=1}^{N} \alpha_{d}^{(t)} \int p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} | \mathbf{y}) f(\mathbf{x}^{(i,t)}) \times \\ \prod \left\{ \left| \frac{1}{N} \frac{p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} | \mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_{d}(k^{*(i,t-1)}, k^{(i,t)}) Q_{\boldsymbol{\theta}_{k^{(i,t)}}} (\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})} \right| > \delta \right\} d\mathbf{x}^{(i,t)}$$

We simplify the equation by denoting $p(\mathbf{x}^{(i,t)}|\mathbf{y})f(\mathbf{x}^{(i,t)})$ as $F(\cdot)$ and $\frac{1}{N} \frac{p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)}|\mathbf{y})f(\mathbf{x}^{(i,t)})}{Q_d(k^{*(i,t-1)}, k^{(i,t)})Q_{\boldsymbol{\theta}_{k^{(i,t)}}}(\hat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})}$ as $M(\cdot)$.

Notice that $F(\cdot)$ is a function of $\mathbf{x}^{(i,t)}$, and $M(\cdot)$ is a function

 $= \lim_{N \to \infty} \sum_{k=1} h(k) p(k|\mathbf{y})$ $= \lim_{N \to \infty} \sum_{k=1}^{n} \int F(\cdot) \ \Pi\{|M(\cdot)| > \delta\} d\mathbf{x}^{(i,t)} p(k|\mathbf{y})$ $\leq \lim_{N \to \infty} \int |F(\cdot)| \sum_{l=1}^{K} \Pi\{|M(\cdot)| > \delta\} p(k|\mathbf{y}) d\mathbf{x}^{(i,t)}$ $= \lim_{N \to \infty} \int_{E_1} |F(\cdot)| \sum_{i=1}^{K} \Pi\{|M(\cdot)| > \delta\} p(k|\mathbf{y}) d\mathbf{x}^{(i,t)}$

$$+ \int_{E_2} |F(\cdot)| \sum_{k=1}^{K} \Pi\{|M(\cdot)| > \delta\} p(k|\mathbf{y}) d\mathbf{x}^{(i,t)}$$
(51)

where $E_1 \triangleq \{\mathbf{x}^{(i,t)} : |f(\mathbf{x}^{(i,t)})| < \infty\}$, and $E_2 \triangleq \{\mathbf{x}^{(i,t)} : |f(\mathbf{x}^{(i,t)})| = \infty\}$ and we have $E_1 \cap E_2 = \emptyset$. As N goes to

infinity, it is clear that the first integral goes to 0, as follows:

$$\lim_{N \to \infty} \int_{E_1} F(\cdot) \sum_{k=1}^{K} \Pi\{|M(\cdot)| > \delta\} p(k|\mathbf{y}) d\mathbf{x}^{(i,t)}$$

$$= \lim_{N \to \infty} \int_{E_1} p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} |\mathbf{y}) f(\mathbf{x}^{(i,t)}) \times$$

$$P_k \left(\left| \frac{1}{N} \frac{p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} |\mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_d(k, k^{(i,t)}) Q_{\boldsymbol{\theta}_{k^{(i,t)}}} (\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})} \right| > \delta \right) d\mathbf{x}^{(i,t)}$$

We have that both $Q_d(k, k^{(i,t)})$ and $Q_{\theta_k(i,t)}(\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})$ are different from 0 and that $f(\cdot)$ is finite. Thus, $\left|\frac{p(k^{(i,t)}, \theta_{k^{(i,t)}}^{(i,t)} | \mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_d(k, k^{(i,t)}) Q_{\theta_k(i,t)}(\theta_{k^{(i,t)}}^{(t-1)}, \theta_{k^{(i,t)}}^{(i,t)})}\right|$ is finite, and as N goes to

$$\begin{aligned} &|Q_{d}(k,k^{(i,t)})Q_{\theta_{k}(i,t)}(\theta_{k}^{*}(i,t),\theta_{k}^{*}(i,t))| \\ &\text{infinity, } P_{k}\left(\left|\frac{1}{N}\frac{p(k^{(i,t)},\boldsymbol{\theta}_{k}^{(i,t)}|\mathbf{y})f(\mathbf{x}^{(i,t)})}{Q_{d}(k^{*},k^{(i,t)})Q_{\theta_{k}(i,t)}(\widehat{\boldsymbol{\theta}}_{k}^{(t-1)},\boldsymbol{\theta}_{k}^{(i,t)})}\right| > \delta\right) \text{ goes } \\ &\text{to zero, for all } \delta \end{aligned}$$

to zero, for all δ .

We now examine the second integral. We use the equation (45).

$$0 \leq \lim_{N \to \infty} \int_{E_2} |F(\cdot)| \sum_{k=1}^{K} \Pi\{|M(\cdot)| > \delta\} p(k|\mathbf{y}) d\mathbf{x}^{(i,t)} \\ = \lim_{N \to \infty} \int_{E_2} |p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} |\mathbf{y}) f(\mathbf{x}^{(i,t)})| \times \\ P_k \left(\left| \frac{1}{N} \frac{p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} |\mathbf{y}) f(\mathbf{x}^{(i,t)})}{Q_d(k, k^{(i,t)}) Q_{\boldsymbol{\theta}_{k^{(i,t)}}} (\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)})} \right| > \delta \right) d\mathbf{x}^{(i,t)} \\ \leq \lim_{N \to \infty} \int_{E_2} |p(k^{(i,t)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)} |\mathbf{y}) f(\mathbf{x}^{(i,t)})| \times \\ P_k \left(\left| \frac{1}{N} \frac{C}{Q_{1}(k, k^{(i,t)}) Q_{2}} - (\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)}, \boldsymbol{\theta}_{k^{(i,t)}}^{(i,t)}) \right| > \delta \right) d\mathbf{x}^{(i,t)} \\ \leq \frac{1}{N} \sum_{Q_{2}(k, k^{(i,t)}) Q_{2}} \left(\frac{\widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(t-1)} - \widehat{\boldsymbol{\theta}}_{k^{(i,t)}}^{(i,t)}} \right) \right) d\mathbf{x}^{(i,t)}$$

$$\bigvee_{k} Q_d(k, k^{(i,t)}) Q_{\theta_{k}(i,t)} (\theta_{k^{(i,t)}}^{*}, \theta_{k^{(i,t)}}^{(t,t)})^{*} = 0.$$

The last equation is true because, again, $Q_d(k, k^{(i,t)})$ and $Q_{\theta_{k^{(i,t)}}}(\widehat{\theta}_{k^{(i,t)}}^{(t-1)}, \theta_{k^{(i,t)}}^{(i,t)})$ both do not equal to 0. The second to the last inequality is valid because of the following fact. If $f(x) \geq g(x)$ pointwise, then we have

$$p(\{x : f(x) > \delta\}) \ge p(\{x : g(x) > \delta\}), \quad (52)$$

or, equivalently

$$E_{\delta 1} \triangleq \{x : f(x) > \delta\} \supseteq E_{\delta 2} \triangleq \{x : g(x) > \delta\}.$$
 (53)

Finally, we have the following convergence:

$$\frac{1}{N}\sum_{i=1}^{N}F(\cdot) \,\Pi\{|M(\cdot)| > \delta\}d\mathbf{x}^{(i,t)} \to_{p} 0 \tag{54}$$

and, we have

$$\frac{1}{N}\sum_{i=1}^{N}\sum_{d=1}^{D}\alpha_{d}^{(t)}\int F(\cdot) \,\,\Pi\{|M(\cdot)| > \delta\}d\mathbf{x}^{(i,t)} \to_{p} 0 \quad (55)$$

Thus, we verified the third condition, and the conclusion follows, namely,

$$\lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} \tilde{w}^{(i,t)} f(\mathbf{x}^{(i,t)}) = \int f(\mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

REFERENCES

- P. Stoica and Y. Selen, "Model order selection: A review of information criterion rules," *IEEE Signal Processing Magazine*, vol. 21, no. 4, pp. 36–47, 2004.
- [2] P. M. Djurić, "Simultaneous detection and frequency estimation of sinusoidal signals," in the Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, Minneapolis, MN, 1993.
- [3] P. M. Djurić, "A model selection rule for sinusoids in white Gaussian noise," *IEEE Transactions on Signal Processing*, vol. 44, pp. 1744–1751, July 1996.
- [4] P. M. Djurić, "Asymptotic map criteria for model selection," *IEEE Transactions on Signal Processing*, vol. 46, pp. 2726–2735, Oct 1998.
- [5] C.-M Cho and P. M. Djurić, "Bayesian detection and estimation of cisoids in colored noise," *IEEE Transactions on Signal Processing*, vol. 43, pp. 2943–2951, 1996.
- [6] P. J. Green, "Reversible jump Markov chain Monte Carlo computation and Bayesian model determination," *Biometrika*, vol. 4, pp. 711–732, 1995.
- [7] C. Andrieu, P. M. Djurić, and A. Doucet, "Model selection by Markov chain Monte Carlo computations," *Signal Processing*, vol. 81, pp. 19–37, 2001.
- [8] C. Andrieu and A. Doucet, "Joint Bayesian model selection and estimation of noise sinusoids via reversible jump MCMC," *IEEE Transaction on Signal Processing*, vol. 47, no. 10, pp. 2667–2676, 1999.
- [9] J. R. Larocque and James P. Reilly, "Reversible jump MCMC for joint detection and estimation of sources in colored noise," *IEEE Transactions (i,t) on Signal Processing*, vol. 50, pp. 231–240, 2002.
- [10] R. Iba, "Population-based Monte Carlo algorithms," Journal of Computational and Graphical Statistics, vol. 7, pp. 175–193, 2000.
- [11] O. Cappé, J. M. Marin, and C. P. Robert, "Population Monte Carlo," *Journal of Computational and Graphical Statistics*, vol. 13, no. 4, pp. 907–929, 2003.
- [12] D. Rubin, "Comment on 'The calculation of posterior distributions by data augmentation' by M. A. Tanner and W. H. J. Wong," *Journal of the American Statistical Association*, vol. 82, pp. 543, 1987.
- [13] C. P. Robert and G. Casella, *Monte Carlo Statistical Methods*, Springer, New York, 2005.
- [14] R. Douc, A. Guillin, J.-M. Marin, and C. P. Robert, "Minimum variance importance sampling via Population Monte Carlo," *ESAIM: Probability* and Statistics, vol. 11, pp. 427–447, 2007.
- [15] R. Douc, A. Guillin, J.-M. Marin, and C. P. Robert, "Convergence of adaptive mixtures of importance sampling schemes," *The Annals of Statistics*, vol. 35, 2007.
- [16] G. Celeux, J.-M. Marin, and C. P. Robert, "Iterated importance sampling in missing data problems," *Computational Statistics & Data Analysis*, vol. 50, pp. 3386–3404, 2006.
- [17] M. F. Bugallo, M. Hong, and P. M. Djurić, "Marginalized Population Monte Carlo," in the Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, Taipei, Taiwan, 2009.
- [18] A. Doucet, N. de Freitas, and N. Gordon, Eds., Sequential Monte Carlo Methods in Practice, Springer, New York, 2001.
- [19] M. Shafi P. J. Smith and H. Gao, "Quick simulation: A review of importance sampling techniques in communications systems," *Journal On Selected Areas of Communications*, vol. 15, pp. 397–412, 1997.
- [20] M. West, "Approximating posterior distribution by mixtures," J. R. Statist. Soc. B, vol. 55, no. 2, pp. 409–422, 1993.
- [21] O. Cappé, A. Guillin, J.-M. Marin, and C. P. Robert, "Population Monte Carlo for ion channel restoration," *Journal of Computational* & *Graphical Statistics*, 2004.
- [22] P. M. Djurić, J. H. Kotecha, J. Zhang, Y. Huang, T. Ghirmai, M. F. Bugallo, and J. Miguez, "Particle filtering," *IEEE Signal Processing Magazine*, vol. 20, no. 5, pp. 19–38, 2003.
- [23] D. B. Rubin, *Bayesian Statistics 3*, chapter Using the SIR algorithm to simulate posterior distributions, pp. 395–402 (with discussion), Oxford: University Press, 1988.
- [24] J. S. Liu, Monte Carlo Strategies in Scientific Computing, Springer, New York, 2001.

because $\alpha_d^{(t)} <= 1$ and $\sum_{d=1}^D \alpha_d^{(t)} = 1$.

- [25] K.M.Wong, J.P. Reilly, Q. Wu, and S.Qiao, "Estimation of the directions of arrival of signals in unknown correlated noise, part i: The map approach and its implementation," *IEEE Transactions on Signal* Processing, vol. 40, pp. 2007-2017, 1992.
- [26] G. E. P. Box an G. C. Tiao, Bayesian Inference in Statistical Analysis, Wiley, 1973.
- [27] S. M. Kay, Modern Spectral Estimation, Prentice Hall, Englewood Cliffs, NJ, 1988.
- [28] S. Hong and P. M. Djurić, "High-throughput scalable parallel resampling [20] D. Hong and F. M. Djane, "Ingr-throughput scatable parallel resampling mechanism for effective redistribution of particles," *IEEE Transactions on Signal Processing*, vol. 54, no. 3, pp. 1144–1155, 2006.
 [29] R. Douc and E.Moulines, "Limit theorems for properly weighted samples with applications to sequential Monte Carlo," *Technical report*, *TSL Telecom Paris*, 2005.
- TSI, Telecom Paris, 2005.