

# Recursive Estimation in Hidden Markov Models\*

François LeGland and Laurent Mevel  
IRISA / INRIA and IRMAR  
Campus de Beaulieu  
35042 Rennes Cédex, France  
e-mail : {legland,lmevel}@irisa.fr

## Abstract

We consider a hidden Markov model (HMM) with multidimensional observations, and where the coefficients (transition probability matrix, and observation conditional densities) depend on some unknown parameter. We study the asymptotic behaviour of two recursive estimators, the recursive maximum likelihood estimator (RMLE), and the recursive conditional least squares estimator (RCLSE), as the number of observations increases to infinity. Firstly, we exhibit the contrast functions associated with the two non-recursive estimators, and we prove that the recursive estimators converge a.s. to the set of stationary points of the corresponding contrast function. Secondly, we prove that the two recursive estimators are asymptotically normal.

## 1 Introduction

We consider the problem of identification of a partially observed finite-state Markov chain, based on observations in  $\mathbf{R}^d$ . We study the asymptotic behaviour (convergence and asymptotic normality) of two recursive estimators, the recursive maximum likelihood estimator (RMLE), and the recursive conditional least squares estimator (RCLSE), as the number of observations increases to infinity. Assuming stationarity of the nonobserved Markov chain, the non-recursive MLE has been studied by Petrie [10] in the case of observations in a finite set, and by Leroux [9], and Bickel, Ritov and Rydén [3], in the case of observations in  $\mathbf{R}^d$ . These results have been obtained in LeGland and Mevel [8] without any stationarity assumption. The convergence of the RCLSE has been proved in Arapostathis and Marcus [1] for a special case of observations in a finite set — see also Le Gland and Mevel [6] for a study of the RMLE in the general case of observations in a finite set.

\*This work was partially supported by the Commission of the European Communities, under the SCIENCE project *System Identification*, project number SC1\*-CT92-0779, and under the HCM project *Statistical Inference for Stochastic Processes*, project number CHRX-CT92-0078, and by the Army Research Office, under grant DAAH04-95-1-0164.

The convergence and the asymptotic normality have been obtained by Rydén in [11] for the recursive maximum split data likelihood estimator (MSDLE) defined therein, under stationarity assumption. In this paper, we prove the convergence and the asymptotic normality of the RMLE and the RCLSE, in the case of observations in  $\mathbf{R}^d$ , without any stationary assumption, using geometric ergodicity results on the *approximate* prediction filter obtained in LeGland and Mevel [7]. We prove also that the RMLE is asymptotically efficient, i.e. its asymptotic covariance matrix is the Fisher information matrix.

### 1.1 Statistical model

Let  $\{X_n, n \geq 0\}$  and  $\{Y_n, n \geq 0\}$  be two sequences, defined on the probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ , with values in the finite set  $S = \{1, \dots, N\}$  and in  $\mathbf{R}^d$  respectively. On the space  $(\Omega, \mathcal{F})$  we consider a family  $(\mathbf{P}^\theta, \theta \in \Theta)$  of probability measures, with  $\Theta$  convex compact subset of  $\mathbf{R}^p$ , such that under  $\mathbf{P}^\theta$ :

- The unobserved state sequence  $\{X_n, n \geq 0\}$  is a Markov chain with *primitive* transition probability matrix  $Q_\theta = (q_\theta^{i,j})$ , i.e. for any  $i, j \in S$

$$q_\theta^{i,j} = \mathbf{P}^\theta[X_{n+1} = j \mid X_n = i],$$

and initial probability distribution  $p_0 = (p_0^i)$  independent of  $\theta \in \Theta$ , and possibly different of the *true* initial probability distribution  $p_\bullet = (p_\bullet^i)$  of  $X_0$ , i.e. for any  $i \in S$

$$p_0^i = \mathbf{P}^\theta[X_0 = i] \neq \mathbf{P}[X_0 = i] = p_\bullet^i.$$

- The observations  $\{Y_n, n \geq 0\}$  are mutually independent given the sequence of states of the Markov chain, i.e.

$$\mathbf{P}^\theta[Y_n \in dy_n, \dots, Y_0 \in dy_0 \mid X_n = i_n, \dots, X_0 = i_0]$$

$$= \prod_{k=0}^n \mathbf{P}^\theta[Y_k \in dy_k \mid X_k = i_k].$$

For any  $n \geq 0$ , and for any  $i \in S$ , the conditional probability distribution of the observation

$Y_n$  given that  $(X_n = i)$ , is absolutely continuous with respect to a positive and  $\sigma$ -finite measure  $\lambda$  on  $\mathbf{R}^d$ , i.e.

$$\mathbf{P}^\theta[Y_n \in dy \mid X_n = i] = b_\theta^i(y) \lambda(dy),$$

with a  $\lambda$ -a.e. positive density. For any  $y \in \mathbf{R}^d$ , let

$$b_\theta(y) = [b_\theta^1(y), \dots, b_\theta^N(y)]^*,$$

$$B_\theta(y) = \text{diag}[b_\theta^1(y), \dots, b_\theta^N(y)].$$

Here and throughout the paper, the notation  $*$  denotes the transpose of a matrix.

**Example 1.1** [conditionally Gaussian observations] Assume that for any  $\theta \in \Theta$ , the observations are of the form

$$Y_n = h_\theta(X_n) + V_n^\theta,$$

for all  $n \geq 0$ , where  $\{V_n^\theta, n \geq 0\}$  is a Gaussian white noise sequence under  $\mathbf{P}^\theta$ , with identity covariance matrix. The mapping  $h_\theta$  from  $S$  to  $\mathbf{R}^d$  is equivalently defined as  $h_\theta = (h_\theta^i)$  where  $h_\theta^i \in \mathbf{R}^d$  for all  $i \in S$ . In this case, the mutual independence condition is satisfied, and

$$b_\theta^i(y) = (2\pi)^{-d/2} \exp\left\{-\frac{1}{2}\|y - h_\theta^i\|^2\right\},$$

for any  $i \in S$ . Here and throughout the paper, the notation  $\|\cdot\|$  denotes the Euclidean norm.

Throughout the paper, the *true* value of the parameter will be denoted by  $\alpha$ , and we make the following assumption :

**Assumption A :** For the *true* value  $\alpha \in \Theta$ , the transition probability matrix  $Q_\alpha = (q_\alpha^{i,j})$  is a positive matrix, i.e.  $q_\alpha^{i,j} \geq \varepsilon$ , for all  $i, j \in S$ , and for some *known*  $\varepsilon > 0$ .

**Assumption A' :** The mapping  $\theta \mapsto Q_\theta$  is two-times differentiable, with bounded first and second derivatives, and Lipschitz continuous second derivative.

**Assumption B :** For any  $y \in \mathbf{R}^d$ , the mapping  $\theta \mapsto b_\theta(y)$  is three-times differentiable, and we define

$$\delta^{(s)}(y) = \max_{\theta \in \Theta} \max_{k_1, \dots, k_s=1, \dots, p} \frac{\max_{i \in S} |\partial_{k_1, \dots, k_s}^s b_\theta^i(y)|}{\min_{i \in S} b_\theta^i(y)}.$$

**Definition 1.2** For any  $p \geq 0$ , and any  $s = 0, 1, 2, 3$ ,

$$\Delta_p^{(s)} = \max_{\theta \in \Theta} \max_{i \in S} \int_{\mathbf{R}^d} [\delta^{(s)}(y)]^p b_\theta^i(y) \lambda(dy),$$

$$\Gamma_p = \max_{i \in S} \max_{\theta \in \Theta} \int_{\mathbf{R}^d} [\max_{j \in S} |\log b_\theta^j(y)|]^p b_\theta^i(y) \lambda(dy),$$

$$\bar{Y}^p = \max_{\theta \in \Theta} \max_{i \in S} \int_{\mathbf{R}^d} |y|^p b_\theta^i(y) \lambda(dy).$$

We use also the alternate notation  $\Delta_p = \Delta_p^{(0)}$ ,  $\Delta'_p = \Delta_p^{(1)}$ ,  $\Delta''_p = \Delta_p^{(2)}$ , and  $\Delta'''_p = \Delta_p^{(3)}$ .

**Assumption B' :**  $\bar{Y}^2$ ,  $\Delta'_2$ ,  $\Delta''_2$ , and  $\Delta'''_2$  are finite.

## 1.2 Prediction filters

For all  $n \geq 1$ , let  $p_n^*$  (= $p_n^i$ ) denote the *prediction filter*, i.e. the conditional probability distribution under  $\mathbf{P}$  of the state  $X_n$  given observations  $(Y_0, \dots, Y_{n-1})$  : for any  $i \in S$

$$p_n^i = \mathbf{P}[X_n = i \mid Y_0, \dots, Y_{n-1}].$$

The random sequence  $\{p_n^*, n \geq 0\}$  takes values in the set  $\mathcal{P}(S)$  of probability distributions over the finite set  $S$ , and satisfies the forward Baum equation

$$p_{n+1}^* = \frac{Q_\alpha^* B_\alpha(Y_n) p_n^*}{b_\alpha^*(Y_n) p_n^*},$$

for all  $n \geq 0$ . Notice that the initial condition  $p_0^*$  is the probability distribution  $p_\bullet$  of  $X_0$ .

The true value of the parameter, and the initial probability distribution of  $X_0$  are unknown, and for any  $\theta \in \Theta$  we consider the following equation for the corresponding *approximate* prediction filter

$$p_{n+1}^\theta = \frac{Q_\theta^* B_\theta(Y_n) p_n^\theta}{b_\theta^*(Y_n) p_n^\theta} \triangleq f_\theta[Y_n, p_n^\theta], \quad (1)$$

for all  $n \geq 0$ , where  $Q_\theta = (q_\theta^{i,j})$  is the stochastic matrix generated by the parameter  $\theta$ .

Differentiating (1) w.r.t. the  $k$ -th component of the  $p$ -dimensional parameter  $\theta$  yields

$$\partial_k p_{n+1}^\theta = \Phi_\theta[Y_n, p_n^\theta] \partial_k p_n^\theta + \partial_k f_\theta[Y_n, p_n^\theta],$$

where for any  $y \in \mathbf{R}^d$ , and any  $p \in \mathcal{P}(S)$

$$\begin{aligned} \Phi_\theta[y, p] &= \frac{Q_\theta^* B_\theta(y)}{b_\theta^*(y) p} \left[ I - \frac{p b_\theta^*(y)}{b_\theta^*(y) p} \right] \\ &= Q_\theta^* \left[ I - \frac{B_\theta(y) p e^*}{b_\theta^*(y) p} \right] \frac{B_\theta(y)}{b_\theta^*(y) p}, \end{aligned}$$

$$\begin{aligned} \partial_k f_\theta[y, p] &= Q_\theta^* \left[ I - \frac{B_\theta(y) p e^*}{b_\theta^*(y) p} \right] \frac{\partial_k B_\theta(y) p}{b_\theta^*(y) p} \\ &\quad + \frac{\partial_k Q_\theta^* B_\theta(y) p}{b_\theta^*(y) p}. \end{aligned}$$

The sequence  $\partial p_n^\theta = (\partial_k p_n^\theta)$  belongs to  $\Sigma^p$ , where

$$\Sigma = \{w \in \mathbf{R}^N : e^* w = 0\},$$

and  $e = (1, \dots, 1)^*$ .

It has been proved in LeGland and Mevel [7], that under the probability measure corresponding to the true value  $\alpha$  of the parameter, the extended Markov chain  $\{Z_n^\theta = (X_n, Y_n, p_n^\theta, \partial p_n^\theta), n \geq 0\}$  with values in  $E \triangleq S \times \mathbf{R}^d \times \mathcal{P}(S) \times \Sigma^p$ , is geometrically ergodic under suitable integrability assumptions. The existence and uniqueness of the invariant measure  $\mu_\theta$  for the extended Markov chain  $\{Z_n^\theta, n \geq 0\}$  is used below to prove the convergence of the two recursive estimators.

## 2 Recursive identification algorithm

In LeGland and Mevel [8], we have proved the consistency and the asymptotic normality of the non-recursive MLE and CLSE, under suitable integrability conditions stated therein. In this section, we will present the recursive estimators, and prove their consistency.

**Definition 2.1** Let  $L_\Theta$  denote the set of functions  $g_\theta = (g_\theta^k)$  defined on  $\mathbf{R}^d \times \mathcal{P}(S) \times \Sigma^p$ , such that for any  $k = 1, \dots, p$ , and any  $y \in \mathbf{R}^d$ , the partial mapping  $(p, w) \mapsto g_\theta^k(y, p, w)$  is locally Lipschitz continuous, in the sense that

$$|g_\theta^k(y, p, w) - g_\theta^k(y, p', w')| \leq \text{Lip}(g, y) [\|w_k - w'_k\| + \|p - p'\| (1 + \|w_k\| + \|w'_k\|)]$$

for any  $p, p' \in \mathcal{P}(S)$ , and any  $w = (w_k), w' = (w'_k) \in \Sigma^p$ , and

$$|g_\theta^k(y, p, w)| \leq K(g, y) (1 + \|w_k\|),$$

for any  $p \in \mathcal{P}(S)$ , and any  $w = (w_k) \in \Sigma^p$ .

Moreover, for any  $k = 1, \dots, p$ , any  $y \in \mathbf{R}^d$ , any  $p \in \mathcal{P}(S)$ , and any  $w = (w_k) \in \Sigma^p$ , the partial mapping  $\theta \mapsto g_\theta^k(y, p, w)$  is locally Lipschitz continuous, in the sense that for any  $\theta, \theta' \in \Theta_\varepsilon$ ,

$$|g_\theta^k(y, p, w) - g_{\theta'}^k(y, p, w)| \leq \text{Lip}_\Theta(g, y) (1 + \|w_k\|) \|\theta - \theta'\|.$$

With the above notations, we consider the recursive algorithm defined as follows : for all  $k = 1, \dots, p$

$$\hat{\theta}_{n+1}^k = \pi_\varepsilon(\hat{\theta}_n^k + \gamma_{n+1} \Psi_n(H_{\hat{\theta}_n^k}^k(Y_n, \hat{p}_n, \hat{w}_n))), \quad (2)$$

$$\bar{\theta}_{n+1} = \bar{\theta}_n + \frac{1}{n} (\hat{\theta}_{n+1} - \bar{\theta}_n),$$

where  $\gamma_n = n^{-2/3}$ ,  $\pi_\varepsilon(\cdot)$  is the projection on the set  $\Theta_\varepsilon = \{\theta \in \Theta : q_\theta^{i,j} \geq \varepsilon, \text{ for all } i, j \in S\}$ ,  $\Psi_n(\cdot)$  is a projection on a ball of growing radius (hence made only a finite number of times), and

$$\hat{p}_{n+1} = f_{\hat{\theta}_{n+1}}[Y_n, \hat{p}_n],$$

$$\hat{w}_{n+1}^k = \Phi_{\hat{\theta}_{n+1}}[Y_n, \hat{p}_n] \hat{w}_n^k + \partial_k f_{\hat{\theta}_{n+1}}[Y_n, \hat{p}_n].$$

We define  $\hat{Z}_n = (X_n, Y_n, \hat{p}_n, \hat{w}_n)$ . Then

$$\mathbf{P}[\hat{Z}_{n+1} \in B \mid \hat{\theta}_0, \dots, \hat{\theta}_{n+1}, \hat{Z}_0, \dots, \hat{Z}_n] = \Pi_{\hat{\theta}_{n+1}}(\hat{Z}_n, B),$$

i.e. the algorithm (2) belongs to the class of stochastic algorithms with Markovian dynamics, see Benveniste, Métivier and Priouret [2].

### 2.1 The recursive MLE

It has been proved in LeGland and Mevel [8] that the log-likelihood function (suitably normalized) for the estimation of the parameter  $\theta$  based on observations  $(Y_0, \dots, Y_n)$  can be expressed as an additive functional of the extended Markov chain  $\{Z_n^\theta, n \geq 0\}$  as follows

$$\ell_n(\theta) = \frac{1}{n} \sum_{k=0}^n \log[b_\theta^*(Y_k) p_k^\theta].$$

In addition, we have proved that the following strong law of large numbers holds :

**Proposition 2.2** Under Assumption A, and if  $\Delta_1$  and  $\Gamma_1$  are finite, then for any  $\theta \in \Theta$  there exists a finite constant  $\ell(\theta)$  such that

$$\ell_n(\theta) \longrightarrow \ell(\theta), \quad \mathbf{P}\text{-a.s.}$$

as  $n \rightarrow \infty$ , where

$$\ell(\theta) = \int_{\mathbf{R}^d \times \mathcal{P}(S)} \log[b_\theta^*(y) p] \nu_\theta(dy, dp),$$

and where  $\nu_\theta$  denotes the marginal on  $\mathbf{R}^d \times \mathcal{P}(S)$  of the invariant measure  $\mu_\theta$ .

Under the above assumptions, we define for any  $\theta \in \Theta$  the Kullback Leibler information

$$K_{\text{ML}}(\theta) = -[\ell(\theta) - \ell(\alpha)] \geq 0,$$

and we have proved that the true value  $\alpha$  of the parameter belongs to the set  $M_{\text{ML}}$  of global minima of  $K_{\text{ML}}$

$$M_{\text{ML}} = \underset{\theta \in \Theta}{\text{argmin}} K_{\text{ML}}(\theta) \supseteq \{\alpha\}.$$

The maximum likelihood estimator (MLE) is defined as

$$\hat{\theta}_n^{\text{ML}} \in \underset{\theta \in \Theta_\varepsilon}{\text{argmax}} \ell_n(\theta).$$

**Theorem 2.3** Under Assumptions A, and B, and if  $\Delta_2, \Delta'_2, \Gamma_1$  and  $\Gamma_2$  are finite, then any MLE sequence  $\{\hat{\theta}_n^{\text{ML}}, n \geq 0\}$  converges in probability to the deterministic set  $M_{\text{ML}}$  as  $n \rightarrow \infty$ .

The recursive maximum likelihood estimator (RMLE)  $\hat{\theta}_n^{\text{RML}} = (\hat{\theta}_n^k)$ , is defined by the following recursive algorithm : for all  $k = 1, \dots, p$

$$\hat{\theta}_{n+1}^k = \pi_\varepsilon(\hat{\theta}_n^k + \gamma_{n+1} \Psi_n(H_{\hat{\theta}_n^k}^k(Y_n, \hat{p}_n, \hat{w}_n))),$$

where for any  $\theta \in \Theta$ , and any  $k = 1, \dots, p$ , the function  $H_\theta^k$  is defined by

$$H_\theta^k(y, p, w) = \frac{b_\theta^*(y) w_k}{b_\theta^*(y) p} + \frac{\partial_k b_\theta^*(y) p}{b_\theta^*(y) p},$$

for any  $y \in \mathbf{R}^d$ , any  $p \in \mathcal{P}(S)$ , and any  $w = (w_k) \in \Sigma^p$ , and belongs to the set  $L_\Theta$ .

For any  $\theta \in \Theta$ , and any  $k = 1, \dots, p$ , we define

$$\begin{aligned} h_{\text{ML}}^k(\theta) &= \int_{\mathbf{R}^d \times \mathcal{P}(S) \times \Sigma^p} H_\theta^k(y, p, w) \lambda_\theta(dy, dp, dw) \\ &= \partial_k K_{\text{ML}}(\theta), \end{aligned}$$

where  $\lambda_\theta$  denotes the marginal on  $\mathbf{R}^d \times \mathcal{P}(S) \times \Sigma^p$  of the invariant measure  $\mu_\theta$ . We define the set of stationary points of the Kullback–Leibler information  $K_{\text{ML}}$

$$L_{\text{ML}} = \{\theta \in \Theta : h_{\text{ML}}(\theta) = 0\}.$$

Under the following two additional assumptions :

**Assumption C :** The function  $K_{\text{ML}}$  is of class  $C^p$ .

**Assumption D :** There exists  $\varepsilon'$ , such that  $\varepsilon > \varepsilon' > 0$  and :

- (i)  $L_{\text{ML}} \subset \Theta_\varepsilon$ ,
- (ii) For any  $\theta \in \Theta_{\varepsilon'} \setminus \Theta_\varepsilon$

$$h_{\text{ML}}^*(\theta) (\pi_\varepsilon(\theta) - \theta) \geq \delta |\pi_\varepsilon(\theta) - \theta|,$$

for some  $\delta > 0$ .

the following result holds, which is based on Delyon [4], and Delyon and Iouditski [5].

**Theorem 2.4** Under Assumptions A, A', B, B', C, and D, and if  $\Delta_8, \Delta'_4, \Delta''_2, \Delta'''_1$ , and  $\Gamma_2$  are finite, then the RMLE sequence  $\{\hat{\theta}_n^{\text{RML}}, n \geq 0\}$  converges  $\mathbf{P}$ -a.s. to the deterministic set  $L_{\text{ML}}$  as  $n \rightarrow \infty$ , and assuming that  $\hat{\theta}_n^{\text{RML}}$  converges to  $\alpha$ , then

$$\sqrt{n}(\hat{\theta}_n^{\text{RML}} - \alpha) \Rightarrow \mathcal{N}(0, I_F^{-1}),$$

where

$$I_F = \partial^2 K_{\text{ML}}(\alpha),$$

is the Fisher information matrix.

## 2.2 The recursive CLSE

We define the following (suitably normalized) prediction error criterion

$$e_n(\theta) \triangleq \frac{1}{2n} \sum_{k=1}^n |Y_k - \mathbf{E}^\theta[Y_k | \mathcal{Y}_{k-1}]|^2.$$

Introducing the  $N \times d$  matrix  $\phi_\theta = (\phi_\theta^i)$ , where for any  $i \in S$ , the mean vector  $\phi_\theta^i \in \mathbf{R}^d$  is defined by

$$\phi_\theta^i = \int_{\mathbf{R}^d} y b_\theta^i(y) \lambda(dy),$$

(assuming the integral exist), the prediction error criterion can be expressed as an additive functional of the extended Markov chain  $\{Z_n^\theta, n \geq 0\}$  as follows

$$e_n(\theta) = \frac{1}{2n} \sum_{k=1}^n |Y_k - \phi_\theta^* p_k^\theta|^2.$$

In addition, we have proved in LeGland and Mevel [8] that the following law of large numbers holds :

**Proposition 2.5** Under Assumptions A, B, and B', and if  $\Delta_1$  and  $\bar{Y}^2$  are finite, then for any  $\theta \in \Theta$  there exists a finite constant  $e(\theta)$  such that

$$e_n(\theta) \longrightarrow e(\theta), \quad \mathbf{P}\text{-a.s.}$$

as  $n \rightarrow \infty$ , where

$$e(\theta) = \frac{1}{2} \int_{\mathbf{R}^d \times \mathcal{P}(S)} |y - \phi_\theta^* p|^2 \nu_\theta(dy, dp),$$

and where  $\nu_\theta$  denotes the marginal on  $\mathbf{R}^d \times \mathcal{P}(S)$  of the invariant measure  $\mu_\theta$ .

Under the above assumptions, we define for any  $\theta \in \Theta$  the prediction error contrast function

$$K_{\text{CLS}}(\theta) = e(\theta) - e(\alpha) \geq 0,$$

and we have proved that the true value  $\alpha$  of the parameter belongs to the set  $M_{\text{CLS}}$  of global minima of  $K_{\text{CLS}}$

$$M_{\text{CLS}} = \underset{\theta \in \Theta}{\operatorname{argmin}} K_{\text{CLS}}(\theta) \supseteq \{\alpha\}.$$

The conditional least square estimator (CLSE) is defined as

$$\hat{\theta}_n^{\text{CLS}} \in \underset{\theta \in \Theta_\varepsilon}{\operatorname{argmin}} e_n(\theta).$$

**Theorem 2.6** Under Assumptions A, B, and B', and if  $\bar{Y}^4$  is finite, then any CLSE sequence  $\{\hat{\theta}_n^{\text{CLS}}, n \geq 0\}$  converges in probability to the deterministic set  $M_{\text{CLS}}$  as  $n \rightarrow \infty$ .

The recursive conditional least squares estimator (RCLSE)  $\hat{\theta}_n^{\text{RCLS}} = (\hat{\theta}_n^k)$ , is defined by the following recursive algorithm : for any  $k = 1, \dots, p$

$$\hat{\theta}_{n+1}^k = \pi_\varepsilon(\hat{\theta}_n^k + \gamma_{n+1} \Psi_n(H_{\hat{\theta}_n^k}^k(Y_n, \hat{p}_n, \hat{w}_n))),$$

where for any  $\theta \in \Theta$ , and any  $k = 1, \dots, p$ , the function  $H_\theta^k$  is defined by

$$H_\theta^k(y, p, w) = [\phi_\theta(y - \phi_\theta^* p)]^* w_k + [\partial_k \phi_\theta(y - \phi_\theta^* p)]^* p.$$

for any  $y \in \mathbf{R}^d$ , any  $p \in \mathcal{P}(S)$ , and any  $w = (w_k) \in \Sigma^p$ , and belongs to the set  $L_\Theta$ .

For any  $\theta \in \Theta$ , and any  $k = 1, \dots, p$ , we define

$$\begin{aligned} h_{\text{CLS}}^k(\theta) &= \int_{\mathbf{R}^d \times \mathcal{P}(S) \times \Sigma^p} H_\theta^k(y, p, w) \lambda_\theta(dy, dp, dw) \\ &= \partial_k K_{\text{CLS}}(\theta), \end{aligned}$$

where  $\lambda_\theta$  denotes the marginal on  $\mathbf{R}^d \times \mathcal{P}(S) \times \Sigma^p$  of the invariant measure  $\mu_\theta$ . We define the set of stationary points of the prediction error contrast function  $K_{\text{CLS}}$

$$L_{\text{CLS}} = \{\theta \in \Theta : h_{\text{CLS}}(\theta) = 0\}.$$

Adapting Assumptions C and D to the case of the RCLSE, we obtain the counterpart of Theorem 2.4.

**Theorem 2.7** *Under Assumptions A, B, B', C, and D, and if  $\Delta_8, \Delta'_4, \Delta''_2, \Delta'''_1$ , and  $\bar{Y}^4$  are finite, then the RCLSE sequence  $\{\hat{\theta}_n^{\text{RCLS}}, n \geq 0\}$  converges P-a.s. to the deterministic set  $L_{\text{CLS}}$  as  $n \rightarrow \infty$ , and assuming that  $\hat{\theta}_n^{\text{RCLS}}$  converges to  $\alpha$ , then*

$$\sqrt{n}(\hat{\theta}_n^{\text{RCLS}} - \alpha) \Rightarrow \mathcal{N}(0, V^{-1}),$$

where

$$V = \partial^2 K_{\text{CLS}}(\alpha) R^{-1} \partial^2 K_{\text{CLS}}(\alpha),$$

and

$$R = \lim_{n \rightarrow \infty} n \mathbf{E}[\partial e_n(\alpha) \partial e_n^*(\alpha)].$$

### 3 Numerical example

Our example consists in a three-state Markov chain observed in white noise, i.e.

$$Y_n = h(X_n) + W_n,$$

where  $\{W_n, n \geq 0\}$  is an  $\mathcal{N}(0, 1)$  i.i.d. sequence, and

$$h = (5 \ 3 \ 7), \quad Q = \begin{pmatrix} 0.85 & 0.08 & 0.07 \\ 0.4 & 0.5 & 0.1 \\ 0.35 & 0.05 & 0.6 \end{pmatrix}.$$

The initial estimates of the transition probability matrix entries are set uniformly at  $1/N$ , and the initial estimate of the vector  $h$  are randomly chosen in an appropriate compact set, depending on prior knowledge about the true value of the parameter. Firstly, assume we know  $h$  exactly. In Figure 1 and Figure 2, we show the respective performance of both estimators. We display the evolution of both estimators, on the basis of the same observation sequence (1 unit = 10 observations). The theoretical better rate of convergence of

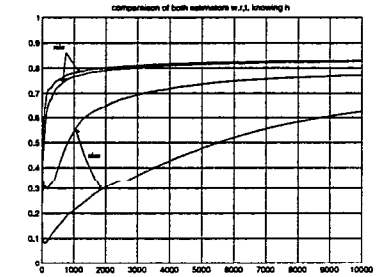
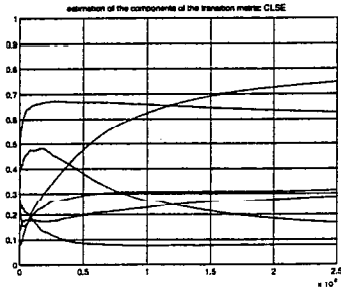
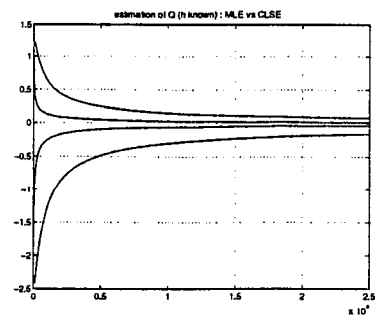
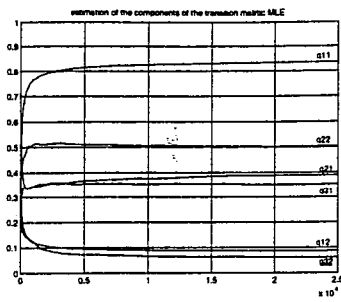
the RMLE over the RCLSE is obviously demonstrated, by comparison of the evolution of both estimators. In Figure 3, we construct, for each estimator, an empirical confidence tube, based on a 50-experiment average. The tube is obtained by using the estimates of the covariance matrix of each estimator. The tube is centered on the difference between  $\theta_n$  and  $\alpha$ . The fact that the zero value belongs to both tubes confirms the convergence of both estimator to  $\alpha$ . The diameter of both tubes goes to zero, with a rate which is proportionnal to the estimate of the corresponding asymptotic covariance matrix. The RMLE tube is much thinner than the RCLSE tube, which shows the asymptotic covariance matrix of the RMLE to be the smallest.

Secondly, we relax the assumed knowledge about  $h$ . In the experience displayed in Figure 4, we investigate the effect of knowing the vector  $h$  or not, on the estimation of the entries of the transition matrix, for both estimators. In the case of the RCLSE, the lack of knowledge about  $h$  results in a clear degradation of the performance of the estimator. On the contrary, the RMLE is less sensible to the lack of knowledge about  $h$ , and seems to be a more robust algorithm in this respect (1 unit = 10 observations). Notice that the fast estimation of  $h$  by the RMLE as seen in Figure 5, illustrates the robustness of the RMLE (compare with the very slow estimation of  $h$  by the RCLSE). More precisely, using the same confidence tube construction as above, we compare the behaviour of the RMLE, for the estimation of  $Q$ . Zero belongs to both confidence tubes, and their diameters go to zero. In that worst situation, the RMLE seems more robust, see Figure 6. On the contrary, the behaviour of the RCLSE is less robust.

In Figure 7, we compare the adaptive prediction filter  $\hat{p}_n$  and the prediction filter  $p_n^\alpha$  for the true value  $\alpha$  of the parameter, and we see that the difference between both probability distributions goes to zero as  $n$  tends to infinity.

### References

- [1] A. ARAPOSTATHIS and S.I. MARCUS. Analysis of an identification algorithm arising in the adaptive estimation of Markov chains. *Mathematics of Control, Signals, and Systems*, 3(1):1-29, 1990.
- [2] A. BENVENISTE, M. MÉTIVIER, and P. PRIOURET. *Adaptive Algorithms and Stochastic Approximations*, volume 22 of *Applications of Mathematics*. Springer Verlag, New York, 1990.
- [3] P.J. BICKEL, Y. RITOV, and T. RYDÉN. Asymptotic normality of the maximum-likelihood estimator for general hidden Markov models. Research Report 1997:2, Department of Mathematical Statistics, Lund University, 1997.



[4] B. DELYON. General results on the convergence of stochastic algorithms. *IEEE Transactions on Automatic Control*, AC-41(9):1245–1255, September 1996.

[5] B. DELYON and A. IOUDITSKI. Stochastic approximation with averaging. Publication Interne 952, IRISA, October 1995. <ftp://ftp.irisa.fr/techreports/1995/PI-952.ps.gz>.

[6] F. LE GLAND and L. MEVEL. Recursive identification of HMM's with observations in a finite set. In *Proceedings of the 34th Conference on Decision and Control, New Orleans 1995*, pages 216–221. IEEE-CSS, December 1995.

[7] F. LE GLAND and L. MEVEL. Geometric ergodicity in hidden Markov models. Publication Interne 1028, IRISA, July 1996. <ftp://ftp.irisa.fr/techreports/1996/PI-1028.ps.gz>.

[8] F. LE GLAND and L. MEVEL. Asymptotic behaviour of the maximum likelihood estimator and the conditional least squares estimator in hidden Markov models. Publication interne, IRISA, 1997. (to appear).

[9] B.G. LEROUX. Maximum-likelihood estimation for hidden Markov models. *Stochastic Processes and their Applications*, 40(1):127–143, 1992.

[10] T. PETRIE. Probabilistic functions of finite state Markov chains. *The Annals of Mathematical Statistics*, 40(1):97–115, 1969.

[11] T. RYDÉN. On recursive estimation for hidden Markov models. *Stochastic Processes and their Applications*, 66(1):79–96, 1997.

