

A PARTICLE FILTER FOR MODEL BASED AUDIO SOURCE SEPARATION

C. Andrieu, S.J. Godsill

Signal Processing Lab., University of Cambridge
Department of Engineering, Trumpington Street
CB2 1PZ Cambridge, UK
Email: {ca226, sjg}@eng.cam.ac.uk

ABSTRACT

In this paper we present an original modelling of the source separation problem that takes into account all the non-stationarities of the underlying processes. The estimation of the sources then reduces to that of a filtering/fixed-lag smoothing algorithm, for which we propose an efficient numerical solution, relying on particle filter techniques.

1 Introduction

The problem of separating convolutively mixed sources is of interest to many applications such as speech enhancement, crosstalk removal in multichannel communications and multipath channel identification. In this paper we consider the problem of separating audio signals modelled as autoregressive processes. The desired sources are clearly non-stationary as the vocal tract is continually changing, sometimes slowly, sometimes rapidly and the coupling filters may be time varying as the propagation conditions in the medium or more simply the sources/sensors geometry change over time. Here we explicitly model the non-stationarities of both the sources and the propagation medium, allowing us to take into account different stationarity time scales: speech will typically be stationary over periods of 20-40 milliseconds, whereas the transmission channel is expected to be stationary over longer periods. Taking into account these non-stationarities clearly removes some of the problems associated with framed based algorithms that assume piecewise stationarity and result in a delayless algorithm.

The modelling we adopt for the system facilitates a state space representation, and the problem of estimating the sources from observed mixtures then reduces to that of a filtering/fixed-lag smoothing estimation problem. The filtering and smoothing problems, that is the recursive estimation of the sources conditional upon the currently available data, require the evaluation of integrals that do not admit closed-form analytic solutions and approximate methods must be employed. Classical methods to obtain approximations to the desired distributions include analytical approximations, such as the extended Kalman filter [1] the Gaussian sum filter [3], and deterministic numerical integration techniques (see *e.g.* [5]). The extended Kalman filter and Gaussian sum filter are computationally cheap, but fail in difficult circumstances.

Instead of hand-crafting such algorithms, we propose here the use of an adaptive stochastic grid approximation, which introduces a so-called system of particles. These particles evolve randomly in time in correlation with one another, and either give birth to offspring particles or die according to their ability to represent the

different zones of interest of the state space dictated by the observation process and the dynamics of the underlying system. This class of Monte Carlo methods was introduced in the late 60's, but were overlooked mostly because of their computational complexity, see [9] for a review. The advantage of these methods is that they get around the problems faced by classical methods: the adaptivity of the stochastic grid results in an accurate representation of the density of interest at a rate independent of the dimension of the problem [6].

Particle techniques are here applied to obtain filtered and fixed-lag smoothed estimates of the non-stationary sources modelled as time varying autoregressive processes, observed in additive white Gaussian noise. The simulation-based algorithms developed here are not just a straightforward application of the basic methods, but are designed to make efficient use of the analytical structure of the model. At each iteration the algorithm has a computational complexity that is linear in the number of particles, and can easily be implemented on parallel computers, thus facilitating near real-time processing. It is also shown how an efficient fixed-lag smoothing algorithm may be obtained by combining the filtering algorithm with Markov chain Monte Carlo (MCMC) methods (see [12] for an introduction to MCMC methods). Note that the power and flexibility of these techniques can allow for far more complex models than those considered here.

The remainder of the paper is organized as follows. The model specification and estimation objectives are stated in Section 2. In Section 3 sequential simulation-based methods are developed to solve the filtering and fixed-lag smoothing problem and determine the model adequacy. In this section we first introduce Monte Carlo basics, secondly show how to take advantage of the structure of the model and reduce dramatically the dimensionality of the problem, thirdly explain why a selection scheme for the particles is required and why diversity must be introduced, here using MCMC steps. Section 4 presents and discusses simulation results.

2 Model of the data

The problem addressed is the problem of source separation, the n sources being modelled as autoregressive processes, from which we have at each time t , m observations which are convolutive mixtures of the n sources.

2.1 Model for the sources

Source i can be modelled for $t = 1, \dots$ as:

$$\mathbf{s}_{i,t} = \mathbf{a}_{i,t}^T \mathbf{s}_{i,t-1:t-p_i} + \sigma_{v,i,t} v_{i,t}^s \quad (1)$$

C. Andrieu is sponsored by AT&T Lab., Cambridge UK.

2.4 Objectives

Our aim is, given the number of sources m , p_i and $l_{i,j}$ to estimate sequentially the sources $(\mathbf{x}_t)_{t=1,\dots}$ and their parameters $\boldsymbol{\theta}_t \triangleq \{\mathbf{a}_t, \mathbf{h}_t, \boldsymbol{\sigma}_{t,1:m,v}^2, \boldsymbol{\sigma}_{t,1:n,w}^2\}$ from the observations $y_{j,t}$. More precisely, in the framework of Bayesian estimation, one is interested in the recursive, in time, estimation of posterior distributions of the type $p(d\boldsymbol{\theta}_t, d\mathbf{x}_t | \mathbf{y}_{1:t+L})$: when $L = 0$ this corresponds to the filtering distribution and when $L > 0$ this corresponds to the fixed-lag smoothing distribution. This is a very complex problem that does not admit any analytical solution, and one has to resort to numerical methods. In the next section we develop such a numerical method based on Monte Carlo simulation. Subsequently we will use the following notation: $\boldsymbol{\alpha}_t \triangleq \{\mathbf{a}_t, \mathbf{h}_t\}$, $\boldsymbol{\beta}_t \triangleq \{\boldsymbol{\sigma}_{t,1:m,v}^2, \boldsymbol{\sigma}_{t,1:n,w}^2\}$ and $\boldsymbol{\gamma}_t \triangleq \{\mathbf{x}_t, \boldsymbol{\sigma}_{t,1:m,v}^2, \boldsymbol{\sigma}_{t,1:n,w}^2\}$.

3 A Simulation-Based Optimal Filter/Fixed-lag Smoother

This section develops a simulation-based optimal filter/fixed-lag smoother to obtain filtered/fixed-lag smoothed estimates of the unobserved sources and their parameters of the type

$$I_L(f_t) \triangleq \int f_t(\boldsymbol{\theta}_t, \mathbf{x}_t) p(d\boldsymbol{\theta}_t, d\mathbf{x}_t | \mathbf{y}_{1:t+L}) \quad (18)$$

The standard Bayesian importance sampling method is first described, and then we show how it is possible to take advantage of the analytical structure of the model by integrating out the parameters \mathbf{a}_t and \mathbf{h}_t which can be high dimensional, using Kalman filter related algorithms. This leads to an elegant and efficient algorithm for which the only tracked parameters are the sources and the noise variances. Then a sequential version of Bayesian importance sampling for optimal filtering is presented, and it is shown why it is necessary to introduce selection as well as diversity in the process. Finally, a Monte Carlo filter/fixed-lag smoother for our problem is described.

3.1 Monte Carlo Simulation for Optimal Estimation

For any f_t it will subsequently be assumed that $|I_L(f_t)| < +\infty$. Suppose that it is possible to sample N *i.i.d.* samples, called particles, $(\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)})$ according to $p(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$. Then an empirical estimate of this distribution is given by

$$\begin{aligned} \hat{p}_N(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L}) \\ = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)}}(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L}) \end{aligned} \quad (19)$$

so that a Monte Carlo approximation of the marginal distribution $p(d\mathbf{x}_t, d\boldsymbol{\theta}_t | \mathbf{y}_{1:t+L})$ follows as

$$\hat{p}_N(d\mathbf{x}_t, d\boldsymbol{\theta}_t | \mathbf{y}_{1:t+L}) = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_t^{(i)}, \boldsymbol{\theta}_t^{(i)}}(d\mathbf{x}_t, d\boldsymbol{\theta}_t) \quad (20)$$

Using this distribution, an estimate of $I_L(f_t)$ for any f_t may be obtained as

$$\begin{aligned} \hat{I}_{L,N}(f_t) &= \int f_t(\mathbf{x}_t, \boldsymbol{\theta}_t) \hat{p}_N(d\mathbf{x}_t, d\boldsymbol{\theta}_t | \mathbf{y}_{1:t+L}) \\ &= \frac{1}{N} \sum_{i=1}^N f_t(\mathbf{x}_t^{(i)}, \boldsymbol{\theta}_t^{(i)}) \end{aligned} \quad (21)$$

This estimate is unbiased and from the strong law of large numbers, $\hat{I}_{L,N}(f_t) \xrightarrow[N \rightarrow +\infty]{a.s.} I_L(f_t)$. Under additional assumptions the estimates satisfy a central limit theorem. The advantage of the Monte Carlo method is clear. It is easy to estimate $I_L(f_t)$ for any f_t , and the rate of convergence of this estimate does not depend on t or the dimension of the state space, but only on the number

of particles N and the characteristics of the function f_t . Unfortunately, it is not possible to sample directly from the distribution $p(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$ at any t , and alternative strategies need to be investigated.

One solution to estimate $p(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$ and $I_L(f_t)$ is the well-known Bayesian importance sampling method [9]. This method assumes the existence of an arbitrary importance distribution $\pi(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$ which is easily simulated from, and whose support contains that of $p(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$. Using this distribution $I_L(f_t)$ may be expressed as

$$I_L(f_t) = \frac{\int \pi(d\mathbf{x}_{1:t+L}, d\boldsymbol{\theta}_{1:t+L} | \mathbf{y}_{1:t+L}) f_t(\mathbf{x}_t, \boldsymbol{\theta}_t) w(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L})}{\int \pi(d\mathbf{x}_{1:t+L}, d\boldsymbol{\theta}_{1:t+L} | \mathbf{y}_{1:t+L}) w(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L})}, \quad (22)$$

where the importance weight $w(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L})$ is given by

$$w(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L}) \propto \frac{p(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})}{\pi(\mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})}. \quad (23)$$

The importance weight can normally only be evaluated up to a constant of proportionality, since, following from Bayes' rule,

$$\begin{aligned} p(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L}) \\ = \frac{p(\mathbf{y}_{1:t+L} | \mathbf{x}_{0:t+L}, \boldsymbol{\theta}_{0:t+L}) p(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L})}{p(\mathbf{y}_{1:t+L})}, \end{aligned} \quad (24)$$

where the normalizing constant $p(\mathbf{y}_{1:t+L})$ can typically not be expressed in closed-form.

If N *i.i.d.* samples $(\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)})$ can be simulated according to a distribution $\pi(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$, a Monte Carlo estimate of $I_L(f_t)$ in (22) may be obtained as

$$\begin{aligned} \hat{I}_{L,N}^1(f_t) &= \frac{\sum_{i=1}^N f_t(\mathbf{x}_t^{(i)}, \boldsymbol{\theta}_t^{(i)}) w(\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)})}{\sum_{i=1}^N w(\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)})} \\ &= \sum_{i=1}^N \bar{w}_{0:t+L}^{(i)} f_t(\mathbf{x}_t^{(i)}, \boldsymbol{\theta}_t^{(i)}), \end{aligned} \quad (25)$$

where the normalized importance weights are given by

$$\bar{w}_{0:t+L}^{(i)} = \frac{w(\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)})}{\sum_{j=1}^N w(\mathbf{x}_{0:t+L}^{(j)}, \boldsymbol{\theta}_{0:t+L}^{(j)})}. \quad (26)$$

This method is equivalent to a point mass approximation of the target distribution of the form

$$\begin{aligned} \hat{p}_N(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L}) \\ = \sum_{i=1}^N \bar{w}_{0:t+L}^{(i)} \delta_{\mathbf{x}_{0:t+L}^{(i)}, \boldsymbol{\theta}_{0:t+L}^{(i)}}(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L}), \end{aligned} \quad (27)$$

The perfect simulation case, when $\pi(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L}) = p(d\mathbf{x}_{0:t+L}, d\boldsymbol{\theta}_{0:t+L} | \mathbf{y}_{1:t+L})$, corresponds to $\bar{w}_{0:t+L}^{(i)} = N^{-1}$, $i = 1, \dots, N$. In practice, the importance distribution will be chosen to be as close as possible to the target distribution in a given sense. For finite N , $\hat{I}_{L,N}^1(f_t)$ is biased, since it involves a ratio of estimates, but asymptotically, according to the strong law of large numbers, $\hat{I}_{L,N}^1(f_t) \xrightarrow[N \rightarrow +\infty]{a.s.} I_L(f_t)$. Under additional assumptions a central limit theorem also holds [9].

3.2 Analytical integrations

It is possible to reduce the estimation of $p(d\mathbf{x}_t, d\boldsymbol{\theta}_t | \mathbf{y}_{1:t+L})$ and $I_L(f_t)$ to one of sampling from $p(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$, where we recall that $\boldsymbol{\gamma}_t \triangleq \{\mathbf{x}_t, \boldsymbol{\sigma}_{t,1:m,v}^2, \boldsymbol{\sigma}_{t,1:n,w}^2\}$. Indeed,

$$p(d\boldsymbol{\alpha}_t, d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L}) = p(d\boldsymbol{\alpha}_t | \boldsymbol{\gamma}_{0:t+L}, \mathbf{y}_{1:t+L}) \times p(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L}) \quad (28)$$

where $p(d\boldsymbol{\alpha}_t | \boldsymbol{\gamma}_{0:t+L}, \mathbf{y}_{1:t+L})$ is a Gaussian distribution whose parameters may be computed using Kalman filter type techniques. Thus, given an approximation of $p(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$, an approximation of $p(d\mathbf{x}_t, d\boldsymbol{\theta}_t | \mathbf{y}_{1:t+L})$ may straightforwardly be obtained. Defining the marginal importance distribution and associated importance weight as

$$\begin{aligned} \pi(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L}) &= \int \pi(d\boldsymbol{\alpha}_{0:t+L} d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L}) \\ w(\boldsymbol{\gamma}_{0:t+L}) &\propto \frac{p(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})}{\pi(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})}, \end{aligned} \quad (29)$$

and assuming that a set of samples $\boldsymbol{\gamma}_{0:t+L}^{(i)}$ distributed according to $\pi(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$ is available, an alternative Bayesian importance sampling estimate of $I_L(f_t)$ follows as

$$\begin{aligned} \hat{I}_N^2(f_t) &= \frac{\sum_{i=1}^N p(\boldsymbol{\alpha}_t^{(i)} | \boldsymbol{\gamma}_{0:t+L}^{(i)}, \mathbf{y}_{1:t+L}) f_t(\mathbf{x}_t^{(i)}, \boldsymbol{\theta}_t^{(i)}) w(\boldsymbol{\gamma}_{0:t+L}^{(i)})}{\sum_{i=1}^N w(\boldsymbol{\gamma}_{0:t+L}^{(i)})} \\ &= \sum_{i=1}^N \tilde{w}_{0:t+L}^{(i)} p(\boldsymbol{\alpha}_t^{(i)} | \boldsymbol{\gamma}_{0:t+L}^{(i)}, \mathbf{y}_{1:t+L}) f_t(\mathbf{x}_t^{(i)}, \boldsymbol{\theta}_t^{(i)}), \end{aligned} \quad (30)$$

provided that $p(\boldsymbol{\alpha}_t | \boldsymbol{\gamma}_{0:t+L}, \mathbf{y}_{1:t+L}) f_t(\mathbf{x}_t, \boldsymbol{\theta}_t)$ can be evaluated in a closed-form expression. In (30) the marginal normalized importance weights are given by

$$\tilde{w}_{0:t+L}^{(i)} = \frac{w(\boldsymbol{\gamma}_{0:t+L}^{(i)})}{\sum_{j=1}^N w(\boldsymbol{\gamma}_{0:t+L}^{(j)})}, \quad i = 1, \dots, N. \quad (31)$$

Intuitively, to reach a given precision, $\hat{I}_{L,N}^2(f_t)$ will need a reduced number of samples over $\hat{I}_{L,N}^1(f_t)$, since it only requires samples from the marginal distribution $\pi(d\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})$. It can be proved that the variance of the estimates is subsequently reduced [8]. In our case this is important as at each time instant the number of parameters is (when assuming that all mixing filters and AR processes have the same length),

- $m^2 L - mL$ parameters for the mixing filters, where L can be large.
- m or 1 parameter(s) for the observation noise.
- $nl + n$ parameters for the autoregressive processes.
- n parameters for the sources.

It is not clear which integration will allow for the best variance reduction [8], but at least in terms of search in the parameters space the integration of the mixing filters and autoregressive filters seems preferable.

Given these results, the subsequent discussion will focus on Bayesian importance sampling methods to obtain approximations of $p(d\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})$ and $I_L(f_t)$ using an importance distribution of the form $\pi(d\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})$. The methods described up to now are batch methods. The next section illustrates how a sequential method may be obtained.

3.3 Sequential Bayesian Importance Sampling

The importance distribution at discrete time t may be factorized as

$$\begin{aligned} \pi(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L}) \\ = \pi(d\boldsymbol{\gamma}_0 | \mathbf{y}_{1:t+L}) \prod_{k=1}^{t+L} \pi(d\boldsymbol{\gamma}_k | \boldsymbol{\gamma}_{0:k-1}, \mathbf{y}_{1:k}), \end{aligned} \quad (32)$$

The aim is to obtain at any time t an estimate of the distribution $p(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$ and to be able to propagate this estimate in

time without modifying subsequently the past simulated trajectories $\boldsymbol{\gamma}_{0:t+L}^{(i)}$. This means that $\pi(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$ should admit $\pi(d\boldsymbol{\gamma}_{0:t-1+L} | \mathbf{y}_{1:t-1+L})$ as marginal distribution. This is possible if the importance distribution is restricted to be of the general form

$$\begin{aligned} \pi(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L}) \\ = \pi(d\boldsymbol{\gamma}_0) \prod_{k=1}^{t+L} \pi(d\boldsymbol{\gamma}_k | \boldsymbol{\gamma}_{0:k-1}, \mathbf{y}_{1:k}), \end{aligned} \quad (33)$$

Such an importance distribution allows recursive evaluation of the importance weights, *i.e.* $w(\boldsymbol{\gamma}_{0:t+L}) = w(\boldsymbol{\gamma}_{0:t-1+L}) w_{t+L}$, and in our particular case

$$\begin{aligned} \frac{p(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})}{\pi(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})} &\propto \frac{p(\boldsymbol{\gamma}_{0:t+L-1} | \mathbf{y}_{1:t+L-1})}{\pi(\boldsymbol{\gamma}_{0:t+L-1} | \mathbf{y}_{1:t+L-1})} \\ &\times \frac{p(\boldsymbol{\gamma}_{t+L} | \boldsymbol{\gamma}_{0:t+L}) p(\mathbf{x}_{t+L} | \mathbf{x}_{t+L-1}, \boldsymbol{\beta}_{t+L}) p(\boldsymbol{\beta}_{t+L} | \boldsymbol{\beta}_{t+L-1})}{\pi(\boldsymbol{\gamma}_{t+L} | \boldsymbol{\gamma}_{0:t+L-1}, \mathbf{y}_{1:t+L})} \end{aligned} \quad (34)$$

The quantity $p(d\mathbf{x}_{t+L} | \mathbf{x}_{t+L}, \boldsymbol{\beta}_{t+L})$ can be computed up to a normalizing constant using a one step ahead Kalman filter for the system given by Eq. (16) and $p(\mathbf{y}_{t+L} | \mathbf{x}_{0:t+L}, \boldsymbol{\beta}_{t+L})$ can be computed using a one step ahead Kalman filter of the system given by Eq. (14).

3.3.1 Choice of the Importance Distribution

There is an unlimited number of choices for the importance distribution $\pi(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$, the only restriction being that its support includes that of $p(d\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})$. Two possibilities are considered next. A possible strategy is to choose at time $t+L$ the importance distribution that minimizes the variance of the importance weights given $\boldsymbol{\gamma}_{0:t-1}$ and $\mathbf{y}_{1:t}$. The importance distribution that satisfies this condition is [9] $p(d\boldsymbol{\gamma}_{t+L} | \boldsymbol{\gamma}_{0:t-1+L}, \mathbf{y}_{1:t+L})$, with the associated incremental importance weight given by

$$\begin{aligned} w_{t+L} &\propto p(\mathbf{y}_{t+L} | \boldsymbol{\gamma}_{0:t-1+L}, \mathbf{y}_{1:t+L-1}) \\ &= \int p(\mathbf{y}_{t+L} | \boldsymbol{\gamma}_{0:t+L}, \mathbf{y}_{1:t+L-1}) p(d\boldsymbol{\gamma}_{t+L} | \boldsymbol{\gamma}_{t-1+L}). \end{aligned} \quad (35)$$

Direct sampling from the optimal importance distribution is difficult, and evaluating the importance weight is analytically intractable. The aim is thus in general to mimic the optimal distribution by means of tractable approximations, typically local linearization $p(d\boldsymbol{\gamma}_t | \boldsymbol{\gamma}_{0:t-1}, \mathbf{y}_{1:t})$. Instead, here we describe a mixed suboptimal method. We propose to sample the particles at time t according to two importance distributions π_1 and π_2 with proportions α and $1 - \alpha$ such that the importance weights $w(\boldsymbol{\gamma}_{0:t+L}^{(i)})$ have now the form (note that it would be possible to draw N_1 and N_2 randomly according to a Bernoulli distribution with parameter α , but this would increase the estimator variance)

$$\left\{ \begin{array}{l} \alpha \frac{p(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L}) / \mathbb{E}_{\pi_1} \left[\frac{p(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})}{\pi_1(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})} \right]}{\pi_1(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})} \\ (1 - \alpha) \frac{p(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L}) / \mathbb{E}_{\pi_2} \left[\frac{p(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})}{\pi_2(\boldsymbol{\gamma}_{0:t+L} | \mathbf{y}_{1:t+L})} \right]}{\pi_2(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})} \end{array} \right\} \quad (36)$$

which in practice is estimated as $\bar{w}_{0:t+L}^{(i)}$

$$\left\{ \begin{array}{l} \alpha \frac{p(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L}) / \pi_1(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})}{\sum_{j=1}^{N_1} p(\boldsymbol{\gamma}_{0:t+L}^{(j)} | \mathbf{y}_{1:t+L}) / \pi_1(\boldsymbol{\gamma}_{0:t+L}^{(j)} | \mathbf{y}_{1:t+L})} \\ (1 - \alpha) \frac{p(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L}) / \pi_2(\boldsymbol{\gamma}_{0:t+L}^{(i)} | \mathbf{y}_{1:t+L})}{\sum_{j=\alpha N+1}^N p(\boldsymbol{\gamma}_{0:t+L}^{(j)} | \mathbf{y}_{1:t+L}) / \pi_2(\boldsymbol{\gamma}_{0:t+L}^{(j)} | \mathbf{y}_{1:t+L})} \end{array} \right\} \quad (37)$$

The importance distribution $\pi_1(d\mathbf{x}_{t+L} | \mathbf{x}_{1:t+L-1}, \boldsymbol{\beta}_{1:t+L}, \mathbf{y}_{1:t+L})$ will be taken to be a normal distribution centered around zero with variance $\sigma_{\mathbf{x}}^2$, and $\pi_2(d\mathbf{x}_{t+L} | \mathbf{x}_{1:t+L-1}, \boldsymbol{\beta}_{1:t+L}, \mathbf{y}_{1:t+L})$ is taken to be $p(d\mathbf{x}_{t+L}^{(i)} | \mathbf{m}_{t+L|t+L}^{\mathbf{a}, \mathbf{h}(i)}, \mathbf{P}_{t+L|t+L}^{\mathbf{a}, \mathbf{h}(i)}, \boldsymbol{\beta}_{t+L}^{(i)}, \mathbf{y}_{1:t+L})$ which is a Gaussian distribution obtained from a one step ahead Kalman filter for the state space model described in (11) with $\mathbf{m}_{t+L|t+L}^{\mathbf{a}(i)}$ and $\mathbf{m}_{t+L|t+L}^{\mathbf{h}(i)}$ as values for $\mathbf{a}^{(i)}$ and $\mathbf{h}^{(i)}$ and initial variances $\mathbf{P}_{t+L|t+L}^{\mathbf{a}}$ and $\mathbf{P}_{t+L|t+L}^{\mathbf{h}}$. The variances are sampled from their prior distributions, and expression (34) is used to compute $\bar{w}_{0:t+L}^{(i)}$. Note that other importance distributions are possible, but this approach yields good results and seems to preserve diversity of the samples.

3.3.2 Degeneracy of the Algorithm

For importance distributions of the form specified by (33) the variance of the importance weights can only increase (stochastically) over time, see [9] and references therein. It is thus impossible to avoid a degeneracy phenomenon. Practically, after a few iterations of the algorithm, all but one of the normalized importance weights are very close to zero, and a large computational effort is devoted to updating trajectories whose contribution to the final estimate is almost zero. For this reason it is of crucial importance to include selection and diversity. This is discussed in more detail in the following section.

3.4 Selection and diversity

The purpose of a selection (or resampling) procedure is to discard particles with low normalized importance weights and multiply those with high normalized importance weights, so as to avoid the degeneracy of the algorithm. A selection procedure associates with each particle, say $\tilde{\gamma}_{0:t}^{(i)}$, a number of children $N_i \in \mathbb{N}$, such that $\sum_{i=1}^N N_i = N$, to obtain N new particles $\gamma_{0:t}^{(i)}$. If $N_i = 0$ then $\tilde{\gamma}_{0:t}^{(i)}$ is discarded, otherwise it has N_i children at time $t+1$. After the selection step the normalized importance weights for all the particles are reset to N^{-1} , thus discarding all information regarding the past importance weights. Thus, the normalized importance weight prior to selection in the next time step is proportional to (34). These will be denoted as $\tilde{w}_t^{(i)}$, since they do not depend on any past values of the normalized importance weights. If the selection procedure is performed at each time step, then the approximating distribution before the selection step is given by $\tilde{p}_N(d\gamma_{0:t} | \mathbf{y}_{1:t}) = \sum_{i=1}^N \tilde{w}_t^{(i)} \delta_{\tilde{\gamma}_{0:t}^{(i)}}(d\gamma_{0:t})$, and the one after the selection step follows as $\hat{p}_N(d\gamma_{0:t} | \mathbf{y}_{1:t}) = N^{-1} \sum_{i=1}^N \delta_{\gamma_{0:t}^{(i)}}(d\gamma_{0:t})$. We choose systematic sampling [8] for its good variance properties.

However selection poses another problem. During the resampling stage any particular particle with a high importance weight will be duplicated many times. In particular, when $L > 0$, the trajectories are resampled L times from time $t+1$ to $t+L$ so that very few distinct trajectories remain at time $t+L$. This is the classical problem of depletion of samples. As a result the cloud of particles may eventually collapse to a single particle. This degeneracy leads to poor approximations of the distributions of interest. Several suboptimal methods have been proposed to overcome this problem and introduce diversity amongst the particles. Most of these are based on kernel density methods [9, 10], which approximate the probability distribution using a kernel density estimate based on the current set of particles, and sample a new set

of distinct particles from it. However, the choice and configuration of a specific kernel are not always straightforward. Moreover, these methods introduce additional Monte Carlo variation. In the next subsection it is shown how MCMC methods may be combined with sequential importance sampling to introduce diversity amongst the samples without increasing the Monte Carlo variation.

An efficient way of limiting sample depletion consists of simply adding a MCMC step to the simulation-based filter/fixed-lag smoother (see Berzuini and Gilks in [9], and [12] for an introduction to MCMC methods). This introduces diversity amongst the samples and thus drastically reduces the problem of depletion of samples. Assume that, at time $t+L$, the particles $\gamma_{0:t+L}'^{(i)}$ are marginally distributed according to $p(d\gamma_{0:t+L}' | \mathbf{y}_{1:t+L})$. If a transition kernel $K(\gamma_{0:t+L}' | d\gamma_{0:t+L}')$ with invariant distribution $p(d\gamma_{0:t+L}' | \mathbf{y}_{1:t+L})$ is applied to each of the particles, then the new particles $\gamma_{0:t+L}^{(i)}$ are still distributed according to the distribution of interest. Any of the standard MCMC methods, such as the Metropolis-Hastings (MH) algorithm or Gibbs sampler, may be used. However, contrary to classical MCMC methods, the transition kernel does not need to be ergodic. Not only does this method introduce no additional Monte Carlo variation, but it improves the estimates in the sense that it can only reduce the total variation norm [12] of the current distribution of the particles with respect to the target distribution.

3.5 Implementation Issues

3.5.1 Algorithm

Given at time $t+L-1$, $N \in \mathbb{N}^*$ particles $\gamma_{0:t+L-1}^{(i)}$ distributed approximately according to $p(d\gamma_{0:t+L-1} | \mathbf{y}_{1:t+L-1})$, the Monte Carlo fixed-lag smoother proceeds as follows at time $t+L$.

Monte Carlo filter/fixed-lag smoother

Sequential Importance Sampling Step

- For $i = 1, \dots, N$, $\tilde{\gamma}_{t+L}^{(i)} \sim \pi(d\gamma_{t+L} | \mathbf{x}_{0:t+L-1}^{(i)}, \mathbf{y}_{1:t+L})$ and set $\tilde{\gamma}_{0:t+L}^{(i)} = (\gamma_{0:t+L-1}^{(i)}, \tilde{\gamma}_{t+L}^{(i)})$.
- For $i = 1, \dots, N$, compute the normalized importance weights $\tilde{w}_{t+L}^{(i)}$ using (34) and (37).

Selection Step

- Multiply / discard particles $\tilde{\gamma}_{0:t+L}^{(i)}$ w.r.t. high / low normalized importance weights to obtain N particles $\gamma_{0:t+L}'^{(i)}$, e.g. using systematic sampling.

MCMC Step

- For $i = 1, \dots, N$, apply to $\gamma_{0:t+L}'^{(i)}$ a Markov transition kernel $K(\gamma_{0:t+L}^{(i)} | d\gamma_{0:t+L}'^{(i)})$ with invariant distribution $p(d\gamma_{0:t+L} | \mathbf{y}_{1:t+L})$ to obtain N particles $\gamma_{0:t+L}^{(i)}$.
-

3.5.2 Implementation of the MCMC Steps

There is an unlimited number of choices for the MCMC transition kernel. Here a one-at-a-time MH algorithm is adopted that updates at time $t+L$ the values of the Markov process from time t to $t+L$. More specifically, $\gamma_k^{(i)}$, $k = t, \dots, t+L$, $i = 1, \dots, N$, is sampled according to an MCMC with $p(d\gamma_k | \gamma_{-k}^{(i)}, \mathbf{y}_{1:t+L})$ as target distribution, where $\gamma_{-k}^{(i)} \triangleq (\gamma_{0:t-1}^{(i)}, \gamma_t^{(i)}, \dots, \gamma_{k-1}^{(i)}, \gamma_{k+1}^{(i)}, \dots, \gamma_{t+L}^{(i)})$. Evaluation of $p(d\gamma_k | \gamma_{-k}^{(i)}, \mathbf{y}_{1:t+L})$ can be done efficiently via a backward-forward algorithm of $O(L+1)$ complexity [7]. The algorithm is fully described in [2]. The computational complexity of the whole algorithm at each iteration is clearly $O(N)$. At first glance, it could appear necessary to keep in memory the paths of all the trajectories $\gamma_{0:t}^{(i)}$, so that the storage requirements would increase linearly with time. In fact, the importance distribution $\pi_1, \pi(\gamma_t | \gamma_{0:t-1}, \mathbf{y}_{1:t})$ and the associated importance weights, depend on $\gamma_{0:t-1}$ only via a set of low-dimensional sufficient statistics $\mathbf{m}_{t|t}^{\mathbf{a}, \mathbf{h}(i)}, \mathbf{P}_{t|t}^{\mathbf{a}, \mathbf{h}(i)}$, and only these values need to be kept in memory. Thus, the storage requirements are also $O(N)$ and do not increase over time.

4 Simulations

To illustrate the efficiency of our method we have applied the procedure to two scenarios. In the first case we generated two stationary time invariant order 2 autoregressive processes, mixed by two filters of length 5. In the second case we made the phase of the poles of the AR processes evolve from .1 to .4 from $t = 1, \dots, 250$ and from .4 to .1 for $t = 251, \dots, 500$. The algorithm was run with 100 particles with $L = 0$, the parameters of the system and the sources were initialized at random. The matrices $\mathbf{B}^{\mathbf{a}}$ and $\mathbf{B}^{\mathbf{h}}$ were set to $.01 \times \mathbf{I}$. The results for the two cases are presented on Fig. (1) and Fig. (2), where the two original sources and the sources estimated on-line using the particle filter are displayed. Note that for the first scenario convergence really occurs from iteration 150 approximately, and that source 2 was apparently mistaken with source 1 before. For the second scenario, as expected, the problem is much harder, especially when the two poles are close, or when zero/pole cancellation occurs. Application of the algorithm to real speech signals is currently under investigation.

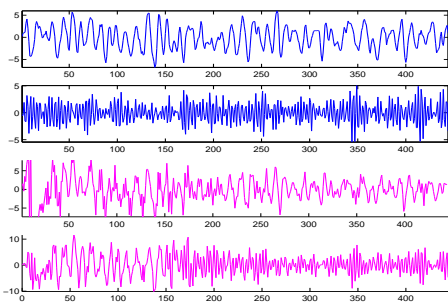


Figure 1: True (two top) and estimated (two bottom) sources.

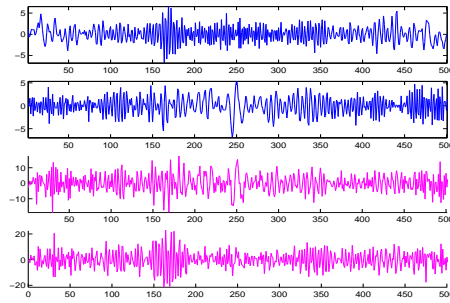


Figure 2: True (two top) and estimated (two bottom) sources.

5 REFERENCES

- [1] B.D.O. Anderson and J.B. Moore, *Optimal Filtering*, Prentice-Hall, Englewood Cliffs, 1979.
- [2] C. Andrieu, S.J. Godsill, "A Particle filter for audio source separation", Tech. Rep. Cambridge University Engineering Department, 2000.
- [3] D.L. Aspach and H.W. Sorenson, "Nonlinear Bayesian estimation using Gaussian sum approximations", *IEEE Trans. Auto. Cont.*, vol. 17, no. 4, pp. 439448, 1972.
- [4] H. Broman, U. Lindgren, H. Sahlin and P. Stoica, "Source Separation: A TITO System Identification Approach", *Signal Processing* 1998.
- [5] R.S. Bucy and K.D. Senne, "Digital synthesis of nonlinear filters", *Automatica*, vol. 7, 1971, pp. 287-298.
- [6] D. Crisan, P. Del Moral and T. Lyons, "Discrete filtering using branching and interacting particle systems", to appear *Markov Processes and Related Fields*, 2000.
- [7] A. Doucet and C. Andrieu, "Iterative algorithms for optimal state estimation of jump Markov linear systems", in *Proc. Conf. IEEE ICASSP*, 1999.
- [8] A. Doucet, J.F.G. de Freitas and N.J. Gordon (eds.), *Sequential Monte Carlo Methods in Practice*, Springer-Verlag, to appear June 2000.
- [9] A. Doucet, S.J. Godsill and C. Andrieu, "On sequential Monte Carlo sampling methods for Bayesian filtering", to appear *Statistics and Computing*, 2000.
- [10] N.J. Gordon, D.J. Salmond and A.F.M. Smith, "Novel approach to nonlinear/non-Gaussian Bayesian state estimation", *IEE Proceedings-F*, vol. 140, no. 2, pp. 107-113, 1993.
- [11] J.E. Handschin and D.Q. Mayne, "Monte Carlo techniques to estimate the conditional expectation in multi-stage nonlinear filtering", *Int. J. Cont.*, vol. 9, no. 5, 1969, pp. 547-559.
- [12] C.P. Robert and G. Casella, *Monte Carlo Statistical Methods*, Springer-Verlag, 1999.
- [13] E. Weinstein, A.V. Oppenheim, M. Feder and J.R. Buck, "Iterative and Sequential Algorithms for Multisensor Signal Enhancement", *IEEE. Trans. SP*, Vol. 42, No. 4, April 1994.