

Identification of jump Markov linear models using particle filters

Andreas Svensson, Thomas B. Schön and Fredrik Lindsten*

August 29, 2018

Abstract

Jump Markov linear models consists of a finite number of linear state space models and a discrete variable encoding the jumps (or switches) between the different linear models. Identifying jump Markov linear models makes for a challenging problem lacking an analytical solution. We derive a new expectation maximization (EM) type algorithm that produce maximum likelihood estimates of the model parameters. Our development hinges upon recent progress in combining particle filters with Markov chain Monte Carlo methods in solving the nonlinear state smoothing problem inherent in the EM formulation. Key to our development is that we exploit a conditionally linear Gaussian substructure in the model, allowing for an efficient algorithm.

1 Introduction

Consider the following *jump Markov linear model* on state space form

$$s_{t+1} \mid s_t \sim p(s_{t+1} \mid s_t), \quad (1a)$$

$$z_{t+1} = A_{s_{t+1}} z_t + B_{s_{t+1}} u_t + w_t, \quad (1b)$$

$$y_t = C_{s_t} z_t + D_{s_t} u_t + v_t, \quad (1c)$$

where \sim means distributed according to and the (discrete) variable s_t takes values in $\{1, \dots, K\}$ (which can be thought of as different *modes* which the model is *jumping* between) and the (continuous) variable z_t lives in \mathbb{R}^{n_z} . Hence, the state variable consists of $x_t \triangleq (z_t, s_t)$. Furthermore, $v_t \in \mathbb{R}^{n_y}$ and $w_t \in \mathbb{R}^{n_z}$ are zero mean white Gaussian noise and $\mathbb{E}w_t w_t^T = Q_{s_{t+1}}$, $\mathbb{E}v_t v_t^T = R_{s_t}$ and $\mathbb{E}w_t v_t^T \equiv 0$. The output (or measurement) is $y_t \in \mathbb{R}^{n_y}$, the input is $u_t \in \mathbb{R}^{n_u}$. As K is finite, $p(s_{t+1} \mid s_t)$ can be defined via a matrix $\Pi \in \mathbb{R}^{K \times K}$ with entries $\pi_{mn} \triangleq p(s_{t+1} = n \mid s_t = m)$.

We are interested in off-line identification of jump Markov linear models on the form (1) for the case of an *unknown jump sequence*, but the number of modes K is known. More specifically, we will formulate and solve the Maximum Likelihood (ML) problem to compute an estimate of the static parameters θ of a jump Markov linear model based on a batch of measurements $y_{1:T} \triangleq \{y_1, \dots, y_T\}$ and (if available) inputs $u_{1:T}$ by solving,

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in \Theta} p_{\theta}(y_{1:T}). \quad (2)$$

Here $\theta \triangleq \{\{A_n, B_n, C_n, D_n, Q_n, R_n\}_{n=1}^K, \Pi\}$, i.e., all unknown static parameters in model (1). Here, and throughout the paper, the dependence on the inputs $u_{1:T}$ is implicit.

*This work was supported by the project *Probabilistic modelling of dynamical systems* (Contract number: 621-2013-5524) funded by the Swedish Research Council (VR) and the project *Learning of complex dynamical systems* (Contract number: 637-2014-466) funded by the Swedish Research Council (VR). Andreas Svensson and Thomas B. Schön are with the Department of Information Technology, Uppsala University, Sweden [andreas.svensson](mailto:andreas.svensson@it.uu.se), [thomas.schon](mailto:thomas.schon@it.uu.se) and Fredrik Lindsten is with the Department of Engineering, University of Cambridge, UK fredrik.lindsten@eng.cam.ac.uk

Solving (2) is challenging and there are no closed form solutions available. Our approach is to derive an expectation maximization (EM) [10] type of solution, where the strategy is to separate the original problem into two closely linked problems. The first problem is a challenging, but manageable nonlinear state smoothing problem and the second problem is a tractable optimization problem. The nonlinear smoothing problem we can solve using a combination of sequential Monte Carlo (SMC) methods (particle filters and particle smoothers) [11] and Markov chain Monte Carlo (MCMC) methods [27]. More specifically we will make use of particle MCMC (PMCMC), which is a systematic way of exploring the strengths of both approaches by using SMC to construct the necessary high-dimensional Markov kernels needed in MCMC [1, 19].

Our main contribution is a new maximum likelihood estimator that can be used to identify jump Markov linear models on the form (1). The estimator exploits the conditionally linear Gaussian substructure that is inherent in (1) via Rao-Blackwellization. More specifically we derive a Rao-Blackwellized version of the particle stochastic approximation expectation maximization (PSAEM) algorithm recently introduced in [18].

Jump Markov linear models, or switching linear models, is a fairly well studied class of hybrid systems. For recent overviews of existing system identification methods for jump Markov linear models, see [13, 22]. Existing approaches considering the problem under study here include two stage methods, where the data is first segmented (using e.g. change detection type of methods) and the individual models are then identified for each segment, see e.g. [23, 6]. There has also been approximate EM algorithms proposed for identification of hybrid systems [5, 15] and the very recent [3] (differing from our method in that we use stochastic approximation EM and Rao-Blackwellization). There are also relevant relationships to the PMCMC solutions introduced in [33] and the SMC-based on-line EM solution derived in [34].

There are also many approaches considering the more general problem with an unknown number of modes K and an unknown state dimension n_z , see e.g. [12] and [4], making use of Bayesian nonparametric models and mixed integer programming, respectively.

2 Expectation maximization algorithms

The EM algorithm [10] provides an iterative method for computing maximum likelihood estimates of the unknown parameters θ in a probabilistic model involving latent variables. In the jump Markov linear model (1) we observe $y_{1:T}$, whereas the state $x_{1:T}$ is latent.

The EM algorithm maximizes the likelihood by iteratively maximizing the *intermediate quantity*

$$\mathcal{Q}(\theta, \theta') \triangleq \int \log p_\theta(x_{1:T}, y_{1:T}) p_{\theta'}(x_{1:T} | y_{1:T}) dx_{1:T}. \quad (3)$$

More specifically, the procedure is initialized in $\theta_0 \in \Theta$ and then iterates between computing an expected (E) value and solving a maximization (M) problem,

$$\begin{aligned} \text{(E)} \quad & \text{Compute } \mathcal{Q}(\theta, \theta_{k-1}). \\ \text{(M)} \quad & \text{Compute } \theta_k = \arg \max_{\theta \in \Theta} \mathcal{Q}(\theta, \theta_{k-1}). \end{aligned}$$

Intuitively, this can be thought of as ‘selecting the new parameters as the ones that make the given measurements *and* the current state estimate as likely as possible’.

The use of EM type algorithms to identify dynamical systems is by now fairly well explored for both linear and nonlinear models. For linear models, there are explicit expressions for all involved quantities, see e.g. [14, 30]. For nonlinear models the intermediate quantity $\mathcal{Q}(\theta, \theta')$ is intractable and we are forced to approximate solutions; see e.g. [18, 29, 21, 7]. This is the case also for the model (1) under study in this work. Indeed, the maximization step can be solved in closed form for the model (1), but (3) is still intractable in our case.

It is by now fairly well established that we can make use of sequential Monte Carlo (SMC) [11] or particle Markov chain Monte Carlo (PMCMC) [1] methods to approximate the joint smoothing distribution for a

general nonlinear model arbitrarily well according to

$$\widehat{p}(x_{1:T} | y_{1:T}) = \sum_{i=1}^N w_T^i \delta_{x_{1:T}^i}(x_{1:T}), \quad (4)$$

where $x_{1:T}^i$ are random samples with corresponding importance weights w_T^i , δ_x is a point-mass distribution at x and we refer to $\{x_{1:T}^i, w_T^i\}_{i=1}^N$ as a *weighted particle system*. The particle smoothing approximation (4) can be used to approximate the integral in (3). Using this approach within EM, we obtain the particle smoothing EM (PSEM) method [21, 29]. PSEM can be viewed as an SMC-analogue of the well known Monte Carlo EM (MCEM) algorithm [32].

However, it has been recognized that MCEM, and analogously PSEM, makes inefficient use of the generated samples [9]. This is particularly true when the simulation step is computationally expensive, which is the case when using SMC or PMCMC. To address this shortcoming, [9] proposed to use a *stochastic approximation* (SA) [26] of the intermediate quantity instead of a vanilla Monte Carlo approximation, resulting in the stochastic approximation EM (SAEM) algorithm. The SAEM algorithm replaces the intermediate quantity \mathcal{Q} in EM with

$$\widehat{\mathcal{Q}}_k(\theta) = (1 - \gamma_k) \widehat{\mathcal{Q}}_{k-1}(\theta) + \gamma_k \log p_\theta(y_{1:T}, x_{1:T}[k]), \quad (5)$$

with $\{\gamma_k\}_{k=1}^\infty$ being a sequence of step sizes which fulfils $\sum_{k=1}^\infty \gamma_k = \infty$ and $\sum_{k=1}^\infty \gamma_k^2 < \infty$. In the above, $x_{1:T}[k]$ is a sample state trajectory, simulated from the joint smoothing distribution $p_{\theta_k}(x_{1:T} | y_{1:T})$. It is shown by [9] that the SAEM algorithm—which iteratively updates the intermediate quantity according to (5) and computes the next parameter iterate by maximizing this stochastic approximation—enjoys good convergence properties. Indeed, despite the fact that the method requires only a single sample $x_{1:T}[k]$ at each iteration, the sequence $\{\theta_k\}_{k \geq 1}$ will converge to a maximizer of $p_\theta(y_{1:T})$ under reasonably weak assumptions.

However, in our setting it is not possible to simulate from the joint smoothing distribution $p_{\theta_k}(x_{1:T} | y_{1:T})$. We will therefore make use of the particle SAEM (PSAEM) method [18], which combines recent PMCMC methodology with SAEM. Specifically, we will exploit the structure of (1) to develop a Rao-Blackwellized PSAEM algorithm.

We will start our development in the subsequent section by considering the smoothing problem for (1). We derive a PMCMC-based Rao-Blackwellized smoother for this model class. The proposed smoother can, principally, be used to compute (3) within PSEM. However, a more efficient approach is to use the proposed smoother to derive a Rao-Blackwellized PSAEM algorithm, see Section 4.

3 Smoothing using Monte Carlo methods

For smoothing, that is, finding $p_\theta(x_{1:t}|y_{1:t}) = p_\theta(s_{1:T}, z_{1:T}|y_{1:T})$, various Monte Carlo methods can be applied. We will use an MCMC based approach, as it fits very well in the SAEM framework (see e.g. [2, 17]), which together shapes the PSAEM algorithm. The aim of this section is therefore to derive an MCMC-based smoother for jump Markov linear models.

To gain efficiency, the jump sequence $s_{1:T}$ and the linear states $z_{1:T}$ are separated using conditional probabilities as

$$p_\theta(s_{1:T}, z_{1:T}|y_{1:T}) = p_\theta(z_{1:T}|s_{1:T}, y_{1:T}) p_\theta(s_{1:T}|y_{1:T}). \quad (6)$$

This allows us to infer the conditionally linear states $z_{1:T}$ using closed form expressions. Hence, it is only the jump sequence $s_{1:T}$ that has to be computed using approximate inference. This technique is referred to as Rao-Blackwellization [8].

3.1 Inferring the linear states: $p(z_{1:T}|s_{1:T}, y_{1:T})$

State inference in linear Gaussian state space models can be performed exactly in closed form. More specifically, the Kalman filter provides the expressions for the filtering PDF $p_\theta(z_t|s_{1:t}, y_{1:t}) = \mathcal{N}(z_t|\widehat{z}_{f;t}, P_{f;t})$

and the one step predictor PDF $p_\theta(z_{t+1}|s_{1:t+1}, y_{1:t}) = \mathcal{N}(z_t|\widehat{z}_{p;t+1}, P_{p;t+1})$. The marginal smoothing PDF $p_\theta(z_t|s_{1:T}, y_{1:T}) = \mathcal{N}(z_t|\widehat{z}_{s;t}, P_{s;t})$ is provided by the Rauch-Tung-Striebel (RTS) smoother [25]. See, e.g., [16] for the relevant results. Here, we use $\mathcal{N}(x | \mu, \Sigma)$ to denote the PDF for the (multivariate) normal distribution with mean μ and covariance matrix Σ .

3.2 Inferring the jump sequence: $p(s_{1:T}|y_{1:T})$

To find $p(s_{1:T}|y_{1:T})$, an MCMC approach is used. First, the concept of using Markov kernels for smoothing is introduced, and then the construction of the kernel itself follows.

MCMC makes use of ergodic theory for statistical inference. Let \mathcal{K}_θ be a Markov kernel (to be defined below) on the T -fold product space $\{1, \dots, K\}^T$. Note that the jump sequence $s_{1:T}$ lives in this space. Furthermore, assume that \mathcal{K}_θ is *ergodic* with unique stationary distribution $p_\theta(s_{1:T}|y_{1:T})$. This implies that by simulating a Markov chain with transition kernel \mathcal{K}_θ , the marginal distribution of the chain will approach $p_\theta(s_{1:T}|y_{1:T})$ in the limit.

Specifically, let $s_{1:T}[0]$ be an arbitrary initial state with $p_\theta(s_{1:T}[0]|y_{1:T}) > 0$ and let $s_{1:T}[k] \sim \mathcal{K}_\theta(\cdot|s_{1:T}[k-1])$ for $k \geq 1$, then by the ergodic theorem [27]:

$$\frac{1}{n} \sum_{k=1}^n h(s_{1:T}[k]) \rightarrow \mathbb{E}_\theta [h(s_{1:T})|y_{1:T}], \quad (7)$$

as $n \rightarrow \infty$ for any function $h : \{1, \dots, K\}^T \mapsto \mathbb{R}$. This allows a smoother to be constructed as in Algorithm 1.

Algorithm 1 MCMC smoother

- 1: Initialize $s_{1:T}[0]$ arbitrarily
 - 2: **for** $k \geq 1$ **do**
 - 3: Generate $s_{1:T}[k] \sim \mathcal{K}_\theta(\cdot|s_{1:T}[k-1])$
 - 4: **end for**
-

We will use the *conditional particle filter with ancestor sampling* (CPF-AS) [19] to construct the Markov kernel \mathcal{K}_θ . The CPF-AS is similar to a standard particle filter, but with the important difference that one particle trajectory (jump sequence), $s'_{1:T}$, is specified *a priori*.

The algorithm statement for the CPF-AS can be found in, e.g., [19]. Similar to an auxiliary particle filter [11], the propagation of $p_\theta(s_{1:t-1}|y_{1:t-1})$ (approximated by $\{s_{1:t-1}^i, w_{t-1}^i\}_{i=1}^N$) to time t is done using the *ancestor indices* $\{a_t^i\}_{i=1}^N$. To generate s_t^i , the ancestor index is sampled according to $\mathbb{P}(a_t^i = j) \propto w_{t-1}^j$, and s_t^i as $s_t^i \sim p_\theta(s_t|s_{t-1}^{a_t^i})$. The trajectories are then augmented as $s_{1:t}^i = \{s_{1:t-1}^{a_t^i}, s_t^i\}$.

This is repeated for $i = 1, \dots, N-1$, whereas s_t^N is set as $s_t^N = s'_t$. To ‘find’ the history for s_t^N , the ancestor index a_t^N is drawn with probability

$$\mathbb{P}(a_t^N = i) \propto p_\theta(s_{1:t-1}^i|s'_{1:t-1}, y_{1:T}). \quad (8)$$

The probability density in (8) is proportional to

$$p_\theta(y_{t:T}, s'_{t:T}|s_{1:t-1}^i, y_{1:t-1}) p_\theta(s_{1:t-1}^i|y_{1:t-1}), \quad (9)$$

where the last factor is the importance weight w_{t-1}^i .

By sampling $s_{1:T}[k+1] = s_{1:T}^J$ from the rendered set of trajectories $\{s_{1:T}^i, w_T^i\}_{i=1}^N$ with $\mathbb{P}(J = j) = w_T^j$, a Markov kernel \mathcal{K}_θ mapping $s_{1:T}[k] = s'_{1:T}$ to $s_{1:T}[k+1]$ is obtained. For this Markov kernel to be useful for statistical inference we require that (i) it is ergodic, and (ii) it admits $p_\theta(s_{1:T}|y_{1:T})$ as its unique limiting distribution. While we do not dwell on the (rather technical) details here, we note that these requirements are indeed fulfilled; see [19].

3.3 Rao-Blackwellization

Rao-Blackwellization of particle filters is a fusion of the Kalman filter and the particle filter based on (6), and it is described in, e.g., [28]. However, Rao-Blackwellization of a particle smoother is somewhat more involved since the process $x_t|y_{1:T}$ is Markovian, but not $s_t|y_{1:T}$ (with z_t marginalized, see, e.g., [33] and [20] for various ways to handle this).

A similar problem as for the particle smoothers arises in the ancestor sampling (8) in the CPF-AS. In the case of a non-Rao-Blackwellized CPF-AS, (8) reduces to $w_{t-1}^i p(x_t^i|x_{t-1}^i)$ [19]. This does not hold in the Rao-Blackwellized case.

To handle this, (8) can be rewritten as

$$w_{t-1}^i p(y_{t:T}, s'_{t:T}|s_{1:t-1}^i, y_{1:t-1}). \quad (10)$$

Using the results from Section 4.4 in [20] (adapted to model (1)), this can be written (omitting w_{t-1}^i , and with the notation $\|z\|_{\Omega}^2 \triangleq z^T \Omega z$, $P \triangleq \Gamma \Gamma^T$, i.e. the Cholesky factorization, $Q_t \triangleq F_t F_t^T$ and $A_t \triangleq A_{s_t}$ etc.)

$$p(y_{t:T}, s'_{t:T}|s_{1:t-1}^i, y_{1:t-1}) \propto Z_{t-1} |\Lambda_{t-1}|^{-1/2} \exp\left(-\frac{1}{2} \eta_{t-1}\right), \quad (11a)$$

with

$$\Lambda_t = \Gamma_{f;t}^{i,T} \Omega_t \Gamma_{f;t}^i + I, \quad (11b)$$

$$\eta_t = \|\hat{z}_{f;t}^i\|_{\Omega_t}^2 - 2\lambda_t^T \hat{z}_{f;t}^{i,T} - \|\Gamma_{f;t}^n (\lambda_t - \Omega_t \hat{z}_{f;t}^n)\|_{M_t^{-1}}^2, \quad (11c)$$

where

$$\Omega_t = A_{t+1}^T \left(I - \hat{\Omega}_{t+1} F_{t+1} M_{t+1}^{-1} F_{t+1}^T \right) \hat{\Omega}_{t+1} A_{t+1}, \quad (11d)$$

$$\hat{\Omega}_t = \Omega_t + C_t^T R_t^{-1} C_t, \quad (11e)$$

$$M_t = F_t^T \hat{\Omega}_t F_t + I, \quad (11f)$$

$$\lambda_t = A_{t+1}^T \left(I - \hat{\Omega}_{t+1} F_{t+1} M_{t+1}^{-1} F_{t+1}^T \right) m_t, \quad (11g)$$

$$\hat{\lambda}_t = \lambda_t + C_t^T R_t^{-1} (y_t - D_t u_t), \quad (11h)$$

$$m_t = (\hat{\lambda}_{t+1} - \hat{\Omega}_{t+1} B_{t+1} u_{t+1}). \quad (11i)$$

and $\Omega_T = 0$ and $\lambda_T = 0$. The Rao-Blackwellization also includes an RTS smoother for finding $p_{\theta}(z_{1:T}|s_{1:T}, y_{1:T})$.

Summarizing the above development, the Rao-Blackwellized CPF-AS (for the jump Markov linear model (1)) is presented in Algorithm 2, where

$$p_{\theta}(y_t|s_{1:t}^i, y_{1:t-1}) = \mathcal{N}(y_t; C_{s_t^i} \hat{z}_{p;t}^n + D_{s_t^i} u_t, C_{s_t^i} P_{p;t} C_{s_t^i}^T + R_{s_t^i}) \quad (12)$$

is used. Note that the discrete state s_t is drawn from a discrete distribution defined by Π , whereas the linear state z_t is handled analytically. The algorithm implicitly defines a Markov kernel \mathcal{K}_{θ} that can be used in Algorithm 1 for finding $p(s_{1:T}|y_{1:T})$, or, as we will see, be placed in an SAEM framework to estimate θ (both yielding PMCMC [1] constructions).

4 Identification of jump Markov linear models

In the previous section, an ergodic Markov kernel \mathcal{K}_{θ} leaving $p_{\theta}(s_{1:T}|y_{1:T})$ invariant was found as a Rao-Blackwellized CPF-AS summarized in Algorithm 2. This will be used together with SAEM, as it allows us to make one parameter update at each step of the Markov chain smoother in Algorithm 1, as presented as PSAEM in [18]. (However, following [18], we make use of all the particles generated by CPF-AS, and not only $s_{1:T}[k+1]$, to compute the intermediate quantity in the SAEM.)

Algorithm 2 Rao-Blackwellized CPF-AS

Input: $s'_{1:T} = s_{1:T}[k]$

Output: $s_{1:T}[k+1]$ (A draw from $\mathcal{K}_\theta(\cdot|s_{1:T}[k])$ and $\{s_{1:T}^i, w_T^i\}_{i=1}^N$)

- 1: Draw $s_1^i \sim p_1(s_1|y_1)$ for $i = 1, \dots, N-1$.
 - 2: Compute $\{\Omega_t, \lambda_t\}_{t=1}^T$ for $s'_{1:T}$ according to (11d) - (11i).
 - 3: Set $(s_1^N, \dots, s_T^N) = (s'_1, \dots, s'_T)$.
 - 4: Compute $\widehat{z}_{f,1}^i$ and $P_{f,1}^i$ $i = 1, \dots, N$.
 - 5: Set $w_1^i \propto p_\theta(y_1|s_1^i)$ (12) for $i = 1, \dots, N$ s.t. $\sum_i w_1^i = 1$
 - 6: **for** $t = 2$ to T **do**
 - 7: Draw a_t^i with $\mathbb{P}(a_t^i = j) = w_{t-1}^j$ for $i = 1, \dots, N-1$.
 - 8: Draw s_t^i with $\mathbb{P}(s_t^i = n) = \pi_{s_{t-1}^i, n}$ for $i = 1, \dots, N-1$.
 - 9: Compute $\{\Lambda_{t-1}^i, \eta_t^i\}$ according to (11b)-(11c).
 - 10: Draw a_t^N with $\mathbb{P}(a_t^N = i) \propto w_{t-1}^i \pi_{s_{t-1}^i, s_t^N} |\Lambda_{t-1}^i|^{-1/2} \exp(-\frac{1}{2}\eta_{t-1}^i)$.
 - 11: Set $s_{1:t}^i = \{s_{1:t-1}^{a_t^i}, s_t^i\}$ for $i = 1, \dots, N$.
 - 12: Set $\widehat{z}_{f,1:t-1}^i = \widehat{z}_{f,1:t-1}^{a_t^i}$, $P_{f,1:t-1}^i = P_{f,1:t-1}^{a_t^i}$, $\widehat{z}_{p,1:t-1}^i = \widehat{z}_{p,1:t-1}^{a_t^i}$ and $P_{p,1:t-1}^i = P_{p,1:t-1}^{a_t^i}$ for $i = 1, \dots, N$.
 - 13: Compute $\widehat{z}_{p,t}^i$, $P_{p,t}^i$, $\widehat{z}_{f,t}^i$ and $P_{f,t}^i$ for $i = 1, \dots, N$.
 - 14: Set $w_t^i \propto p_\theta(y_t|s_t^i, y_{1:t-1})$ for $i = 1, \dots, N$ s.t. $\sum_i w_t^i = 1$.
 - 15: **end for**
 - 16: **for** $t = T$ to 1 **do**
 - 17: Compute $\widehat{z}_{s,t}^i$, $P_{s,t}^i$ for $i = 1, \dots, N$
 - 18: **end for**
 - 19: Set $s_{1:T}[k+1] = s_{1:T}^J$ with $\mathbb{P}(J = j) = w_T^j$
-

This leads to the approximation (cf. (5))

$$\begin{aligned} \widehat{\mathcal{Q}}_k(\theta) &= (1 - \gamma_k) \widehat{\mathcal{Q}}_{k-1}(\theta) + \\ &\gamma_k \sum_{i=1}^N w_T^i \mathbb{E}_{\theta_k} [\log p_\theta(y_{1:T}, z_{1:T}, s_{1:T}^i) | s_{1:T}^i, y_{1:T}], \end{aligned} \quad (13)$$

where the expectation is w.r.t. $z_{1:T}$. Putting this together, we obtain a Rao-Blackwellized PSAEM (RB-PSAEM) algorithm presented in Algorithm 3. Note that this algorithm is similar to the MCMC-based smoother in Algorithm 1, but with the difference that the model parameters are updated at each iteration, effectively enabling simultaneous smoothing and identification.

Algorithm 3 Rao-Blackwellized PSAEM

- 1: Initialize $\widehat{\theta}_0$ and $s_{1:T}[0]$, and $\widehat{\mathcal{Q}}_0(\theta) \equiv 0$.
 - 2: **for** $k \geq 1$ **do**
 - 3: Run Algorithm 2 to obtain $\{s_{1:T}^i, w_T^i\}_{i=1}^N$ and $s_{1:T}[k]$.
 - 4: Compute $\widehat{\mathcal{Q}}_k(\theta)$ according to (13).
 - 5: Compute $\widehat{\theta}_k = \arg \max_{\theta \in \Theta} \widehat{\mathcal{Q}}_k(\theta)$
 - 6: **end for**
-

(For notational convenience, the iteration number k is suppressed in the variables related to $\{s_{1:T}^i, w_T^i\}_{i=1}^N$.)

With a strong theoretical foundation in PMCMC and Markovian stochastic approximation, the RB-PSAEM algorithm presented here enjoys very favourable convergence properties. In particular, under certain smoothness and ergodicity conditions, the sequence of iterates $\{\theta_k\}_{k \geq 1}$ will converge to a maximizer of

$p_\theta(y_{1:T})$ as $k \rightarrow \infty$, regardless of the number of particles $N \geq 2$ used in the internal CPF-AS procedure (see [18, Proposition 1] together with [17] for details). Furthermore, empirically it has been found that a small number of particles can work well in practice as well. For instance, in the numerical examples considered in Section 5, we run Algorithm 3 with $N = 3$ with accurate identification results.

For the model structure (1), there exists infinitely many solutions to the problem (2); all relevant involved matrices can be transformed by a linear transformation matrix and the modes can be re-ordered, but the input-output behaviour will remain invariant. The model is therefore over-parametrized, or lacks identifiability, in the general problem setting. However, it is shown in [24] that the Cramér-Rao Lower Bound is not affected by the over-parametrization. That is, the estimate quality, in terms of variance, is unaffected by the over-parametrization.

4.1 Maximizing the intermediate quantity

When making use of RB-PSAEM from Algorithm 3, one major question arises from Step 5, namely the maximization of the intermediate quantity $\widehat{Q}_k(\theta)$. For the jump Markov linear model, the expectation in (13) can be expressed using sufficient statistics, as will be shown later, as an inner product

$$\sum_{i=1}^N w_T^i \mathbb{E}_{\theta_k} [\log p_\theta(y_{1:T}, z_{1:T}, s_{1:T}^i) | s_{1:T}^i, y_{1:T}] = \langle S^k, \eta(\theta) \rangle, \quad (14)$$

for a sufficient statistics S and corresponding natural parameter $\eta(\theta)$. Hence \widehat{Q}_k can be written as

$$\widehat{Q}_k(\theta) = (1 - \gamma_k) \widehat{Q}_{k-1}(\theta) + \gamma_k \langle S^k, \eta(\theta) \rangle = \langle \mathbb{S}^k, \eta(\theta) \rangle \quad (15)$$

if the transformation

$$\mathbb{S}^k = (1 - \gamma_k) \mathbb{S}^{k-1} + \gamma_k S^k \quad (16)$$

is used. In detail,

$$\begin{aligned} & \sum_{i=1}^N w_T^i \mathbb{E}_{\theta_k} [\log p_\theta(y_{1:T}, z_{1:T}, s_{1:T}^i) | s_{1:T}^i, y_{1:T}] = \\ & \sum_{n=1}^K \sum_{m=1}^K S_{n,m}^{(1)} \log \pi_{n,m} - \sum_{n=1}^K \frac{1}{2} \left(S_n^{(2)} \log(|Q_n| |R_n|) + \text{Tr}(H_n^\theta S_n^{(3)}) \right) \end{aligned} \quad (17a)$$

neglecting constant terms in the last expression. This can be verified to be an inner product (as indicated in (14)) in $S = \{S^{(1)}, S^{(2)}, S^{(3)}\}$. Here the sufficient statistics

$$S_{n,m}^{(1)} = \sum_{i=1}^N w_T^i \sum_{t=1}^T \mathbf{1}(s_t^i = m, s_{t-1}^i = n), \quad (17b)$$

$$S_n^{(2)} = \sum_{i=1}^N w_T^i \sum_{t=1}^T \mathbf{1}(s_t^i = n), \quad (17c)$$

$$S_n^{(3)} = \sum_{i=1}^N w_T^i \sum_{t=1}^T \mathbf{1}(s_t^i = n) (\widehat{\xi}_t^i \widehat{\xi}_t^{i,T} + M_{t|T}^i), \quad (17d)$$

with

$$\widehat{\xi}_t^i = \left(\widehat{z}_{s;t}^{i,T} \left[\widehat{z}_{s;t-1}^{i,T} \ u_{t-1}^T \right] \ y_t^T \left[\widehat{z}_{s;t}^{i,T} \ u_t^T \right] \right)^T, \quad (17e)$$

and

$$H_n^\theta = \begin{pmatrix} [I \ A_n^T \ B_n^T] Q_n^{-1} \begin{bmatrix} I \\ A_n \\ B_n \end{bmatrix} & 0 \\ 0 & [I \ C_n^T \ D_n^T] R_n^{-1} \begin{bmatrix} I \\ C_n \\ D_n \end{bmatrix} \end{pmatrix} \quad (17f)$$

have been used. Further notation introduced is $\mathbf{1}(\cdot)$ as the indicator function, and

$$M_{t|T}^i = \begin{pmatrix} P_{s;t}^i & P_{s;t,t-1}^i & 0 & 0 & P_{s;t}^i & 0 \\ P_{s;t,t-1}^i & P_{s;t-1}^i & 0 & 0 & P_{s;t,t-1}^i & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ P_{s;t}^i & P_{s;t,t-1}^i & 0 & 0 & P_{s;t-1}^i & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}. \quad (17g)$$

For computing this, the RTS-smoother in step 17 in Algorithm 2 has to be extended by calculation of $P_{s;t+1,t} \triangleq \text{Cov}[\hat{z}_{s,t+1}, \hat{z}_{s;t}^T]$, which can be done as follows [31, Property P6.2]

$$P_{s;t,t-1} = P_{f;t} J_{t-1}^T + J_t (P_{s;t+1,t} - A_{t+1} P_{f;t}) J_{t-1}^T, \quad (18)$$

initialized with $P_{T,T-1|T} = (I - K_T C_T) A_T P_{f;t-1}$.

For notational convenience, we will partition $S_n^{(3)}$ as

$$S_n^{(3)} = \begin{pmatrix} \Phi_n & \Psi_n & & & & \\ \Psi_n^T & \Sigma_n & & & & \\ & & \Omega_n & \Lambda_n & & \\ & & \Lambda_n^T & \Xi_n & & \end{pmatrix}. \quad (19)$$

Lemma 1. *Assume for all modes $n = 1, \dots, K$, that all states z are controllable and observable and $\sum_t \mathbf{1}(s_t = n) u_t^T u_t > 0$. The parameters θ maximizing $\hat{Q}_k(\theta)$ for the jump Markov linear model (1) are then given by*

$$\pi_{n,m}^j = \frac{\mathbb{S}_{n,m}^{(1),k}}{\sum_l \mathbb{S}_{n,l}^{(1),k}}, \quad (20a)$$

$$[A_n \ B_n] = \Psi_n \Sigma_n^{-1}, \quad (20b)$$

$$[C_n \ D_n] = \Lambda_n \Xi_n^{-1}, \quad (20c)$$

$$[Q_n] = (\mathbb{S}_n^{(2),k})^{-1} (\Phi_n - \Psi_n \Sigma_n^{-1} \Psi_n^T), \quad (20d)$$

$$[R_n] = (\mathbb{S}_n^{(2),k})^{-1} (\Omega_n - \Lambda_n \Xi_n^{-1} \Lambda_n^T), \quad (20e)$$

for $n, m = 1, \dots, K$.

Φ_n, Ψ_n, \dots are the partitions of $\mathbb{S}_n^{(3),k}$ indicated in (19), and $\mathbb{S}^{(i)}$ are the ‘SA-updates’ (16) of the sufficient statistics (17b)-(17d).

Remark: If $B \equiv 0$, the first square bracket in (17e) can be replaced by $[\hat{z}_{s;t-1}^{i,T}]$, and (20b) becomes $[A_n] = \Psi_n \Sigma_n^{-1}$. The case with $D \equiv 0$ is fully analogous.

Proof. With arguments directly from [14, Lemma 3.3], the maximization of the last part of (17a) for a given $s_t = n$ (for any sufficient statistics Z in the inner product, and in particular $Z = \mathbb{S}^k$), is found to be (20b)-(20e).

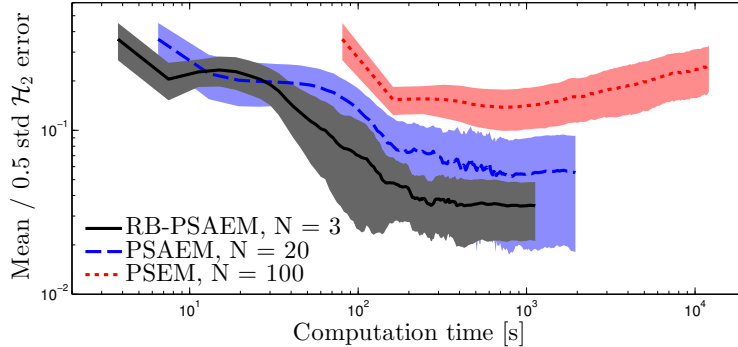


Figure 1: Numerical example 1. Mean (lines) and 0.5 standard deviation (fields) \mathcal{H}_2 error for 7 runs of our RB-PSAEM using $N = 3$ particles (black) PSAEM [18] using $N = 20$ particles (blue) and PSEM [29] using $N = 100$ particles and $M = 20$ backward trajectories (red).

Using Lagrange multipliers and that $\sum_i \pi_{n,m} = 1$, the maximum w.r.t. Π of the first part of (17a) is obtained as

$$\pi_{n,m} = \frac{\mathbb{S}_{n,m}^{(1),k}}{\sum_l \mathbb{S}_{n,l}^{(1),k}}. \quad (21)$$

□

4.2 Computational complexity

Regarding the computational complexity of Algorithm 3, the most important result is that it is linear in the number of measurements T . It is also linear in the number of particles N .

5 Numerical examples

Some numerical examples are given to illustrate the properties of the Rao-Blackwellized PSAEM algorithm. The Matlab code for the examples is available via the homepage of the first author.

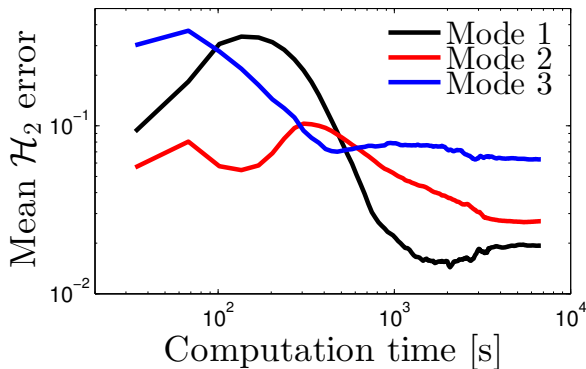
5.1 Example 1 - Comparison to related methods

The first example concerns identification using simulated data ($T = 3000$) for a one-dimensional ($n_z = 1$) jump Markov linear model with 2 modes ($K = 2$) (with parameters randomly generated according to $A_n \sim U_{[-1,1]}$, $B_n \sim U_{[-5,5]}$, $C_n \sim U_{[-5,5]}$, $D_n \equiv 0$, $Q_n \sim U_{[0.01,0.1]}$, $R_n \sim U_{[0.01,0.1]}$) with low-pass filtered white noise as u_t . The following methods are compared:

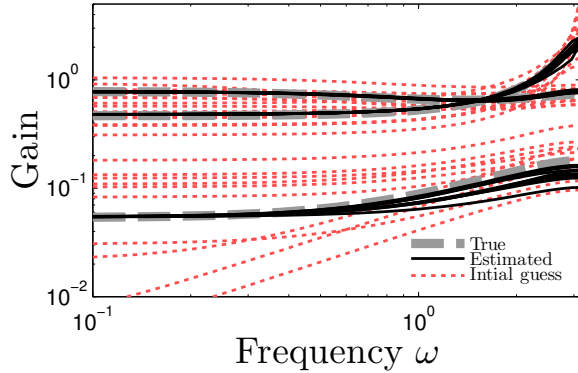
1. RB-PSAEM from Algorithm 3, with (only) $N = 3$ particles,
2. PSAEM as presented in [18] with $N = 20$,
3. PSEM [29] with $N = 100$ forward particles and $M = 20$ backward simulated trajectories.

The initial parameters $\hat{\theta}_0$ are each randomly picked from $[0.5\theta^*, 1.5\theta^*]$, where θ^* is the true parameter value. The results are illustrated in Figure 1, which shows the mean (over all modes and 7 runs) \mathcal{H}_2 error for the transfer function from the input u to the output y .

From Figure 1 (note the log-log scale used in the plot) it is clear that our new Rao-Blackwellized PSAEM algorithm has a significantly better performance, both in terms of mean and in variance between different runs, compared to the previous algorithms.



(a) Mean \mathcal{H}_2 error for each mode.



(b) Bode plots of the estimates (black), true (dashed grey) and the initializations (dotted red).

Figure 2: Plots from Numerical example 2.

5.2 Example 2 - Identification of multidimensional systems

Let us now consider a two-dimensional system ($n_z = 2$) with $K = 3$ modes. The eigenvalues for A_n are randomly picked from $[-1, 1]$. The other parameters are randomly picked as $B_n \sim U_{[-5, 5]}$, $C_n \sim U_{[-5, 5]}$, $D_n \equiv 0$, $Q_n \sim I_2 \cdot U_{[0.01, 0.1]}$, $R_n \sim U_{[0.01, 0.1]}$, and the system is simulated for $T = 8000$ time steps with input u_t being a low-pass filtered white noise. The initialization of the Rao-Blackwellized PSAEM algorithm is randomly picked from $[0.6\theta^*, 1.4\theta^*]$ for each parameter. The number of particles used in the particle filter is $N = 3$. Figure 2a shows the mean (over 10 runs) \mathcal{H}_2 error for each mode, similar to Figure 1. Figure 2b shows the estimated Bode plots after 300 iterations. As is seen from Figure 2b, the RB-PSAEM algorithm has the ability to catch the dynamics of the multidimensional system fairly well.

6 CONCLUSION AND FUTURE WORK

We have derived a maximum likelihood estimator for identification of jump Markov linear models. More specifically an expectation maximization type of solution was derived. The nonlinear state smoothing problem inherent in the expectation step was solved by constructing an ergodic Markov kernel leaving the joint state smoothing distribution invariant. Key to this development was the introduction of a Rao-Blackwellized conditional particle filter with ancestor sampling. The maximization step could be solved in closed form. The experimental results indicate that we obtain significantly better performance both in terms of accuracy and computational time when compared to previous state of the art particle filtering based methods. The ideas underlying the smoother derived in this work have great potential also outside the class of jump Markov linear models and this is something worth more investigation. Indeed, it is quite possible that it can turn out to be a serious competitor also in finding the joint smoothing distribution for general nonlinear state space models.

References

- [1] Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342, 2010.
- [2] Christophe Andrieu, Eric Moulines, and Pierre Priouret. Stability of stochastic approximation under verifiable conditions. *SIAM Journal on control and optimization*, 44(1):283–312, 2005.

- [3] Trevor T Ashley and Sean B Andersson. A sequential monte carlo framework for the system identification of jump markov state space models. In *Proceedings of American Control Conference*, pages 1144–1149, Portland, Oregon, June 2014.
- [4] Alberto Bemporad, Jacob Roll, and Lennart Ljung. Identification of hybrid systems via mixed-integer programming. In *Proceedings of 40th IEEE Conference on Decision and Control*, pages 786–792, Orlando, Florida, 2001.
- [5] Lars Blackmore, Stephanie Gil, Seung Chung, and Brian Williams. Model learning for switching linear systems with autonomous mode transitions. In *Proceedings of 46th IEEE Conference on Decision and Control*, pages 4648–4655, New Orleans, LA, 2007.
- [6] José Borges, Vincent Verdult, Michel Verhaegen, and Miguel Ayala Botto. A switching detection method based on projected subspace classification. In *Proceedings of 44th IEEE Conference on Decision and Control, jointly with European Control Conference*, pages 344–349, Sevilla, Spain, 2005.
- [7] Olivier Cappé, Éric Moulines, and Tobias Rydén. *Inference in Hidden Markov Models*. Springer Series in Statistics. Springer, New York, USA, 2005.
- [8] George Casella and Christian P Robert. Rao-Blackwellisation of sampling schemes. *Biometrika*, 83(1):81–94, 1996.
- [9] Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the EM algorithm. *Annals of Statistics*, 27(1):94–128, 1999.
- [10] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- [11] Arnaud Doucet and Adam M Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. In Dan Crisan and Boris Rozovsky, editors, *Nonlinear Filtering Handbook*, pages 656–704. Oxford University Press, Oxford, 2011.
- [12] Emily Fox, Erik B Sudderth, Michael I Jordan, and Alan Willsky. Bayesian nonparametric inference of switching dynamic linear models. *IEEE Transactions of Signal Processing*, 59(4):1569–1585, April 2011.
- [13] Andrea Garulli, Simone Paoletti, and Antonio Vicino. A survey on switched and piecewise affine system identification. In *Proceedings of 16th IFAC Symposium on System Identification*, volume 16, pages 344–355, Brussels, Belgium, 2012.
- [14] Stuart Gibson and Brett Ninness. Robust maximum-likelihood estimation of multivariable dynamic systems. *Automatica*, 41(10):1667–1682, October 2005.
- [15] Stephanie Gil and Brian Williams. Beyond local optimality: An improved approach to hybrid model learning. In *Proceedings of 48th IEEE Conf Decision and Control, jointly with 28th Chinese Control Conference*, pages 3938–3945, Shanghai, China, 2009.
- [16] Thomas Kailath, Ali H. Sayed, and Babak Hassibi. *Linear estimation*. Prentice Hall, Upper Saddle River, NJ, 2000.
- [17] Estelle Kuhn and Marc Lavielle. Coupling a stochastic approximation version of EM with an MCMC procedure. *ESAIM: Probability and Statistics*, 8:115–131, September 2004.
- [18] Fredrik Lindsten. An efficient stochastic approximation EM algorithm using conditional particle filters. In *Proceedings of 38th IEEE International Conference Acoustics, Speech, and Signal Processing (ICASSP)*, pages 6274–6278, Vancouver, Canada, May 2013.

- [19] Fredrik Lindsten, Michael I Jordan, and Thomas B Schön. Particle Gibbs with ancestor sampling. *Journal of Machine Learning Research*, 15:2145–2184, 2014.
- [20] Fredrik Lindsten and Thomas B Schön. Backward simulation methods for Monte Carlo statistical inference. *Foundations and Trends in Machine Learning*, 6(1):1–143, 2013.
- [21] Jimmy Olsson, Randal Douc, Olivier Cappé, and Éric Moulines. Sequential Monte Carlo smoothing with application to parameter estimation in nonlinear state-space models. *Bernoulli*, 14(1):155–179, 2008.
- [22] Simone Paoletti, Aleksandar Lj Juloski, Giancarlo Ferrari-Trecate, and René Vidal. Identification of hybrid systems a tutorial. *European journal of control*, 13(2):242–260, 2007.
- [23] Komi Midzodzi Pekpe, Gilles Mourot, Komi Gasso, and José Ragot. Identification of switching systems using change detection technique in the subspace framework. In *Proceedings of 43rd IEEE Conference on Decision and Control*, volume 4, pages 3838–3843, Bahamas, 2004.
- [24] Rik Pintelon, Joannes Schoukens, Tomas McKelvey, and Yves Rolain. Minimum variance bounds for overparameterized models. *IEEE Transactions on Automatic Control*, 41(5):719–720, 1996.
- [25] HE Rauch, CT Striebel, and F Tung. Maximum likelihood estimates of linear dynamic systems. *AIAA journal*, 3(8):1445–1450, 1965.
- [26] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- [27] Christian P. Robert and George Casella. *Monte Carlo statistical methods*. Springer, New York, 2. ed. edition, 2004.
- [28] Thomas B Schön, Fredrik Gustafsson, and P-J Nordlund. Marginalized particle filters for mixed linear/nonlinear state-space models. *IEEE Transactions on Signal Proceedings*, 53(7):2279–2289, 2005.
- [29] Thomas B Schön, Adrian Wills, and Brett Ninness. System identification of nonlinear state-space models. *Automatica*, 47(1):39–49, January 2011.
- [30] Robert H Shumway and David S Stoffer. An approach to time series smoothing and forecasting using the EM algorithm. *Journal of Time Series Analysis*, 3(4):253–264, 1982.
- [31] Robert H Sumway and David S Stoffer. *Time series analysis and its applications with R examples*. Springer, New York, 2006.
- [32] Greg CG Wei and Martin A Tanner. A Monte Carlo implementation of the EM algorithm and the poor man’s data augmentation algorithms. *Journal of the American Statistical Association*, 85(411):699–704, 1990.
- [33] Nick Whiteley, Christophe Andrieu, and Arnaud Doucet. Efficient Bayesian inference for switching state-space models using discrete particle Markov chain Monte Carlo methods. *arXiv preprint arXiv:1011.2437*, 2010.
- [34] Sinan Yildirim, Sumeetpal S Singh, and Arnaud Doucet. An online expectation–maximization algorithm for changepoint models. *Journal of Computational and Graphical Statistics*, 22(4):906–926, 2013.