

CONVERGENCE OF A STOCHASTIC APPROXIMATION VERSION OF THE EM ALGORITHM

BY BERNARD DELYON, MARC LAVIELLE AND ERIC MOULINES

*IRISA/INRIA, Université Paris V and Université Paris-Sud and
Ecole National Supérieure des Télécommunications*

The expectation-maximization (EM) algorithm is a powerful computational technique for locating maxima of functions. It is widely used in statistics for maximum likelihood or maximum a posteriori estimation in incomplete data models. In certain situations, however, this method is **not applicable** because the expectation step cannot be performed in closed form. To deal with these problems, a novel method is introduced, the stochastic approximation EM (SAEM), which replaces the expectation step of the EM algorithm by **one iteration of a stochastic approximation procedure**. The convergence of the SAEM algorithm is established under conditions that are applicable to many practical situations. Moreover, it is proved that, under mild additional conditions, the attractive stationary points of the SAEM algorithm correspond to the local maxima of the function. presented to support our findings.

1. Introduction. The EM algorithm [Dempster, Laird and Rubin (1977)] is a very popular tool for maximum-likelihood (or maximum a posteriori) estimation. The common strand to problems where this approach is applicable is a notion of *incomplete data*, which includes the conventional sense of missing data but is much broader than that. The EM algorithm demonstrates its strength in situations where some hypothetical experiment yields (complete) data that are related to the parameters more conveniently than the measurements are. Problems where the EM algorithm has proven to be useful include, among many others, maximum likelihood estimation of the parameters of mixture of densities [see, e.g., Titterington, Smith and Makov (1985)] and of hidden Markov models [MacDonald and Zucchini (1997) and the references therein], maximum a posteriori estimation in censored data models [Little and Rubin (1987), Tanner (1993)]. The EM algorithm has several appealing properties. Because it relies on complete-data computations, it is generally simple to implement: the **E-step** of each iteration only involves taking expectation over complete-data conditional distribution; the **M-step** only involves complete data maximum-likelihood estimation, which is often in simple closed form. Moreover, it is numerically stable, in the sense that it increases the incomplete likelihood at each iteration. When either the maximization step (M-step) or the expectation step (E-step) involves intricate or even infeasible computa-

Received October 1996; revised December 1998.

AMS 1991 subject classifications. Primary 65U05, 62F10; secondary 62M30, 60K35.

Key words and phrases. Incomplete data, optimization, maximum likelihood, missing data, Monte Carlo algorithm, EM algorithm, simulation, stochastic algorithm.

tions, the EM paradigm is no longer directly applicable. The problems raised by the potentially difficult global maximization involved in the M-step have recently been addressed successfully. A possible solution consists in replacing the global optimization by a chain of simpler conditional maximization, leading to the so-called ECM algorithm [see for example, Meng and Rubin (1993), Liu and Rubin (1994)]. Another approach is to use a single iteration of an approximate Newton's method leading to the EM gradient algorithm [see Lange (1995)]. On the contrary, only partial answers have been obtained to deal with problems for which the expectation of the complete likelihood cannot be done in closed form. A possible solution to cope with this problem has been proposed by Wei and Tanner (1990) (see Section 2). The basic idea is to compute the expectation in the E-step by means of a Monte Carlo method. In this contribution, a stochastic version of the EM algorithm (referred to as SAEM, standing for stochastic approximation EM) is presented, as an alternative to the MCEM algorithm. It makes use of a stochastic approximation procedure for estimating the conditional expectation of the complete data log-likelihood. Given the current approximation of the parameters, complete data are simulated under the a posteriori density; these simulated complete data are then used to update the current value of the conditional expectation of the complete data likelihood. A (decreasing) step size is used to control the allowed amount of update during the successive iterations of the algorithm (Section 3). The convergence of the sequence of parameter estimates to a stationary point of the incomplete likelihood is established for a general class of complete data likelihood functions. It is further demonstrated that *only* the local maxima of the incomplete likelihood are attractive for the stochastic approximation algorithm; that is, convergence toward saddle points are avoided with probability 1.

2. The EM and the MCEM algorithms. In this section, we shall review the key properties of the EM algorithm that we shall need, as derived by Dempster, Laird and Rubin (1977). Let μ be a σ -finite positive Borel measure on \mathbb{R}^l and let $\mathcal{F} = \{f(z; \theta), \theta \in \Theta\}$ be a family of positive integrable Borel functions on \mathbb{R}^l , where Θ is a subset of \mathbb{R}^p . Define

$$\begin{aligned} (1) \quad & g(\theta) \triangleq \int_{\mathbb{R}^l} f(z; \theta) \mu(dz), \\ (2) \quad & l(\theta) \triangleq \log g(\theta), \\ (3) \quad & p(z; \theta) \triangleq \begin{cases} f(z; \theta)/g(\theta), & \text{if } g(\theta) \neq 0, \\ 0, & \text{if } g(\theta) = 0. \end{cases} \end{aligned}$$

Note that $p(z; \theta)$ defines a probability density function w.r.t. to the measure μ . In the terminology introduced by Geyer (1994), \mathcal{F} is a family of *unnormalized density*, $\mathcal{P} = \{p(z; \theta), \theta \in \Theta\}$ is the family of *normalized density* and $g: \Theta \rightarrow [0, \infty)$ is the *normalizing function* for the family \mathcal{F} . We wish to find the value $\hat{\theta} \in \Theta$ that maximizes $g(\theta)$ (conditions upon which such maximum

exists are detailed later on). Many statistical inference problems fall into the framework (1). In the standard missing data problem:

1. $g(\theta)$ is the incomplete data likelihood, that is, the likelihood of the observed data y [the dependence of $g(\theta)$ w.r.t y is here implicit];
2. $f(z; \theta)$ is the complete data likelihood, that is, the likelihood of the complete data x obtained by augmenting the observed data y with the missing data z : $x = (y, z)$,
3. $p(z; \theta)$ is the posterior distribution of the missing data z given the observed data y (often referred to as the *predictive distribution*).

Other examples may be found in spatial statistics and stochastic geometry [Geyer (1994, 1996)] or in Bayesian inference. We will find it convenient to use in the sequel the classical terminology of the *missing data problem*, even though the approaches developed here apply to a more general context. Define

$$(4) \quad Q(\theta | \theta') = \int_{\mathbb{R}^l} \log f(z; \theta) p(z; \theta') \mu(dz).$$

The EM algorithm is useful in situations where maximization of $\theta \rightarrow Q(\theta | \theta')$ is much simpler than direct maximization of $\theta \rightarrow l(\theta)$. Indeed, EM is an iterative algorithm, which maximizes $l(\theta)$ by iteratively maximizing $Q(\theta | \theta')$. Each iteration may be formally decomposed into two steps: an E-step and a M-step. At iteration k , the E-step consists of evaluating

$$Q(\theta | \theta_k) = \int_{\mathbb{R}^l} \log f(z; \theta) p(z; \theta_k) \mu(dz).$$

In the M-step, the value of θ maximizing $Q(\theta | \theta_k)$ is found. This yields the new parameter estimate θ_{k+1} . This two-step procedure is repeated until convergence is apparent. The essence of the EM algorithm is that increasing $Q(\theta | \theta_k)$ forces an increase of $l(\theta)$.

In situations where global maximization of $\theta \rightarrow Q(\theta | \theta')$ is not in simple closed form, alternate solutions can be contemplated. Possible solutions are considered in Meng and Rubin (1993), where global maximization is replaced by a cycle of simpler maximization problem and in Lange (1995), where the M-step is replaced by a single iteration of a Newton's method (see Section 8). In certain situations, the expectation step can be numerically involved or even intractable. To deal with these cases, Wei and Tanner (1990) [see also Tanner (1993)] propose to replace the expectation in the computation of $Q(\theta | \theta_k)$ by a Monte Carlo integration, leading to the so-called Monte Carlo EM (MCEM algorithm). Basically, the E-step at iteration k is replaced by the following procedure.

Simulation step (S-step): generate $m(k)$ realizations $z_k(j)$ ($j = 1, \dots, m(k)$) of the missing data vector under the distribution function $p(z; \theta_k)$,
Monte Carlo integration: compute the current approximation of $Q(\theta | \theta_k)$

according to

$$(5) \quad \tilde{Q}_k(\theta) = m(k)^{-1} \sum_{j=1}^{m(k)} \log f(z_k(j); \theta).$$

The maximization step remains unchanged [see Tanner (1993) for implementation details]. Simulation under the posterior density $p(z; \theta)$ is generally easy when $p(z; \theta)$ is a product of low-dimensional marginal distributions; standard stochastic simulation methods such as importance sampling may in such a case be used. When the posterior $p(z; \theta)$ cannot be decomposed, it is often convenient to resort to Markov chain Monte Carlo simulation methods.

3. SAEM: The stochastic approximation EM algorithm. We propose in this contribution an alternative scheme, which shares most of the attractive behavior of the MCEM algorithm. Similar to the MCEM algorithm, the basic idea is to split the E-step into a simulation step and an integration step. Similar to the MCEM, the S-step consists of generating realizations of the missing data vector under the posterior distribution $p(z; \theta)$; the Monte Carlo integration is substituted by a stochastic averaging procedure. The proposed algorithm may thus be summarized as follows.

Simulation: generate $m(k)$ realizations $z_k(j)$ ($j = 1, \dots, m(k)$) of the missing data under the posterior density $p(z; \theta_k)$.

Stochastic approximation: update $\hat{Q}_k(\theta)$ according to

$$(6) \quad \hat{Q}_k(\theta) = \hat{Q}_{k-1}(\theta) + \gamma_k \left(\frac{1}{m(k)} \sum_{j=1}^{m(k)} \log f(z_k(j); \theta) - \hat{Q}_{k-1}(\theta) \right),$$

where $\{\gamma_k\}_{k \geq 1}$ is a sequence of positive step size.

Maximization: maximize $\hat{Q}_k(\theta)$ in the feasible set Θ , that is, find $\theta_{k+1} \in \Theta$ such that

$$\hat{Q}_k(\theta_{k+1}) \geq \hat{Q}_k(\theta) \quad \forall \theta \in \Theta.$$

The SAEM algorithm shares some similarity with the stochastic EM (also referred to as the probabilistic teacher) algorithm developed by Celeux and Diebolt in a series of paper [see Celeux and Diebolt (1988, 1992), Diebolt and Celeux (1993); see also Diebolt and Ip (1996)]. The main difference w.r.t to the SEM approach is the use of a decreasing sequence of step size to approximate (using stochastic approximation) the EM auxiliary function. In the SEM algorithm, the stepsize is set to zero $\gamma_k = 0$ (no approximation) and the number of simulations by iteration is constant [most often, $m(k) = 1$]; hence, under mild assumptions, $\{\theta_k\}_{k \geq 0}$ is a homogeneous Markov chain. Under appropriate conditions [see, e.g., Diebolt and Ip (1996)], it may be shown that this Markov chain is geometrically ergodic and point estimate can be obtained, for example, by computing ergodic averages. The behavior of the stationary distribution of the chain has been characterized in some very simple scenarios; see

Diebolt and Ip (1996). The SAEM algorithm is also related to the stochastic approximation procedure developed by Younes [see Younes (1989, 1992)] for incompletely observed Gibbsian fields (more on this later).

It should be noted that the stochastic approximation is used here in a slightly unusual context, because the amount of data to process (i.e., the incomplete data) is fixed. **The convergence of the SAEM algorithm** depends on the choice of step sizes γ_k and/or the specification of $m(k)$ used in the stochastic approximation. It is inappropriate to start with small values for step size γ_k and/or large values for the number of simulations $m(k)$. Rather, it is recommended that one decrease γ_k and/or increase $m(k)$ as the current approximation of the parameter vector moves closer to a stationary point. These intuitions are formally developed in the next section, where convergence of the sequence $\{\theta_k\}_{k \geq 0}$ is discussed. When the maximization step is straightforward to implement and/or is, from the computational point of view, much faster than the simulation step, one may set the number of simulations $m(k) = 1$ for all the iterations.

In comparison with the MCEM algorithm, the SAEM makes a more efficient use of the imputed missing values. At each new iteration of the MCEM algorithm, a whole set of missing values needs to be simulated and all the missing values simulated during the previous iterations are dropped. In the SAEM algorithm, *all* the simulated missing values contribute to the evaluation of the auxiliary quantity $\hat{Q}_k(\theta)$; they are gradually discounted, with a forgetting factor inversely proportional to the step size. As a result, the SAEM algorithm moves toward modal areas more quickly than the MCEM in terms of the number of simulations. The computational advantage of the SAEM algorithm over the MCEM algorithm is striking in problems where maximization is much cheaper than simulation.

All the acceleration methods devised from the EM paradigm [such as the method proposed by Louis (1982) or the modified scoring developed by Meilijson (1989)] can be adapted to the SAEM algorithm following essentially the same lines as Wei and Tanner (1990) for the MCEM algorithm. Note in particular that the gradient (the Fisher score function) and the Hessian (observed Fisher information) of $l(\theta)$ can be obtained almost directly by using the values of the simulated missing data $z_k(j)$, $1 \leq j \leq m(k)$. Using the so-called Fisher identity, the Jacobian of the incomplete data log-likelihood $l(\theta)$ is equal to the conditional expectation of the complete data log-likelihood,

$$(7) \quad \partial_\theta l(\theta) \triangleq E_\theta[\partial_\theta \log f(Z; \theta)],$$

where ∂_θ denotes the differential with respect to θ and

$$(8) \quad E_\theta[\psi(Z)] \triangleq \int \psi(z) p(z; \theta) \mu(dz).$$

Equation (7) suggests the following stochastic approximation scheme:

$$(9) \quad \Delta_k = \Delta_{k-1} + \gamma_k \left(\frac{1}{m(k)} \sum_{j=1}^{m(k)} \partial_\theta \log f(z_k(j); \theta_k) - \Delta_{k-1} \right).$$

Using Louis's missing information principle [Louis (1982)], the Hessian of l at θ , $\partial_\theta^2 l(\theta)$ (the observed Fisher information matrix) may be expressed as

$$(10) \quad -\partial_\theta^2 l(\theta) = -E_\theta[\partial_\theta^2 \log f(Z; \theta)] - \text{Cov}_\theta[\partial_\theta \log f(Z; \theta)],$$

where $\text{Cov}_\theta[\psi(Z)] \triangleq E_\theta[(\psi(Z) - E_\theta(\psi(Z)))(\psi(Z) - E_\theta(\psi(Z)))^t]$. Using this expression, it is possible to derive the following stochastic approximation procedure to approximate $\partial_\theta^2 l(\theta)$:

$$(11) \quad G_k = G_{k-1} + \gamma_k \left(\frac{1}{m(k)} \sum_{j=1}^{m(k)} (\partial_\theta^2 \log f(z_k(j); \theta_k) + \partial_\theta \log f(z_k(j); \theta_k) \partial_\theta \log f(z_k(j); \theta_k)^t) \right) - G_{k-1},$$

$$H_k = G_k - \Delta_k \Delta_k^t.$$

Provided the SAEM algorithm converges to a limiting value θ^* (see Section 6) and $\theta \rightarrow l(\theta)$ is sufficiently smooth, H_k converges to $\partial_\theta^2 l(\theta^*)$. When $l(\theta)$ is an incomplete data likelihood function of a *regular* statistical experiment [see, e.g., Ibragimov and Has'minski (1981)], the maximum likelihood estimator [i.e., the value of θ that maximizes $l(\theta)$ over the feasible set Θ] is asymptotically normal and the inverse of the observed Fisher information matrix $-\partial_\theta^2 l(\theta^*)$ converges to the asymptotic covariance of the estimator. Hence, the limiting value of G_k can be used to assess the dispersion of the estimator.

4. Convergence of the EM algorithm for curved exponential families. The convergence of the EM algorithm has been addressed by many different authors, starting with the seminal paper by Dempster, Laird and Rubin (1977). Convergence under very general conditions has been established by Wu (1983). Before embarking on the more subtle task of investigating the convergence property of the SAEM algorithm, we briefly recall the basic ingredients needed to prove the convergence of the EM algorithm. We base our discussion mainly on the recent work by Lange (1995). In the sequel, we restrict attention to models for which the unnormalized density $f(z; \theta)$ belongs to the curved exponential family.

(M1) The parameter space Θ is an open subset of \mathbb{R}^p . The complete data likelihood function is given by

$$(12) \quad f(z; \theta) = \exp \left\{ -\psi(\theta) + \langle \tilde{S}(z), \phi(\theta) \rangle \right\},$$

where $\langle \cdot, \cdot \rangle$ denotes the scalar product, $\tilde{S}(\cdot)$ is a Borel function on \mathbb{R}^l taking its values in an open subset \mathcal{S} of \mathbb{R}^m . Moreover, the convex hull of $\tilde{S}(\mathbb{R}^l)$ is included in \mathcal{S} , and, for all $\theta \in \Theta$,

$$(13) \quad \int_{\mathbb{R}^l} |\tilde{S}(z)| p(z; \theta) \mu(dz) < \infty.$$

The choice of $\tilde{S}(z)$ is of course not unique. It is generally guided by considerations on the practical implementation of the algorithm. This assumption is often met in situations where the EM algorithm is employed. It implies that the expectation step reduces to the computation of

$$E_\theta[\tilde{S}(Z)] \triangleq \int_{\mathbb{R}^t} \tilde{S}(z) p(z; \theta) \mu(dz).$$

Note that $E_\theta[\tilde{S}(Z)]$ belongs to the convex hull of $\tilde{S}(\mathbb{R}^t) \subset \mathcal{S}$. To develop the theory further, some regularity conditions are needed. Define $L: \mathcal{S} \times \Theta \rightarrow \mathbb{R}$ as

$$(14) \quad L(s; \theta) = -\psi(\theta) + \langle s, \phi(\theta) \rangle.$$

(M2) The functions, $\psi(\theta)$ and $\phi(\theta)$ are twice continuously differentiable on Θ .

(M3) The function $\bar{s}: \Theta \rightarrow \mathcal{S}$ defined as

$$(15) \quad \bar{s}(\theta) \triangleq \int_{\mathbb{R}^t} \tilde{S}(z) p(z; \theta) \mu(dz)$$

is continuously differentiable on Θ .

(M4) The function $l(\theta)$ is continuously differentiable on Θ and

$$\partial_\theta \int f(z; \theta) \mu(dz) = \int \partial_\theta f(z; \theta) \mu(dz).$$

(M5) There exists a function $\hat{\theta}: \mathcal{S} \rightarrow \Theta$, such that

$$\forall \theta \in \Theta, \quad \forall s \in \mathcal{S}, \quad L(s; \hat{\theta}(s)) \geq L(s; \theta).$$

Moreover, the function $\hat{\theta}(s)$ is continuously differentiable on \mathcal{S} .

In many models of practical interest, the function $\theta \rightarrow L(s; \theta)$ has a unique global maximum, and the existence and the differentiability of $s \rightarrow \hat{\theta}(s)$ is a direct consequence of the implicit function theorem. Using the notations introduced above, the EM-reestimation functional $Q(\theta | \theta')$ may be expressed as

$$(16) \quad Q(\theta | \theta') = L(\bar{s}(\theta'); \theta).$$

Using these assumptions, the iteration $\theta_k \rightarrow \theta_{k+1}$ of the EM algorithm is defined as

$$(17) \quad \theta_{k+1} = T(\theta_k) = \hat{\theta}(\bar{s}(\theta_k)).$$

Adaptations of the algorithm to situations where the maximum of $L(s; \theta)$ cannot be computed explicitly will be considered in Section 8. The following properties are simple consequences of the basic properties of the EM algorithm:

$\theta \rightarrow T(\theta) = \hat{\theta}(\bar{s}(\theta))$ is continuously differentiable on the feasible set Θ . Moreover, $l(T(\theta)) \geq l(\theta)$, with equality occurring only when θ is a fixed point of T ;

The set of fixed points of $T(\theta)$ coincides with the set \mathcal{L} of stationary points of $l(\theta)$:

$$(18) \quad \mathcal{L} = \{\theta \in \Theta; \partial_{\theta} l(\theta) = 0\} = \{\theta \in \Theta; \theta = T(\theta)\}.$$

Note that \mathcal{L} is closed. In the terminology of Dempster, Laird and Rubin (1977), $T(\theta)$ is a continuous EM algorithm. The EM recursion may equivalently be written in that case in terms of expectation of the complete data sufficient statistics $s_k \triangleq \bar{s}(\theta_{k-1})$ as

$$(19) \quad s_{k+1} = \bar{s}(\hat{\theta}(s_k)) \triangleq G(s_k),$$

this form coming up more naturally in the analysis of the SAEM algorithm. Under (M1–M5), G also is continuously differentiable on \mathcal{L} . For the model above, the convergence of the EM algorithm can be established upon the additional condition that the sequence $\theta_{k+1} = T(\theta_k)$ stays within some compact set in Θ . The following theorem is adapted from Lange (1995) [Propositions 4 and 5; see also Theorems 1, 2 and 4 in Wu (1983)]. For $\theta \in \Theta$, we denote $\mathcal{L}(\theta)$ the set of limit points of the sequence $\theta_{k+1} = T(\theta_k)$, starting at $\theta_0 = \theta$. We denote $\text{clos}(A)$ and $\text{int}(A)$ the closure and the interior of the set A , and, when A is a closed set, $d(x, A)$ the distance of x to A .

THEOREM 1. *Assume that (M1)–(M5) hold. Assume in addition that for any $\theta \in \Theta$, $\text{clos}(\mathcal{L}(\theta))$ is a compact subset of Θ . Then, for any initial point $\theta_0 = \theta$, the sequence $l(\theta_k)$ is increasing and $\lim_{k \rightarrow \infty} d(\theta_k, \mathcal{L}(\theta)) = 0$.*

Note that the compactness of the sequence $\{\theta_k\}_{k \geq 0}$ plays an essential role in the proof. In Lange (1995) [see also Wu (1983)] this property is guaranteed by assuming that $l(\theta)$ is continuous and lower compact in the sense that the level set

$$(20) \quad \{\theta \in \Theta: -l(\theta) \leq c\}$$

is compact for any $c \geq 0$.

5. General results on Robbins–Monro type stochastic approximation procedures. The SAEM algorithm is a Robbins–Monro type stochastic approximation procedure. These procedures have received considerable attention in the literature [see, e.g., Kushner and Clark (1978), Duflo (1996), Kushner and Yin (1997) and the references therein]. The most commonly used tool for proving w.p.1 convergence for such algorithm is of far Theorem 2.3.1 of Kushner and Clark (1978) [see Fort and Pagès (1996)]. However, some conditions of this theorem prove intractable, because the mean field (see definition below) associated with the SAEM algorithm is multimodal in many situations of interest. In the sequel, we use an alternate technique to prove the convergence, which extends results obtained earlier in Delyon (1996). The type of assumptions on which these results are based are, of course, well suited for the analysis of the SAEM algorithm. However, these results are applicable

in a much wider context; this is why we present these results for a general Robbins–Monro (RM) stochastic approximation procedure. We will apply these results to the SAEM algorithm in the next section. Throughout this section, we consider the following RM stochastic approximation procedure:

$$(21) \quad s_n = s_{n-1} + \gamma_n h(s_{n-1}) + \gamma_n e_n + \gamma_n r_n,$$

where $\{e_n\}_{n \geq 1}$ and $\{r_n\}_{n \geq 1}$ are random processes defined on the same probability space taking their values in an open subset $\mathcal{X} \subset \mathbb{R}^m$; h is referred to as the mean field of the algorithm; $\{r_n\}_{n \geq 1}$ is a remainder term and $\{e_n\}_{n \geq 1}$ is the stochastic excitation.

THEOREM 2. *Assume that:*

(SA0) *w.p.1, for all $n \geq 0$, $s_n \in \mathcal{X}$;*

(SA1) *$\{\gamma_n\}_{n \geq 1}$ is a decreasing sequence of positive number such that $\sum_{n=1}^{\infty} \gamma_n = \infty$;*

(SA2) *The vector field h is continuous on \mathcal{X} and there exists a continuously differentiable function $V: \mathcal{X} \rightarrow \mathbb{R}$ such that*

(i) *for all $s \in \mathcal{X}$, $F(s) = \langle \partial_s V(s), h(s) \rangle \leq 0$;*

(ii) *$\text{int}(V(\mathcal{L})) = \emptyset$, where $\mathcal{L} \triangleq \{s \in \mathcal{X}: F(s) = 0\}$;*

(SA3) *w.p.1, $\text{clos}(\{s_n\}_{n \geq 0})$ is a compact subset of \mathcal{X} ;*

(SA4) *w.p.1, $\lim_{p \rightarrow \infty} \sum_{n=1}^p \gamma_n e_n$ exists and is finite, $\lim_{n \rightarrow \infty} r_n = 0$.*

Then, w.p.1, $\limsup d(s_n, \mathcal{L}) = 0$.

The proof of this theorem is given in Appendix A. Note that this theorem shares most of the assumptions of Theorem 2.3.1 of Kushner and Clark (1978) and Theorem 2.1 in Kushner and Yin (1997). The main difference lies in assumption (SA2), which replaces the following recurrence assumption: “the sequence $\{s_k\}_{k \geq 0}$ is w.p. 1 (or with a probability greater than or equal to ρ) infinitely often in some compact set in the domain of attraction of an asymptotically stable point s^* (in the sense of Lyapunov) of h .” [Note that a stationary point of the vector field h is a point s^* such that $h(s^*) = 0$. The *domain of attraction* of s^* is the set $D(s^*) \subset \mathcal{X}$, such that for every $x \in D(s^*)$ the ODE $ds(t)/dt = h(s(t))$, with initial condition $s(0) = x$ has a solution for $t \geq 0$ verifying $\lim_{t \rightarrow \infty} s(t) = s^*$.] This assumption is almost impossible to verify in situations where there are more than one stationary point (which is by far the most common situation when dealing with likelihoods). The boundedness of the sequence $\{s_k\}_{k \geq 0}$ does not imply in that case that $\{s_k\}_{k \geq 0}$ is infinitely often in a compact subset of the domain of attraction of a stationary point. Assumption (SA3) implies that, along every trajectory of (21), the sequence $\{s_k\}_{k \geq 0}$ stays in a compact subset (depending upon the trajectory) of \mathcal{X} . This assumption (compactness) is often difficult to check, because the behavior of the logarithm of the normalizing function $l(\theta)$ on the boundary of the feasible set is most often pathological. This assumption can be replaced by a recurrence condition, provided there exists a Lyapunov function controlling the excursion outside the compact sets of \mathcal{X} .

LEMMA 1. Assume (SA0)–(SA2). Assume in addition:

(STAB1) There exists a continuously differentiable function $W: \mathcal{X} \rightarrow \mathbb{R}$ and a compact set $\mathcal{K} \subset \mathcal{X}$ such that

- (i) For all $c \geq 0$, the level set $\{s \in \mathcal{X}: W(s) \leq c\}$ is a compact subset of \mathcal{X} , and for $c < c'$, $\{s \in \mathcal{X}: W(s) \leq c\} \subset \text{int}(\{s \in \mathcal{X}: W(s) \leq c'\})$;
- (ii) $\langle \partial_s W(s), h(s) \rangle < 0$, for all $s \in \mathcal{X} \setminus \mathcal{K}$;

(STAB2) For any positive integer M , the sequence $\{\varepsilon_n\}_{n \geq 1}$ defined as

$$\varepsilon_n = \sum_{k=1}^n \gamma_k e_k \mathbb{1}(W(s_{k-1}) \leq M)$$

converges w.p.1 to a finite limit. In addition, w.p.1,

$$\limsup_{k \rightarrow \infty} |r_k| \mathbb{1}(W(s_{k-1}) \leq M) = 0;$$

(STAB3) w.p.1, the sequence $\{s_k\}_{k \geq 0}$ admits a limit point in \mathcal{X} .

Then $\limsup W(s_k) < \infty$ w.p.1 and assumptions (SA3) and (SA4) are satisfied.

The proof of this result may be adapted from the proof of Theorem 1 in Delyon (1996) (it is omitted for brevity). The function W does not necessarily coincide with the Lyapunov function V used in Theorem 2 [assumption (SA2)]. The condition (STAB3) is a recurrence condition: w.p.1, the sequence $\{s_k\}_{k \geq 0}$ returns infinitely often in a compact set. In some cases, the recurrence condition (STAB3) is not verified. A possible solution to overcome this problem consists in implementing an explicit stabilization device. To that purpose, we will present a simple modification of (21), adapting the procedure presented in Chen, Guo and Gao (1988) [see also Andradottir (1995)]. The procedure consists in truncating the original RM recursion (21): every time s_k is outside a specific set, it is re-initialized at a point chosen at random in a compact set of \mathcal{X} . In the technique proposed by Chen, Guo and Gao (1988), the truncation bounds are random functions of the recursion index k . The advantage of this approach is that the truncation *does not* modify the mean field (and hence, the set of stationary points of the recursion) of the original RM recursion; this is in contrast with the traditional projection or truncation approach, where the truncation set is fixed. This is also advantageous from the practical point of view, because the truncation set is selected automatically. The procedure goes as follows. Choose a sequence of compact subsets such that

$$(22) \quad \mathcal{K}_n \subset \text{int}(\mathcal{K}_{n+1}), \quad \mathcal{X} = \bigcup_{n=0}^{\infty} \mathcal{K}_n.$$

Every time s_k wanders out of the compact subset \mathcal{K}_{n_k} , the sequence is reset at an arbitrary point inside \mathcal{K}_0 , and the index n_k is increased (n_k is thus the

number of projections up to the k th iteration). More precisely,

$$\begin{aligned}
 \hat{s}_k &= s_{k-1} + \gamma_k h(s_{k-1}) + \gamma_k e_k + \gamma_k r_k, \\
 (23) \quad & \text{if } \hat{s}_k \in \mathcal{K}_{n_{k-1}}: \begin{cases} s_k = \hat{s}_k, \\ n_k = n_{k-1}, \end{cases} \\
 & \text{if } \hat{s}_k \notin \mathcal{K}_{n_{k-1}}: \begin{cases} s_k = s'_k, \\ n_k = n_{k-1} + 1, \end{cases}
 \end{aligned}$$

where $\{s'_k\}_{k \geq 0}$ is an arbitrary sequence of random variables taking their values in \mathcal{K}_0 .

The modified recursion (23) automatically satisfies the recurrence condition (STAB3), because by construction, all the sequences that do not stay in a compact set of \mathcal{X} are infinitely often in the set \mathcal{K}_0 and hence have a limit point in that set, since that set is compact. We have the following convergence result.

THEOREM 3. *Consider $\{s_k\}_{k \geq 0}$ as given by the truncated RM procedure (23) and assume that (SA0)–(SA2) and (STAB1) and (STAB2) hold. Then, w.p.1, $\lim d(s_k, \mathcal{L}) = 0$.*

5.1. Rate of convergence, step size selection and averaging of iterates. The difficulty of selecting an appropriate step size sequence has long been considered as a serious handicap for practical applications. In a path-breaking paper, Polyak (1990) [see also Polyak and Juditski (1992)] showed that if the step size γ_n goes to zero slower than $1/n$ (yet fast enough to ensure convergence at a given rate), then the averaged sequence, defined as $n^{-1} \sum_{i=1}^n s_i$ converges to its limit at an *optimum* rate (see the comments below). This result implies that we should use *large* step size (larger than $\gamma_n = 1/n$; typically $\gamma_n = n^{-2/3}$) and an off-line averaging controls the increased noise effect. The practical value of the averaging method has been reported for many different stochastic approximation procedures [see Kushner and Yin (1997), Chapter 11, for a thorough investigation of the averaging method; see also Delyon and Juditski (1992)]. In particular, it happens that this approach tends to robustify the overall procedure, in the sense that the *primary* approximation algorithm (which uses large step size) is less likely to get stuck at an early stage in a local minimum and more likely to have a faster initial convergence.

The use of an averaging procedure makes sense when the stochastic approximation procedure converges to a *regular stable stationary point*, defined as follows.

DEFINITION. We say that s^* is a *regular stable stationary point* of the stochastic approximation procedure (21), if (i) $h(s^*) = 0$, (ii) h is twice differentiable in a neighborhood of s^* and (iii) $H(s^*)$, the Jacobian matrix of h at s^* , is a Hurwitz matrix; that is, the real parts of the eigenvalues of $H(s^*)$ are negative.

As evidenced by Kushner and Yin (1997), averaging improves the asymptotic convergence whenever s_n converges to s^* with a given rate of convergence. We say that the rate of convergence of s_n about s^* is $\gamma_n^{1/2}$ if the sequence $\gamma_n^{-1/2}(s_n - s^*)$ is bounded in probability, that is,

$$(24) \quad \lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} P\left(\gamma_n^{-1/2} \|s_n - s^*\| \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s^*\| = 0\right) \geq M\right) = 0.$$

To that purpose, we need to strengthen the assumptions made on the step size sequence $\{\gamma_n\}_{n \geq 1}$ and on the disturbance processes $\{e_n\}_{n \geq 1}, \{r_n\}_{n \geq 1}$. Denote $\mathcal{F} = \{\mathcal{F}_n, n \geq 1\}$, the increasing family of σ -algebra generated by $s_0, e_1, \dots, e_n, r_1, \dots, r_n$. For $\rho > 0$, we define $e_n(\rho) \triangleq e_n \mathbb{1}(\|s_{n-1} - s^*\| \leq \rho)$. The process $\{e_n(\rho)\}$ is adapted to \mathcal{F} .

(AVE1) For some $\rho > 0$ such that $\{s \in \mathcal{X} : \|s - s^*\| \leq \rho\} \subset \mathcal{X}$, there exists a constant $C(\rho) < \infty$, such that, for any deterministic sequence of matrices $\{\alpha_k\}_{k \geq 1}$, and any $n \geq 1$, we have

$$(25) \quad E\left(\left\|\sum_{k=1}^n \alpha_k e_k(\rho)\right\|^2\right) \leq C(\rho) \sum_{k=1}^n \|\alpha_k\|^2.$$

In addition, w.p.1, $\lim_{p \rightarrow \infty} \sum_{n=1}^p \gamma_n e_n(\rho)$ exists and is finite.

(AVE2) w.p.1, $\limsup_{n \rightarrow \infty} \gamma_n^{-1/2} \|r_n\| \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s^*\| = 0) < \infty$;

(AVE3(α)) ($1/2 < \alpha < 1$) $\{\gamma_n\}_{n \geq 1}$ is a nonincreasing sequence of positive real numbers; in addition, $\lim_{n \rightarrow \infty} n^\alpha \gamma_n = \gamma_*$ and $\gamma_{n+1}/\gamma_n = 1 + O(n^{-1})$.

If $\{e_k\}_{k \geq 1}$ is a L^2 -martingale increment (w.r.t \mathcal{F}), then $\{e_k(\rho)\}_{k \geq 0}$ also is a L^2 -martingale increment and (AVE1) is verified with $C(\rho) = \sup_{k \geq 0} E(\|e_k(\rho)\|^2)$, provided that $\sup_{k \geq 0} E(\|e_k(\rho)\|^2) < \infty$. Assumption (AVE1) is also verified for a large class of weak dependent disturbance process $\{e_k\}$ (including Markov chains), by application of Rosenthal’s inequality [see, e.g., Doukhan (1994)].

Under these assumptions, it may be shown that $\gamma_n^{-1/2}(s_n - s^*) \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s^*\| = 0)$ is bounded in probability, provided s^* is a regular stable stationary point (see Lemma 6). Next, we define the averaged version of the algorithm,

$$(26) \quad \begin{aligned} s_n &= s_{n-1} + \gamma_n h(s_{n-1}) + \gamma_n e_n + \gamma_n r_n, \\ \bar{s}_n &= \bar{s}_{n-1} + n^{-1}(s_{n-1} - \bar{s}_{n-1}), \end{aligned}$$

where $\bar{s}_0 = 0$. The basic result for the convergence of the average of the iterates is given below. We denote \rightarrow_P and $\rightarrow_{\mathcal{L}}$ the convergences in probability and in distribution, and $\mathcal{N}(\mu, \Gamma)$ the multivariate normal distribution with mean μ and covariance matrix Γ .

THEOREM 4. Assume that (AVE1)–(AVE3(α)) hold, with $1/2 < \alpha < 1$. Assume in addition that s^* is a stable regular stationary point. Then, $\gamma_n^{-1/2}(s_n - s^*)$

$\mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_\star\| = 0)$ is bounded in probability and

$$(27) \quad \sqrt{n} \left(\bar{s}_n - s^\star + n^{-1} \sum_{k=1}^n \nu_k \right) \mathbb{1} \left(\lim_{n \rightarrow \infty} \|s_n - s_\star\| = 0 \right) \rightarrow_P 0,$$

where $\nu_k \triangleq e_k + r_k$. If, in addition,

$$(28) \quad 1/\sqrt{n} \sum_{k=1}^n \nu_k \mathbb{1} \left(\lim_{n \rightarrow \infty} \|s_n - s_\star\| = 0 \right) \rightarrow_{\mathcal{L}} \mathcal{N}(0, S) \mathbb{1} \left(\lim_{n \rightarrow \infty} \|s_n - s_\star\| = 0 \right),$$

then

$$(29) \quad \begin{aligned} &\sqrt{n}(\bar{s}_n - s^\star) \mathbb{1} \left(\lim_{n \rightarrow \infty} \|s_n - s_\star\| = 0 \right) \\ &\rightarrow_{\mathcal{L}} \mathcal{N}(0, H(s^\star)^{-1} S H(s^\star)^{-1}) \mathbb{1} \left(\lim_{n \rightarrow \infty} \|s_n - s_\star\| = 0 \right). \end{aligned}$$

REMARK 1. As shown in Chickin and Poznyak (1984) and Chickin (1988), the rate obtained in (29) is the best asymptotic rate of convergence. This rate may be achieved using a Gauss–Newton type stochastic approximation algorithm; however, in this latter case, knowledge of $H(s^\star)$ is needed; this is not required when the averaging of iterates is used.

REMARK 2. The central limit theorem is checked under a wide variety of conditions. We apply the result in the case where $\{e_n\}_{n \geq 0}$ is a martingale increment; the CLT is thus verified under standard Lindeberg’s type conditions [see Hall and Heyde (1980)].

REMARK 3. When coupled with the stabilization device ((23)), the averaging procedure should be modified. The modification consists in resetting to zero the averaged value at every time instants the primary stochastic approximation algorithm is restarted and computing averages from that point. As shown in Theorem 3, the number of times the algorithm is restarted is finite w.p.1; (27) still holds for the modified algorithm.

The proof of Theorem 4 is given in Appendix B. It should be stressed that our result is obtained under much less stringent conditions than those used in Delyon and Juditski (1992) (this contribution deals uniquely with gradient type algorithms; moreover, assumption (A3) of this contribution involves conditions on the growth of the Lyapunov function and of the mean field outside the compact sets of \mathcal{X}). It also weakens some of the assumptions of Theorem 3.1, page 338, in Kushner and Yin (1997); in particular, it avoids the tightness condition (A3.1), which is most often difficult to check (see Remark 4 in Appendix B).

6. Convergence of the SAEM algorithm. To avoid cumbersome notations, we set $m(k) = 1$. The k th iteration of the SAEM algorithm boils down to

$$(30) \quad S_k = S_{k-1} + \gamma_k(\tilde{S}(Z_k) - S_{k-1}), \quad \theta_k = \hat{\theta}(S_k),$$

where $\{Z_k\}_{k \geq 1}$ is the missing value simulated at step k , under the posterior density $p(z; \theta_{k-1})$ [see (4)]. It is assumed that the random variables (r.v.) S_0, Z_1, Z_2, \dots are defined on the same probability space (Ω, \mathcal{A}, P) ; we denote $\mathcal{F} = \{\mathcal{F}_n\}_{n \geq 0}$ the increasing family of σ -algebra generated by the r.v. $S_0, Z_1, Z_2, \dots, Z_n$. In addition, we assume that:

- (SAEM1) For all $k > 0, 0 \leq \gamma_k \leq 1, \sum_{k=1}^\infty \gamma_k = \infty$ and $\sum_{k=1}^\infty \gamma_k^2 < \infty$;
- (SAEM2) $l: \Theta \rightarrow \mathbb{R}$ and $\hat{\theta}: \mathcal{S} \rightarrow \Theta$ are m times differentiable;
- (SAEM3) For all positive Borel function ϕ ,

$$E[\phi(Z_{k+1}) | \mathcal{F}_k] = \int \phi(z)p(z; \theta_k)\mu(dz);$$

(SAEM4) For all $\theta \in \Theta, \int \|\tilde{S}(z)\|^2 p(z; \theta)\mu(dz) < \infty$, and the function

$$(31) \quad \Gamma(\theta) \triangleq \text{Cov}_\theta(\tilde{S}(Z))$$

is continuous w.r.t θ .

Since the convex hull of $\tilde{S}(\mathbb{R}^l) \subset \mathcal{S}$, (SAEM1) implies that, for all $k \geq 0, S_k \in \mathcal{S}$. Assumption (SAEM3) is equivalent to saying that given $\theta_0, \dots, \theta_k$, the simulated missing observations Z_1, \dots, Z_k are conditionally independent. This assumption can be relaxed to allow for Markovian dependence, a situation which is typical when the Markov chain Monte Carlo method is used for simulation. We may express the recursion (30) into the Robbins–Monro form [Kushner and Clark (1978)],

$$(32) \quad S_k = S_{k-1} + \gamma_k h(S_{k-1}) + \gamma_k E_k,$$

where $h(s)$ stands for the *mean field* of the algorithm and E_k is a random perturbation

$$(33) \quad h(s) = E_{\hat{\theta}(s)}(\tilde{S}(Z)) - s = \bar{s}(\hat{\theta}(s)) - s,$$

$$(34) \quad E_k = \tilde{S}(Z_k) - E[\tilde{S}(Z_k) | \mathcal{F}_{k-1}] = \tilde{S}(Z_k) - \bar{s}(\hat{\theta}(S_{k-1})).$$

Proving the convergence amounts to verifying that the assumptions of Theorem 2 are satisfied.

LEMMA 2. *Assume that (M1)–(M5) and (SAEM2) hold. Then (SA2) is satisfied with $V(s) = -l(\hat{\theta}(s))$. Moreover,*

$$(35) \quad \{s \in \mathcal{S}: F(s) = 0\} = \{s \in \mathcal{S}: \partial_s V(s) = 0\},$$

$$(36) \quad \hat{\theta}(\{s \in \mathcal{S}: F(s) = 0\}) = \{\theta \in \Theta: \partial_\theta l(\theta) = 0\} = \mathcal{L},$$

where $F(s) \triangleq \langle \partial_s V(s), h(s) \rangle$.

PROOF. Under (M3), $\bar{s}(\theta)$, is continuously differentiable on Θ . Under (SAEM2), $\hat{\theta}(s)$ is (m times) continuously differentiable on \mathcal{S} and $h(s) = \bar{s}(\hat{\theta}(s)) - s$ is continuously differentiable on \mathcal{S} . Hence, $h(s)$ is bounded on every compact subset of \mathcal{S} . Under assumption (M5), for all s in \mathcal{S} , the function $\hat{\theta}(s)$ satisfies $\partial_\theta L(s; \hat{\theta}(s)) = 0$, thus

$$(37) \quad -\partial_\theta \psi(\hat{\theta}(s)) + s^t \partial_\theta \phi(\hat{\theta}(s)) = 0,$$

where the superscript t denotes the transposition. Under (M2)–(M5), we may differentiate the relation $\partial_\theta L(s; \hat{\theta}(s)) = 0$ with respect to s ,

$$(38) \quad \partial_\theta^2 L(s; \hat{\theta}(s)) \partial_s \hat{\theta}(s) = -\phi(\hat{\theta}(s))^t,$$

where ∂_θ^2 denotes the second-order derivative w.r.t. θ . On the other hand, for all $\theta \in \Theta$, assumption (M4) implies

$$\partial_\theta l(\theta) = -\partial_\theta \psi(\theta) + \bar{s}(\theta)^t \partial_\theta \phi(\theta).$$

Plugging (37) into the previous expression yields

$$\begin{aligned} \partial_\theta l(\hat{\theta}(s)) &= (-s + \bar{s}(\hat{\theta}(s)))^t \partial_\theta \phi(\hat{\theta}(s)) \\ &= h(s)^t \partial_\theta \phi(\hat{\theta}(s)) \\ &= -h(s)^t \partial_s \hat{\theta}(s)^t \partial_\theta^2 L(s; \hat{\theta}(s)), \\ \partial_s l(\hat{\theta}(s)) &= -h(s)^t \partial_s \hat{\theta}(s)^t \partial_\theta^2 L(s; \hat{\theta}(s)) \partial_s \hat{\theta}(s). \end{aligned}$$

Since, under (M5), $\partial_\theta^2 L(s; \hat{\theta}(s)) \leq 0$, we have

$$(39) \quad \begin{aligned} F(s) &\triangleq \langle \partial_s V(s), h(s) \rangle = -\langle \partial_s l(\hat{\theta}(s)), h(s) \rangle \\ &= h(s)^t \partial_s \hat{\theta}(s)^t \partial_\theta^2 L(s; \hat{\theta}(s)) \partial_s \hat{\theta}(s) h(s) \leq 0, \end{aligned}$$

which proves (SA2i). Note that obviously $\{s: \partial_s V(s) = 0\} \subset \{s: F(s) = 0\}$. On the other hand, let $s^* \in \mathcal{S}$ be such that $F(s^*) = 0$, since under (M5) $\partial_\theta^2 L(s^*; \hat{\theta}(s^*))$ is a nonpositive matrix, (39) implies that $\partial_\theta l(\hat{\theta}(s^*)) = 0$, and thus $\partial_s V(s^*) = 0$, showing the reverse inclusion $\{s: F(s) = 0\} \subset \{s: \partial_s V(s) = 0\}$. By (M4) and (M5), function V is m times continuously differentiable. Sard's theorem [Brocker (1975)] implies that $V(\{s: \partial_s V(s) = 0\}) = 0$ has zero Lebesgue measure, showing (SA2ii). \square

We may now formulate a theorem which is the stochastic counterpart of basic convergence Theorem 1 for a continuous EM algorithm.

THEOREM 5. *Assume that the assumptions (M1)–(M5) and (SAEM1)–(SAEM4) hold. Assume in addition, that w.p.1, (A) $\text{clos}(\{S_k\}_{k \geq 1})$ is a compact subset of \mathcal{S} . Then, w.p.1, $\lim_{k \rightarrow \infty} d(S_k, \{s \in \mathcal{S}: \partial_s V(s) = 0\}) = 0$ and $\lim_{k \rightarrow \infty} d(\theta_k, \mathcal{S}) = 0$.*

PROOF. (SA0) is verified under (M1) and (SAEM1) because $0 < \gamma_k < 1$ and the convex hull of $\tilde{S}(\mathbb{R}^m)$ is included in \mathcal{S} . (SA1) is implied by (SAEM1) and (SA3) by (A). Note that under (A), there exists w.p.1 a compact set \mathcal{X} , such that $S_k \in \mathcal{X}$ for all $k \geq 0$. Denote $M_n = \sum_{k=1}^n \gamma_k E_k$. Then $\{M_n\}_{n \geq 1}$ is a \mathcal{F} -martingale which satisfies, under (SAEM1)–(SAEM3),

$$\begin{aligned} & \sum_{n=1}^{\infty} E[\|M_{n+1} - M_n\|^2 \mid \mathcal{F}_n] \\ & \leq \sum_{n=1}^{\infty} \gamma_{n+1}^2 \int \|\tilde{S}(z)\|^2 p(z; \hat{\theta}(S_n)) \mu(dz) < \infty \quad \text{w.p.1} \end{aligned}$$

since, by (A(i)–A(ii)) and (M5), w.p.1 $\hat{\theta}(S_n)$ is in a compact set $\hat{\theta}(\mathcal{X})$ of Θ . This property implies that w.p.1, $\lim_{n \rightarrow \infty} M_n$ exists [see Hall and Heyde (1980), Theorem 2.15, page 33], proving (SA4). Finally, Lemma 2 shows that (SA2) is satisfied and Theorem 2 holds. Since the function $\hat{\theta}: \mathcal{S} \rightarrow \Theta$ is continuous, the proof is concluded by applying Lemma 2, equation (35). \square

Note that, since the function $\hat{\theta}: \mathcal{S} \rightarrow \Theta$ is continuous, Theorem 5 and Lemma 2, equation (35), imply that

$$(40) \quad d(\theta_k, \mathcal{L}) = d(\hat{\theta}(S_k), \hat{\theta}(\{s \in \mathcal{S} : F(s) = 0\})) \rightarrow 0 \quad \text{w.p.1.}$$

In certain models, the compactness condition (A) of Theorem 5 is trivially satisfied. This is the case when the complete data are bounded and the unnormalized density is well behaved. In others situations, checking (A) may prove intractable; we apply in these cases the stabilization procedure presented in Section 5. The truncated SAEM algorithm takes the form:

$$(41) \quad \begin{aligned} \hat{S}_k &= S_{k-1} + \gamma_k (\tilde{S}(Z_k) - S_{k-1}), \\ \text{if } \hat{S}_k \in \mathcal{X}_{n_{k-1}} &: \begin{cases} S_k = \hat{S}_k, & \theta_k = \hat{\theta}(S_k), \\ n_k = n_{k-1}, \end{cases} \\ \text{if } \hat{S}_k \notin \mathcal{X}_{n_{k-1}} &: \begin{cases} S_k = S'_k, & \theta_k = \hat{\theta}(S_k), \\ n_k = n_{k-1} + 1, \end{cases} \end{aligned}$$

where $\{\mathcal{X}_n\}_{n \geq 0}$ is a sequence of compact sets, such that

$$(42) \quad \mathcal{X}_n \subset \text{int}(\mathcal{X}_{n+1}) \quad \text{and} \quad \bigcup_{n=0}^{\infty} \mathcal{X}_n = \mathcal{S}$$

and S'_k is an arbitrary sequence of random variables taking their values in \mathcal{X}_0 . The convergence of the truncated SAEM algorithm follows from Theorems 3 and 6.

7. Convergence to local maxima. The results obtained in the previous section demonstrate that, under appropriate conditions, the sequence $\{\theta_k\}_{k \geq 1}$ [defined in (30) or (41)] converges to a connected component of the set \mathcal{L}

of stationary points of $l(\theta)$. Assume that the connected components \mathcal{L} are reduced to points $\mathcal{L} = \bigcup\{\theta^*\}$, so that $\theta_k \rightarrow \theta^*$. Depending upon the values of the Hessian of l , these stationary points may correspond either to local maxima, local minima or saddle points. Of course, it would be of interest to find conditions upon which the convergence toward local maxima is guaranteed. Such conditions are worked out in this section.

(LOC1) The stationary points of $l(\theta)$ are isolated: any compact subset of \mathcal{L} contains only a finite number of such points.

(LOC2) For every stationary point $\theta^* \in \mathcal{L}$, the matrices

$$E_{\theta^*}[\partial_\theta L(\tilde{S}(Z); \theta^*)^t \partial_\theta L(\tilde{S}(Z); \theta^*)] \quad \text{and} \quad \partial_\theta^2 L(E_{\theta^*}[\tilde{S}(Z)]; \theta^*)$$

are positive definite.

Assumption (LOC1) is of course satisfied if the Hessian of $l(\theta)$ at θ^* is nonsingular. Under the conditions of Theorems 5 or 3, the sequence $\{\theta_k\}_{k \geq 1}$ converges w.p.1 to a compact and connected subset of \mathcal{L} . Under (LOC1), the connected components of \mathcal{L} are reduced to points so that the sequence $\{\theta_k\}_{k \geq 1}$ converges (w.p.1) pointwise to some point $\theta^* \in \mathcal{L}$. We may associate to θ^* the point in \mathcal{L} ,

$$s^* = \bar{s}(\theta^*) = E_{\theta^*}[\tilde{S}(Z)].$$

Heuristically, since $\theta_k \rightarrow \theta^*$, the missing data Z_k (for large enough k) are all simulated under approximately the same posterior distribution $p(z; \theta^*)$, it is expected that $\{S_k\}_{k \geq 1}$ converges w.p.1 to s^* .

THEOREM 6. *Under the assumptions of Theorem 5 and (LOC1), the sequence $\{S_k\}_{k \geq 1}$ defined by (30) converges to $s^* = \bar{s}(\theta^*)$, where $\theta^* = \lim_{k \rightarrow \infty} \theta_k$. Then θ^* and s^* are fixed points of the EM-mappings $T: \Theta \rightarrow \Theta$ (17) and $G: \mathcal{L} \rightarrow \mathcal{L}$ (19).*

PROOF. Denote $\delta_k = S_k - s^*$. We have

$$(43) \quad \delta_{k+1} = (1 - \gamma_{k+1})\delta_k + \gamma_{k+1}(\tilde{S}(Z_{k+1}) - s^*),$$

$$(44) \quad = \delta_k - \gamma_{k+1}\delta_k + \gamma_{k+1}(\bar{s}(\theta_k) - s^*) + \gamma_{k+1}E_{k+1},$$

where $E_{k+1} = \tilde{S}(Z_{k+1}) - \bar{s}(\theta_k)$. Since $\{\theta_k\}_{k \geq 1}$ converges w.p.1, the sequence $\{\theta_k\}_{k \geq 1}$ is w.p.1 bounded and w.p.1, the martingale $\sum_{k=0}^n \gamma_k E_k$ converges. On the other hand, the continuity of $\bar{s}(\cdot)$ implies that

$$\lim_{k \rightarrow \infty} \bar{s}(\theta_k) - s^* = 0.$$

The sequence $\{\delta_k\}$ converges to zero w.p.1 by application of Theorem 2 with the Lyapunov function $V(x) = -\|x\|^2$. This proves the convergence of S_k to s^* . The relation $\theta_k = \hat{\theta}(S_k)$ implies that $\theta^* = \hat{\theta}(s^*)$ and since $\bar{s}(\theta^*) = s^*$, θ^* and s^* are fixed points of T and G . \square

A similar result can be obtained for the stabilized version of the SAEM algorithm. The regularity of the stationary points (see the definition above) is connected with the stability of the EM mappings $T: \Theta \rightarrow \Theta$ and $G: \mathcal{L} \rightarrow \mathcal{L}$ at $\theta^* \in \mathcal{L}$ and $s^* = \bar{s}(\theta^*)$.

LEMMA 3. Assume that (M1)–(M5) and (LOC1) and (LOC2) hold. Let θ^* be any fixed point of T and let $s^* = \bar{s}(\theta^*)$. Then:

- (i) s^* is a fixed point of the mapping G ;
- (ii) $\partial_\theta T(\theta^*)$ is diagonalizable and its eigenvalues are positive real numbers;
- (iii) the fixed point θ^* is stable if it is a proper maximizer of the normalizing function $l(\theta)$. It is hyperbolic if it is a saddle point of $l(\theta)$. It is unstable if it is a proper local maximizer of $l(\theta)$;
- (iv) s^* is a regular stable stationary point if and only if θ^* is a proper maximizer of $\theta \rightarrow l(\theta)$.

PROOF. (i) Follows from

$$G(s^*) = \bar{s}(\hat{\theta}(\bar{s}(\theta^*))) = \bar{s}(T(\theta^*)) = s^*.$$

(ii) The differential of the mapping T at θ^* is given by [see Lange (1995), equation 7],

$$(45) \quad \partial_\theta T(\theta^*) = \partial_\theta^2 L(s^*; \theta^*)^{-1} \partial_\theta^2 L(s^*; \theta^*) - \partial_\theta^2 l(\theta^*),$$

$$(46) \quad \partial_\theta^2 l(\theta^*) - \partial_\theta^2 L(s^*; \theta^*) = \theta^* [\partial_\theta L(\tilde{S}(Z); \theta^*)^t \partial_\theta L(\tilde{S}(Z); \theta^*)].$$

Because matrices $-\partial_\theta^2 L(s^*; \theta^*)$ and $\partial_\theta^2 l(\theta^*) - \partial_\theta^2 L(s^*; \theta^*)$ are positive definite, $\partial_\theta T(\theta^*)$ is diagonalizable with strictly positive eigenvalues.

(iii) Here $\partial_\theta T(\theta^*)$ has the same eigenvalues (counting multiplicities) as the matrix A^* ,

$$A^* = I + B^*, \quad B^* = (-\partial_\theta^2 L(s^*; \theta^*))^{-1/2} \partial_\theta^2 l(\theta^*) (-\partial_\theta^2 L(s^*; \theta^*))^{-1/2}.$$

The proof is concluded by applying the Sylvester law of inertia [Horn and Johnson (1985)], showing that B^* has the same inertia (number of positive, negative and zero eigenvalues) as $\partial_\theta^2 l(\theta^*)$.

(iv) Using definitions (17) and (19) of the mappings T and G , we have

$$(47) \quad \partial_\theta T(\theta^*) = \partial_s \hat{\theta}(s^*) \partial_\theta \bar{s}(\theta^*), \quad \partial_s G(s^*) = \partial_\theta \bar{s}(\theta^*) \partial_s \hat{\theta}(s^*).$$

It follows from Theorem 1.3.20 in Horn and Johnson (1985) that $\partial_s G(s^*)$ has the same eigenvalues as $\partial_\theta T(\theta^*)$ counting multiplicities, together with additional $(m - p)$ eigenvalues equal to 0. The proof is concluded by noting that $\partial_s h(s^*) = \partial_s G(s^*) - I$. \square

Hence, θ^* is a stable stationary point of T if and only if θ^* is a proper maximizer of l . Otherwise, θ^* is hyperbolic or unstable. This is why most often the EM sequences do not converge toward saddle points or local minima [see Murray (1977)]. Of course, in the stochastic context, such ill-convergences are avoided automatically; roughly speaking, the stochastic approximation noise

prevents the convergence toward hyperbolic or unstable points. The following additional assumption is needed.

(LOC3) The minimum eigenvalue of the covariance matrix

$$R(\theta) = E_{\theta}[(\tilde{S}(Z) - \bar{s}(\theta))(\tilde{S}(Z) - \bar{s}(\theta))^t]$$

is bounded away from zero for θ in any compact subset $\mathcal{K} \subset \Theta$.

By application of Brandiere and Duflo (1996), the sequence $\{\theta_k\}$ converges w.p.1 to a proper maximizer of $l(\theta)$, under the assumptions of Theorem 5 and (LOC1)–(LOC3).

In the situations where the S_k converge to a regular stable stationary point [which is guaranteed under (LOC1)–(LOC3)], it makes sense to use an averaging procedure. In the application considered here, it is more interesting to compute the average on the parameters themselves. The averaging procedure takes the form

$$(48) \quad S_k = S_{k-1} + \gamma_k(\tilde{S}(Z_k) - S_{k-1}), \quad \theta_k = \hat{\theta}(S_k),$$

$$(49) \quad \bar{\theta}_{k+1} = \bar{\theta}_k + k^{-1}(\theta_k - \bar{\theta}_k),$$

where γ_k is a sequence such that $\lim_{k \rightarrow \infty} k^\alpha \gamma_k = \gamma_*$ and $\gamma_k/\gamma_{k+1} = 1 + O(k^{-1})$, for some $1/2 < \alpha < 1$. Adaptations of the procedure for the stabilized algorithm are along the same lines. To apply Theorem 4, we need to strengthen the assumptions on the stochastic perturbation $\{E_k\}$. Under (SAEM3), $\{E_k\}_{k \geq 1}$ is a martingale increment; CLT for martingale increments hold under Lyapunov type assumptions.

(SAEM4') For some $\alpha > 0$, $\theta \in \Theta$, $E_{\theta}[\|\tilde{S}(Z)\|^{2+\alpha}] < \infty$ and $\Gamma(\theta) \triangleq \text{Cov}_{\theta}(\tilde{S}(Z))$ is a continuous function of θ .

Under (SAEM1) and (SAEM4'), $n^{-1/2} \sum_{k=1}^n E_k \rightarrow_{\mathcal{L}} \mathcal{N}(0, \Gamma(\theta^*))$, and, by a direct application of Theorem 4, it holds that

$$(50) \quad \sqrt{n}(\bar{\theta}_n - \theta^*) \rightarrow_{\mathcal{L}} \mathcal{N}(0, \Sigma(\theta^*)),$$

$$(51) \quad \Sigma(\theta^*) = \partial_s \hat{\theta}(s^*) H(s^*)^{-1} \Gamma(\theta^*) H(s^*)^{-1} \partial_s \hat{\theta}(s^*)^t,$$

where $H(s^*) = \partial_s h(s^*)$. Some insight may be gained by working out the previous expression. To that purpose, note that, by application of (38) and (45),

$$\begin{aligned} \partial_s \hat{\theta}(s^*) H(s^*)^{-1} &= \partial_s \hat{\theta}(s^*) (\partial_{\theta} \bar{s}(\theta^*) \partial_s \hat{\theta}(s^*) - I)^{-1}, \\ &= (\partial_s \hat{\theta}(s^*) \partial_{\theta} \bar{s}(\theta^*) - I)^{-1} \partial_s \hat{\theta}(s^*), \\ &= -[\partial_{\theta}^2 l(\theta^*)]^{-1} \partial_{\theta}^2 L(s^*, \theta^*) \partial_s \hat{\theta}(s^*), \\ &= [\partial_{\theta}^2 l(\theta^*)]^{-1} \partial_{\theta} \phi(\theta^*)^t. \end{aligned}$$

Note that, by (46),

$$(52) \quad \text{Cov}_{\theta^*}(\partial_{\theta} L(\tilde{S}(Z); \theta^*)) = \partial_{\theta} \phi(\theta^*)^t \text{Cov}_{\theta^*}(\tilde{S}(Z)) \partial_{\theta} \phi(\theta^*),$$

$$(53) \quad = \partial_{\theta}^2 l(\theta^*) - \partial_{\theta}^2 L(s^*; \theta^*),$$

which finally implies that

$$(54) \quad \Sigma(\theta^*) = [\partial_\theta^2 l(\theta^*)]^{-1} [\partial_\theta^2 l(\theta^*) - \partial_\theta^2 L(s^*; \theta^*)] [\partial_\theta^2 l(\theta^*)]^{-1}.$$

In the missing data terminology, the asymptotic covariance of the averaged estimate is thus related to the *missing information*, the difference between the incomplete versus complete data Fisher information matrix. Note that an estimate of $\Sigma(\theta^*)$ can be recursively obtained from (11). We may summarize the above results as follows.

THEOREM 7. *Assume that (SAEM1)–(SAEM4') and (LOC1)–(LOC3) hold. Then, the sequence $\{\theta_n\}$ converges to some proper maximizer θ^* of $\theta \rightarrow l(\theta)$. If in addition, $\lim_{k \rightarrow \infty} k^\alpha \gamma_k = \gamma_*$ and $\gamma_k / \gamma_{k+1} = 1 + O(k^{-1})$, then $\sqrt{n}(\bar{\theta}_n - \theta^*) \mathbb{1}(\lim_{n \rightarrow \infty} \|\theta_n - \theta^*\| = 0)$ has a limiting distribution,*

$$(55) \quad \sqrt{n}(\bar{\theta}_n - \theta^*) \mathbb{1}\left(\lim_{n \rightarrow \infty} \|\theta_n - \theta^*\| = 0\right) \rightarrow_{\mathcal{L}} \mathcal{N}(0, \Sigma(\theta^*)) \mathbb{1}\left(\lim_{n \rightarrow \infty} \|\theta_n - \theta^*\| = 0\right),$$

where $\Sigma(\theta^*) = [\partial_\theta^2 l(\theta^*)]^{-1} [\partial_\theta^2 l(\theta^*) - \partial_\theta^2 L(s^*; \theta^*)] [\partial_\theta^2 l(\theta^*)]^{-1}$.

In situations where the dispersion of the estimator θ^* is given by $\partial_\theta^2 l(\theta^*)$ [this quantity can be approximated by (11)], the expression (55) gives a practical stopping rule, consisting of comparing an estimate of the simulation variance $\Sigma(\theta^*)/n$ (where n is the number of stochastic approximation cycles) with $\partial_\theta^2 l(\theta^*)$.

8. Extensions and discussions.

8.1. A stochastic approximation version of the ECM algorithm. It has long been noticed that there is a variety of problems where the maximum of $\theta \rightarrow L(s; \theta)$ is hard to compute. In many cases, however, the maximum of $L(s; \theta)$ restricted to particular subsets of Θ are in closed form and relatively easy to obtain. Motivated by this observation, Meng and Rubin (1993) have proposed a modified version of the EM algorithm, named ECM, which maintains the E-step but replaces the M-step by a cycle of conditional maximizations (CM-steps). For $c = 1, \dots, C$ and $s \in \mathcal{S}$ and $\eta \in \Theta$, define $\tilde{\theta}_c(s)$ as the value of θ that maximizes $\theta \rightarrow L(s; \theta)$ subject to the constraint $g_c(\theta) = g_c(\eta)$ (it is admitted at this point that such value exists and is unique). A full-cycle of CM-step is recursively defined by composing the conditional maximization: starting from a point $(s, \theta) \in \mathcal{S} \times \Theta$, this iteration takes the form for $1 \leq c \leq C$,

$$(56) \quad \hat{\theta}_c(s; \theta) \triangleq \tilde{\theta}_c(s; \hat{\theta}_{c-1}(s; \theta)), \quad \hat{\theta}_0(s; \theta) = \theta.$$

We denote $\hat{\theta}'(s; \theta) \triangleq \hat{\theta}_C(s; \theta)$; that is, the value of θ after a full cycle of conditional maximization. From a global point of view, we have replaced the M-step consisting of the determination of the global maximum $\hat{\theta}(s)$ of the function $L(s; \theta)$ by a cycle of conditional maximizations, which, starting from some point; $(s; \theta) \in \mathcal{S} \times \Theta$ leads us to some point $\hat{\theta}'(s; \theta) \in \Theta$. Note that, contrary

to $\hat{\theta}'(s; \theta)$, $\hat{\theta}(s)$ does not depend on the *starting point*; this simplifies the analysis of the EM algorithm. Not surprisingly, this has some implications for the specification and for the convergence analysis of the stochastic version of the EM algorithm. The ECM algorithm for the curved exponential model may be compactly written as

$$(57) \quad \theta_{k+1} = T'(\theta_k) = \hat{\theta}'(\bar{s}(\theta_k); \theta_k).$$

The ECM algorithm belongs to the class of generalized EM algorithms. The ECM algorithm is monotone in the sense that at each step the incomplete data likelihood function is increased: $l(\theta_{k+1}) \geq l(\theta_k)$. The convergence of the ECM has been thoroughly analyzed by Meng (1994); it is shown in these contributions that the ECM sequence is convergent as soon as the set of constraints verify a *space-filling* property; at any point θ , the convex hull of all feasible directions determined by the constraint spaces is the whole Euclidean space \mathbb{R}^p (the resulting maximization by repeated conditional maximizations is over the whole space and not a subspace of it). This space-filling condition in fact guarantees that the sequence of iterates $\hat{\theta}^{(i)}(s; \theta)$, defined recursively as

$$(58) \quad \hat{\theta}^{(i)}(s; \theta) = \hat{\theta}'(s; \hat{\theta}^{(i-1)}(s; \theta)), \quad \hat{\theta}^{(0)}(s; \theta) = \theta,$$

converges to $\hat{\theta}(s)$, and that the order of convergence is linear. To state the results in full generality, we start from this property, rather than from the *space-filling* assumption. We will consider the following assumptions.

(MAX1) For all $s \in \mathcal{S}$, $\theta \rightarrow L(s; \theta)$ has a unique global maximum on Θ , denoted $\hat{\theta}(s)$; moreover, $\hat{\theta}(s)$ is continuously differentiable on \mathcal{S} .

(MAX2) There exists a function $\hat{\theta}': \mathcal{S} \times \Theta \rightarrow \Theta$, a function $C: \mathcal{S} \times \Theta \rightarrow [0, \infty)$ and a function $\rho: \mathcal{S} \rightarrow [0, 1)$, such that, for all $(s; \theta) \in \mathcal{S} \times \Theta$, and all $i \in \mathbb{N}$,

$$(59) \quad \|\hat{\theta}^{(i)}(s; \theta) - \hat{\theta}(s)\| \leq C(s; \theta)\rho(s)^i,$$

where $\{\hat{\theta}^{(i)}(s; \theta)\}$ are the iterates of $\hat{\theta}(s; \theta)$ defined recursively as

$$(60) \quad \hat{\theta}^{(i)}(s; \theta) = \hat{\theta}'(s; \hat{\theta}^{(i-1)}(s; \theta)) \quad \text{for } i > 0, \quad \hat{\theta}^{(0)}(s; \theta) = \theta.$$

Moreover, for all $i \geq 0$, the functions $\hat{\theta}^{(i)}: \mathcal{S} \times \Theta \rightarrow \Theta$ are continuous and the functions C and ρ are bounded on the compact subsets of $\mathcal{S} \times \Theta$ and \mathcal{S} , respectively.

In many models of practical interests, the function L is for a given $s \in \mathcal{S}$ strictly convex w.r.t. θ , and the feasible set Θ also is convex. In that case, (MAX1) is generically verified, while (MAX2) holds, under mild assumptions, for a large variety of iterative maximization procedures, including generalized coordinate ascent methods (as above), steepest ascent methods and/or Gauss–Newton type algorithms. In the latter case, $\hat{\theta}^{(i)}(s; \theta)$ are the iterates of the function $\hat{\theta}'(s; \theta)$ defined as

$$(61) \quad \hat{\theta}'(s; \theta) = \theta + \alpha(s; \theta) d(s; \theta), \quad d(s; \theta) \triangleq [\partial_\theta^2 L(s; \theta)]^{-1} \partial_\theta L(s; \theta),$$

where $\alpha(s; \theta)$ maximizes $\alpha \rightarrow L(s + \alpha d(s; \theta))$ on the interval $[0, 1]$. Gauss–Newton type iterates have been considered, in incomplete data problems by Lange (1995), who proposed a gradient-type algorithm equivalent to the EM algorithm using basically (61).

In such a context, the SAEM algorithm may be adapted as follows:

$$(62) \quad S_k = S_{k-1} + \gamma_k(\tilde{S}(Z_k) - S_{k-1}), \quad \theta_k = \hat{\theta}(S_k; \theta_{k-1}),$$

where $\{Z_k\}$ is simulated under $p(z; \theta_{k-1})$ conditionally to θ_k independently from the past. The previous equation may be reexpressed [see (33) and (33)] as

$$(63) \quad S_k = S_{k-1} + \gamma_k h(S_{k-1}; \theta_{k-1}) + \gamma_k E_k,$$

where

$$(64) \quad h(s; \theta) = E_\theta[\tilde{S}(Z)] - s = \bar{s}(\theta) - s,$$

$$(65) \quad E_k = \tilde{S}(Z_k) - E_{\theta_{k-1}}(\tilde{S}(Z)) = \tilde{S}(Z_k) - \bar{s}(\theta_{k-1}).$$

Note that the mean field $h(s)$ of the original SAEM algorithm [see (33)] is given by $h(s) \triangleq h(s; \hat{\theta}(s))$. Because of the dependence in θ , (63) is not in the standard RM form, and the results obtained previously cannot be directly applied. We will use, to prove the convergence of this scheme, a *state perturbation approach*. Define, for $(s, \theta) \in \mathcal{S} \times \Theta$ and $m \geq n \geq 0$,

$$(66) \quad v_{n,m}(s; \theta) \triangleq \sum_{i=0}^m \gamma_{i+n+1} [h(s; \hat{\theta}^{(i)}(s; \theta)) - h(s)]$$

and denote $v_n(s; \theta) = \lim_{m \rightarrow \infty} v_{n,m}(s; \theta)$ (the convergence of the series is shown in Lemma 7). Define now the *perturbed state* as

$$(67) \quad \tilde{S}_n = S_n + v_n(S_n; \theta_n).$$

By the definitions of (63) and (66), the *perturbed state algorithm* may be written as

$$(68) \quad \tilde{S}_k = \tilde{S}_{k-1} + \gamma_k h(S_{k-1}; \theta_{k-1}) + \gamma_k E_k + v_k(S_k; \theta_k) - v_{k-1}(S_{k-1}; \theta_{k-1}).$$

The last term on the right can be expanded as

$$(69) \quad \begin{aligned} &v_k(S_k; \theta_k) - v_{k-1}(S_{k-1}; \theta_{k-1}) \\ &= -\gamma_k [h(S_{k-1}; \theta_{k-1}) - h(S_{k-1})] + v_k(S_k; \theta_k) \\ &\quad - v_k(S_{k-1}; \hat{\theta}(S_{k-1}; \hat{\theta}(S_{k-1}; \theta_{k-1}))). \end{aligned}$$

Using this latter expression, we may express the perturbed state equation (68) as

$$(70) \quad \tilde{S}_k = \tilde{S}_{k-1} + \gamma_k h(\tilde{S}_{k-1}) + \gamma_k E_k + \gamma_k R_k,$$

where the remainder term R_k is defined as follows:

$$(71) \quad R_k = h(S_{k-1}) - h(\tilde{S}_{k-1}) + \gamma_k^{-1} \left[v_k(S_k; \theta_k) - v_k(S_{k-1}; \hat{\theta}(S_{k-1}; \theta_{k-1})) \right].$$

This form illustrates the value of the perturbed state approach. The use of the perturbation removes $h(S_{k-1}; \theta_{k-1})$ and replaces it by $h(\tilde{S}_{k-1})$. If we may prove that $R_k = o(1)$, then the perturbed state equation will have the same mean field as the original SAEM algorithm and will thus converge to the same set of stationary points. Then, by adapting Theorem 6.2 (using Theorem 5.1), we will show that

$$(72) \quad \lim_{n \rightarrow \infty} d(\tilde{S}_n, \{s \in \mathcal{S} : \partial_s V(s) = 0\}) = 0 \quad \text{w.p.1.}$$

This property will in turn imply the convergence of $\{S_n\}$, because Lemma 7 implies that $\tilde{S}_n - S_n = O(\gamma_n)$ w.p.1 (under the compactness assumption). The discussion is summarized in the following theorem.

THEOREM 8. *Assume that assumptions (M1)–(M4), (SAEM1)–(SAEM4) and (MAX1)–(MAX2) hold. Assume in addition that w.p.1, (A) $\text{clos}(\{S_k\}_{k \geq 1})$ is a compact subset of \mathcal{S} . Then, w.p.1, $\lim_{k \rightarrow \infty} d(S_k, \{s \in \mathcal{S} : \partial_s V(s) = 0\}) = 0$ and $\lim_{k \rightarrow \infty} d(\hat{\theta}(S_k), \mathcal{L}) = 0$.*

The proof is in Appendix C. All the issues on the stability, the convergence and the control of the algorithm can be adapted (using the state perturbation approach). We do not pursue this discussion here.

8.2. Stochastic approximation EM versus stochastic gradient approaches.

One of the oldest and most widely known methods for maximizing a function of several variables is the method of steepest ascent (often referred to as the gradient method). The method of steepest ascent is defined by the iterative algorithm,

$$(73) \quad \theta_k = \theta_{k-1} + \alpha_k g_k, \quad g_k = \partial_\theta l(\theta_{k-1}),$$

where α_k is a nonnegative scalar [generally chosen as the maximum of the function $l(\theta_k + \alpha g_k)$]. In the incomplete data models, the basic steepest ascent method (73) cannot be directly applied, because $\partial_\theta l(\theta)$ is not in closed form and/or is hard to compute. A stochastic approximation version of the steepest ascent has been proposed for incomplete data models by Younes (1989, 1992); the basic version of this algorithm may be written as

$$(74) \quad \theta_k = \theta_{k-1} + \gamma_k \partial_\theta \log f(Z_k; \theta_{k-1}),$$

where the *missing* data Z_k is imputed from the current predictive distribution $p(z; \theta_{k-1})$ [conditionally to θ_{k-1} independently from the past; see (SAEM3)]. Algorithm (74) shares with the SAEM algorithm **the same imputation step**. It **differs** however in the way this imputation step is used to update the param-

eters. Algorithm (74) may be written in RM form as

$$(75) \quad \theta_k = \theta_{k-1} + \gamma_k h(\theta_{k-1}) + \gamma_k E_k,$$

with the following definitions:

$$h(\theta) = E_\theta[\partial_\theta \log f(Z; \theta)],$$

$$E_k = \partial_\theta \log f(Z_k; \theta_{k-1}) - h(\theta_{k-1}).$$

Note that (under the stated assumptions), E_k is a martingale difference. Under standard regularity assumptions, we have

$$(76) \quad h(\theta) = \partial_\theta l(\theta),$$

a relation which is often referred as **Fisher's identity**. In that sense, algorithm (74) is a *stochastic gradient algorithm*, since its mean field is nothing but the gradient of the objective function to maximize. The analysis of the stochastic gradient algorithm (74) can be done using the tools presented in this contribution. The Lyapunov function $V(\theta)$ [needed to check (SA2)] trivially is the negated incomplete data log-likelihood $V(\theta) \triangleq -l(\theta)$,

$$\langle \partial_\theta V(\theta), h(\theta) \rangle = -\|\partial_\theta l(\theta)\|^2.$$

The set of stationary points of (74), $\{\theta \in \Theta: \langle \partial_\theta V(\theta), h(\theta) \rangle = 0\}$, coincides with the set of stationary points of the incomplete data log-likelihood: $\mathcal{L} = \{\theta \in \Theta: \partial_\theta l(\theta) = 0\}$. By application of Sard's theorem, $\text{int}(l(\mathcal{L})) = \emptyset$ provided that $l: \Theta \rightarrow \mathbb{R}$ is p times continuously differentiable. A basic convergence theorem similar to Theorem 5 can be adapted from Theorem 2. The other issues (stability, reprojection, rate of convergence...) can be addressed along the same lines as above.

When the stationary points of $\theta \rightarrow l(\theta)$ are isolated, the procedure (74) converges pointwise to a proper maximizer of $\theta \rightarrow l(\theta)$, under mild assumptions [similar to (LOC1)–(LOC3)]. In such a case, it makes sense to consider the averaged sequence

$$(77) \quad \bar{\theta}_n = \bar{\theta}_{n-1} + n^{-1}(\theta_n - \bar{\theta}_{n-1}),$$

where θ_n is given by (74), the step size sequence $\{\gamma_k\}$ being chosen in such a way that $\lim_{k \rightarrow \infty} k^\alpha \gamma_k = \gamma_*$ and $\gamma_k/\gamma_{k+1} = 1 + O(k^{-1})$. It follows under basically the same assumptions as Theorem 7 that $\sqrt{n}(\bar{\theta}_n - \theta^*) \mathbb{1} \lim \|\theta_n - \theta^*\| = 0$) has a limiting distribution. This limiting distribution can be deduced almost directly from Theorem 4. Note that $\partial_\theta h(\theta^*) = \partial_\theta^2 l(\theta^*)$; (46) implies that $\text{Cov}_{\theta^*}(\partial_\theta L(\tilde{S}(Z); \theta^*)) = \partial_\theta^2 l(\theta^*) - \partial_\theta^2 L(s^*; \theta^*)$ and thus the asymptotic covariance of the averaged estimates is identical to (54). Thus, the averaged version of the MCEM algorithm and of the gradient algorithm have the same asymptotic rate of covariance \sqrt{n} and the same asymptotic covariance matrix. The choice between these two algorithms is dictated by about the same considerations as for their deterministic counterparts.

APPENDIX A

Convergence of stochastic approximation.

PROOF OF THEOREM 2. We prove the results on a sample path by sample path basis. In all the statements below, the qualifier “w.p.1” is implicit. The function F is upper semicontinuous and nonpositive; the set $\mathcal{L} \triangleq \{s \in \mathcal{S}: F(s) = 0\}$ is closed. Define

$$s'_n = s_n + \sum_{i=n+1}^{\infty} \gamma_i e_i.$$

Rewriting (21) yields

$$(78) \quad s'_n = s'_{n-1} + \gamma_n h(s_{n-1}) + \gamma_n r_n.$$

Condition (SA3) implies that there exists a compact set $\mathcal{K} \subset \mathcal{X}$ such that $s_n \in \mathcal{K}$, for all $n \geq 0$ (note that this compact set depends upon the trajectory). For $\alpha > 0$, define $\mathcal{K}_\alpha = \{s \in \mathbb{R}^m: d(s, \mathcal{K}) \leq \alpha\}$, where $d(s, \mathcal{K})$ denotes the distance from s to \mathcal{K} . Since \mathcal{X} is an open set, there exists $\rho > 0$ such that $\mathcal{K}_\rho \subset \mathcal{X}$. Since (i) under (SA4), $\sum_{i=n+1}^{\infty} \gamma_i e_i \rightarrow 0$; and $r_n \rightarrow 0$; (ii) under (SA2), the mean field h is bounded on \mathcal{K} , this implies

$$(79) \quad \|s'_n - s'_{n-1}\| = O(\gamma_n)$$

and the segment $[s'_n, s'_{n-1}] \subset \mathcal{K}_\rho$, for n sufficiently large. Under (SA2), the Lyapunov function V is continuously differentiable on \mathcal{X} . For n sufficiently large, there exists $s''_{n-1} \in [s'_n, s'_{n-1}]$ such that

$$\begin{aligned} V(s'_n) &= V(s'_{n-1}) + \langle \partial_s V(s''_{n-1}), (s'_n - s'_{n-1}) \rangle \\ &= V(s'_{n-1}) + \gamma_n F(s'_{n-1}) + \gamma_n r'_n \end{aligned}$$

with

$$\begin{aligned} r'_n &= \langle \partial_s V(s'_{n-1}), h(s_{n-1}) - h(s'_{n-1}) \rangle + \langle \partial_s V(s''_{n-1}) - \partial_s V(s'_{n-1}), h(s_{n-1}) \rangle \\ &\quad + \langle \partial_s V(s''_{n-1}), r_n \rangle. \end{aligned}$$

Since h and $\partial_s V(s)$ is continuous on \mathcal{K}_ρ and \mathcal{K}_ρ is compact h and $\partial_s V(s)$ are uniformly continuous on \mathcal{K}_ρ ; thus, $r'_n = o(1)$ and

$$(80) \quad V(s'_n) = V(s'_{n-1}) + \gamma_n F(s'_{n-1}) + o(\gamma_n).$$

Note that the function F is bounded on \mathcal{K}_ρ and (80) in particular implies that there exists $C_V < \infty$ such that for n sufficiently large,

$$(81) \quad |V(s'_n) - V(s'_{n-1})| \leq C_V \gamma_n.$$

The proof proceeds in two steps: we first show that the sequence $\{V(s'_n)\}_{n \geq 0}$ converges to some point of $V(\mathcal{L})$ (Step 1). We then show that the sequence $\{s'_n\}_{n \geq 0}$ converges to \mathcal{L} (Step 2).

STEP 1. For $\alpha > 0$, define $\mathcal{A}_\alpha = \{x \in \mathbf{R}: d(x, V(\mathcal{L} \cap \mathcal{K}_\rho)) \leq \alpha\}$. Here $V(\mathcal{L} \cap \mathcal{K}_\rho)$ is a compact subset of \mathbb{R} and \mathcal{A}_α is compact; \mathcal{A}_α is a finite union of disjoint closed intervals (of length greater than 2α),

$$(82) \quad \mathcal{A}_\alpha = \bigcup_{i=1}^{N_\alpha} [a_\alpha^{(i)}, b_\alpha^{(i)}] \quad \text{where } a_\alpha^{(1)} < b_\alpha^{(1)} < a_\alpha^{(2)} < \dots < b^{(N_\alpha)}.$$

Set $\mathcal{N}_\alpha = V^{-1}[\text{int}(\mathcal{A}_\alpha)]$. Since V is continuous, \mathcal{N}_α is an open neighborhood of the set $\mathcal{L} \cap \mathcal{K}_\rho$. Since F is upper semicontinuous (any upper semicontinuous function reaches its maximum on any compact set), there exists $\varepsilon_\alpha > 0$ such that for n sufficiently large,

$$(83) \quad s'_{n-1} \in \mathcal{K}_\rho \setminus \mathcal{N}_\alpha \implies F(s'_{n-1}) \leq -2\varepsilon_\alpha.$$

This implies, using (80) and (81),

$$(84) \quad V(s'_n) \leq V(s'_{n-1}) - \gamma_n \varepsilon_\alpha + \gamma_n (C_V + \varepsilon_\alpha) \mathbb{1}(V(s'_{n-1}) \in \text{int}(\mathcal{A}_\alpha))$$

Since $\sum \gamma_n = \infty$, (84) implies that $\{V(s'_n)\}$ is infinitely often in \mathcal{A}_α . Here \mathcal{A}_α is a finite union of disjoint intervals, thus $\{V(s'_n)\}$ is infinitely often in a given interval of the partition (82), say $[a_\alpha, b_\alpha]$ (the superscript is implicit).

Let $\delta > 0$ be such that $[a_\alpha - 2\delta, b_\alpha + 2\delta] \cap \mathcal{A}_\alpha = [a_\alpha, b_\alpha]$ (such δ exists because \mathcal{A}_α is a finite union of closed disjoint intervals), and let N_V be such that, for all $n \geq N_V$, $\gamma_n \leq \delta/C_V$. Finally, let N be the first index greater than N_V such that $V(s'_N) \in [a_\alpha, b_\alpha]$ [such index exists because $V(s'_n)$ is i.o. in $[a_\alpha, b_\alpha]$]. Equation (84) implies that, for all $n \geq N$, $V(s'_n) \in [a_\alpha - \delta, b_\alpha + \delta]$. Hence, δ being arbitrary, the set of limit points \mathcal{J} of the sequence $\{V(s'_n)\}_{n \geq 0}$ is included in the interval $[a_\alpha, b_\alpha]$. Take now $\alpha' < \alpha$. Proceeding as above, we may show that $\mathcal{J} \subset [a_{\alpha'}, b_{\alpha'}]$, where $[a_{\alpha'}, b_{\alpha'}]$ is a given element of the partition (82) of $\mathcal{A}_{\alpha'}$. Of course, we must have $[a_{\alpha'}, b_{\alpha'}] \subset [a_\alpha, b_\alpha]$. Applying the above result with a nonincreasing sequence $\alpha_n \rightarrow 0$, we may show that \mathcal{J} is included in the intersection of a decreasing sequence of closed intervals, and \mathcal{J} is thus itself a closed interval. Since all the intervals $[a_{\alpha_n}, b_{\alpha_n}]$ are included in \mathcal{A}_{α_n} and $\bigcap_n \mathcal{A}_{\alpha_n} = V(\mathcal{L} \cap \mathcal{K}_\rho)$, \mathcal{J} also is a subset of $V(\mathcal{L} \cap \mathcal{K}_\rho)$. The set $V(\mathcal{L} \cap \mathcal{K}_\rho)$ has an empty interior, and the connected components of $V(\mathcal{L} \cap \mathcal{K}_\rho)$ are reduced to points. Since \mathcal{J} is connected, it must be reduced to a point, which implies that $\lim_{n \rightarrow \infty} V(s'_n)$ exists. The proof of Step 1 is concluded by noting that $\lim_{n \rightarrow \infty} V(s'_n) = \lim_{n \rightarrow \infty} V(s_n)$.

STEP 2. Let \mathcal{N} be an arbitrary open neighborhood of \mathcal{L} . By (80), (83), there exists $\varepsilon_{\mathcal{N}} > 0$ such that for n sufficiently large,

$$(85) \quad V(s'_n) \leq V(s'_{n-1}) - \gamma_n \varepsilon_{\mathcal{N}} + \gamma_n (C_V + \varepsilon_{\mathcal{N}}) \mathbb{1}(s'_{n-1} \in \mathcal{N}).$$

Since $V(s'_n)$ converges and $\gamma_n = o(1)$, for any $\tau > 0$, there exists $N_\tau < \infty$ such that

$$(86) \quad |V(s'_n) - V(s'_p)| \leq \tau, \quad N_\tau \leq n \leq p \quad \text{and} \quad \gamma_k < \tau/\varepsilon_{\mathcal{N}}, N_\tau \leq k.$$

Equations (85) and (86) imply that

$$(87) \quad \varepsilon_{\mathcal{N}} \sum_{k=n+1}^p \gamma_k - (C_V + \varepsilon_{\mathcal{N}}) \sum_{k=n+1}^p \gamma_k \mathbb{1}(s'_{k-1} \in \mathcal{N}) \leq \tau, \quad N_\tau \leq n \leq p.$$

Choose $n > N_\tau$ and let $p(n)$ be the first integer larger than n such that

$$\frac{\tau}{\varepsilon_{\mathcal{N}}} < \sum_{k=n+1}^{p(n)} \gamma_k \leq 2 \frac{\tau}{\varepsilon_{\mathcal{N}}}.$$

Equation (87) implies that there exists $k(n) \in [n, p(n)]$ such that $s'_{k(n)} \in \mathcal{N}$. By (79), there exists a finite constant K such that, for all $n \geq 1$, $\|s'_n - s'_{n-1}\| \leq K\gamma_k$, which implies

$$\|s'_{k(n)} - s'_n\| \leq K \sum_{l=n+1}^{k(n)} \gamma_l \leq K \sum_{l=n+1}^{p(n)} \gamma_l \leq 2K \frac{\tau}{\varepsilon_{\mathcal{N}}}.$$

By construction, $s'_{k(n)} \in \mathcal{N}$, and the later relation implies that $d(s'_{k(n)}, \text{clos}(\mathcal{N})) \leq 2\tau K/\varepsilon_{\mathcal{N}}$. Since τ is arbitrary, this implies that

$$\lim_{n \rightarrow \infty} d(s'_n, \mathcal{N}) = 0.$$

Since \mathcal{N} is arbitrary, $\{s'_n\}$ tends to \mathcal{L} and so does $\{s_n\}$. The set of limit points of $\{s_n\}$ is compact because it is bounded and closed; it is connected because by (79), $\|s_{n+1} - s_n\| \rightarrow 0$. \square

PROOF OF THEOREM 3. The truncated RM procedure (23) can be written in the RM form (21),

$$(88) \quad s_k = s_{k-1} + \gamma_k h(s_{k-1}) + \gamma_k e_k + \gamma_n r'_k,$$

where the remainder term r'_k is defined as

$$(89) \quad \text{if } \hat{s}_k \in \mathcal{K}_{n_{k-1}}: \begin{cases} r'_k = r_k, \\ n_k = n_{k-1}, \end{cases}$$

$$(90) \quad \text{if } \hat{s}_k \notin \mathcal{K}_{n_{k-1}}: \begin{cases} r'_k = \gamma_k^{-1}(s'_k - s_{k-1}) - h(s_{k-1}) - e_k, \\ n_k = n_{k-1} + 1. \end{cases}$$

We use Lemma 1 to check that the assumptions of Theorem 2 are satisfied. Note that the stochastic approximation procedure defined (88) satisfies (SA0)–(SA2), (STAB1) (because the projected algorithm has the same mean field h as the original procedure) and also (STAB3) (because the algorithm is by construction recurrent in a compact subset). Also, the first condition (STAB2) is directly verified; we have only to check that, for any integer M ,

$$(91) \quad \lim_{k \rightarrow \infty} \|r'_k\| \mathbb{1}(W(s_{k-1}) \leq M) = 0.$$

If $\hat{s}_k \in \mathcal{X}_{n_{k-1}}$, then $r'_k = r_k$, and (89) shows that

$$\limsup_{k \rightarrow \infty} \|r'_k\| \mathbb{1}(\{W(s_{k-1}) \leq M\}) \mathbb{1}(\hat{s}_k \in \mathcal{X}_{n_{k-1}}) = 0.$$

Thus, we must show that $\mathbb{1}(\hat{s}_k \notin \mathcal{X}_{n_{k-1}}) \mathbb{1}(W(s_{k-1}) \leq M) = 0$ for all but a finite number of indexes k . Since $\{s: W(s) \leq (M + 1)\}$ is a compact subset and $\bigcup \mathcal{X}_n = \mathcal{X}$, there exists $N_M < \infty$ such that $\{s: W(s) \leq (M + 1)\} \subset \mathcal{X}_n$ for all $n \geq N_M$. Under the stated assumptions, $\{s: W(s) \leq M\}$ is a compact set; thus, there exists w.p.1 an index K_M such that, for all $k \geq K_M$,

$$W(s_{k-1}) \leq M \rightarrow W(\hat{s}_k) > (M + 1).$$

Since, for all $n \geq N_M$, we have $\{s \notin \mathcal{X}_n\} \subset \{W(s) > (M + 1)\}$, the previous relation implies that, w.p.1, for all $k \geq K_M$ and all $n \geq N_M$,

$$W(s_{k-1}) \leq M \rightarrow \hat{s}_k \notin \mathcal{X}_n.$$

Thus, w.p.1, $n_k \leq N_M$ and $\mathbb{1}(\{\hat{s}_k \notin \mathcal{X}_{n_{k-1}}\}) \mathbb{1}(W(s_{k-1}) \leq M) = 0$ only for a finite index number. This completes the proof of (91). \square

APPENDIX B

Rate of convergence and averaging of estimates. This Appendix is devoted to some technical results needed to derive the limiting distribution of the averaging estimates. These topics are intimately related, because basically, the averaging procedure can be fruitfully applied as soon as the *primary algorithm* converges to a regular stationary point with a given rate of convergence. Rates of convergence and averaging of estimates are extensively studied in Kushner and Yin (1997), Chapter 10 and 11 (see also the references therein). Our main contribution is to weaken some of the assumptions used in this reference and in particular the tightness condition [(A3.1), Kushner and Yin (1997), Chapter 11, page 336]. The improvement mainly stems from the decomposition of the error used in Lemma 6, (92) and (93). We preface this Appendix by some additional notations, definitions and technical results.

DEFINITIONS AND NOTATION. For A a square matrix, denote $\|A\|_2$ the spectral norm. Let $\{a_n\}_{n \geq 0}$ be a sequence of positive real numbers. We say that $X_n = O_P(a_n)$ if $a_n^{-1} X_n$ is bounded in probability and $X_n = o_P(a_n)$ if $a_n^{-1} X_n$ converges to zero in probability. We say that $X_n = O_q(a_n)$ (where $q > 0$) if $\sup_{n \in \mathbb{N}} \|a_n^{-1} X_n\|_q < \infty$. We say that $X_n = O(a_n)$ w.p.1. [resp., $X_n = o(a_n)$ w.p.1] if $\sup_{n \in \mathbb{N}} a_n^{-1} \|X_n\|$ is finite w.p.1 (resp., $\limsup a_n^{-1} \|X_n\| = 0$). Similarly, let Y_n be a sequence of random variables. We say that $X_n = O(Y_n)$ w.p.1 if $Y_n^{-1} X_n = O(1)$ w.p.1.

The two following lemmas play a key role in the sequel (proofs are omitted for brevity).

LEMMA 4. Let $\{\gamma_k\}_{k \geq 0}$ be a positive sequence such that $\lim_{k \rightarrow \infty} \gamma_k = 0$, $\lim_{k \rightarrow \infty} \gamma_k^{-1} - \gamma_{k-1}^{-1} = 0$. Let $\{\varepsilon_k\}_{k \geq 0}$ be a nonnegative sequence. For $b > 0$ and $p \geq 0$, define

$$\sigma_n = \sum_{i=1}^n \gamma_i^{p+1} \exp(-b\bar{\gamma}(i, n))\varepsilon_i,$$

where $\bar{\gamma}(i, n) = \sum_{j=i+1}^n \gamma_j$ for $0 \leq i < n$. Then, $\limsup \gamma_k^{-p} \sigma_k \leq b^{-1} \limsup \varepsilon_k$.

LEMMA 5. Let $\{A_k\}_{k \geq 0}$ be a sequence of square matrix such that $\lim_{k \rightarrow \infty} \|A_k - A\|_2 = 0$ where A is a Hurwitz matrix. Let $\{\gamma_k\}_{k \geq 0}$ be a sequence of positive number such that $\lim_{k \rightarrow \infty} \gamma_k = 0$. Then there exist $\beta > 0$ and a constant $C_\beta < \infty$, such that, for all $m > n \geq 0$, it holds that

$$\left\| \prod_{k=m+1}^n (I + \gamma_k A_k) \right\|_2 \leq C_\beta \exp(-\beta\bar{\gamma}(m, n)).$$

The crux for obtaining results on the rate of convergence and on the averaged sequence is to approximate the original nonlinear difference equation by a linear one and to bound the error incurred by the linearization. This is the purpose of the following lemma.

LEMMA 6. Assume that (AVE1), (AVE2) and (AVE3(α)) hold for some $0 < \alpha < 1$, and let s^* be a regular stable stationary point. Define $v_n = e_n + r_n$ and

$$(92) \quad \mu_n = (I + \gamma_n H(s^*))\mu_{n-1} + \gamma_n v_n, \quad \mu_0 = 0,$$

$$(93) \quad \rho_n = s_n - s^* - \mu_n.$$

Then

$$(94) \quad \mu_n \mathbb{1} \left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0 \right) =_{w.p.1} O(\gamma_n^{1/2})O_2(1) \quad \text{and} \quad \rho_n =_{w.p.1} O(\gamma_n)O_1(1).$$

REMARK 4. The expression $\mu_n \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0) = O(\gamma_n^{1/2})O_2(1)$ w.p.1 is a shorthand for $\|\mu_n\| \leq XY_n$ where X is a (nonnegative) random variable w.p.1 finite and $\{Y_n\}_{n \geq 0}$ is a sequence of r.v. bounded in L^2 . The same notational convention holds for ρ_n . Note that μ_n can be seen as a leading term in the expansion of the error $s_n - s^*$, whereas ρ_n is a remainder term.

PROOF LEMMA 6. By Taylor’s expansion at s^* and since $h(s^*) = 0$, it holds, for n sufficiently large,

$$s_n - s^* = s_{n-1} - s^* + \gamma_n H^*(s_{n-1} - s^*) + \gamma_n v_n + \gamma_n \varepsilon_n,$$

where $H^* \triangleq H(s^*)$ and $\varepsilon_n \triangleq [\varepsilon_{n,1}, \dots, \varepsilon_{n,m}]^t$ is defined component-wise as

$$\varepsilon_{n,i} = \sum_{k,l=1}^m R_{n,i}(k,l)(s_{n-1,k} - s_k^*)(s_{n-1,l} - s_l^*),$$

$$R_{n,i}(k,l) = \int_0^1 \frac{(1-t)^2}{2!} \frac{\partial^2 h_i}{\partial s_k \partial s_l}(s_{n-1} + t(s^* - s_{n-1})) dt.$$

Using definition (92), we have

$$\rho_n = (I + \gamma_n H^*)\rho_{n-1} + \gamma_n \varepsilon_n = (I + \gamma_n H_n)\rho_{n-1} + \gamma_n \varepsilon'_n,$$

where $H_n = ([H_n(i,j)])_{1 \leq i,j \leq m}$ and $\varepsilon'_n = [\varepsilon'_{n,1}, \dots, \varepsilon'_{n,m}]^t$ are, respectively, defined as

$$H_n(i,j) \triangleq H^*(i,j) + \sum_{k=1}^m (2R_{n,i}(j,k)\mu_{n-1,k} + R_{n,i}(j,k)\rho_{n-1,k}),$$

$$\varepsilon'_{n,i} \triangleq \sum_{k,l=1}^m R_{n,i}(k,l)\mu_{n-1,k}\mu_{n-1,l}.$$

Denote for $n > k$, $\psi_*(n,k) = (I + \gamma_n H^*) \cdots (I + \gamma_{k+1} H^*)$ ($\psi_*(n,n) = I$ and $\psi_*(n,k) = 0$ if $k > n$). By Lemma 5, there exists $\beta > 0$ and $C'(H^*, \beta) < \infty$, such that, for all $n \geq k \geq 0$,

$$\|\psi_*(n,k)\|_2 \leq C'(H^*, \beta) \exp(-\beta\bar{\gamma}(n,k)).$$

Decompose μ_n [see (92)] as $\mu_n = \mu_n^{(0)} + \mu_n^{(1)}$, where

$$(95) \quad \mu_n^{(0)} = \sum_{k=1}^n \gamma_k \psi_*(n,k) e_k(\rho),$$

$$(96) \quad \mu_n^{(1)} = \sum_{k=1}^n \gamma_k \psi_*(n,k) r'_k,$$

where $e_k(\rho) = e_k \mathbb{1}(\|s_{k-1} - s^*\| \leq \rho)$ and $r'_k = r_k + e_k \mathbb{1}(\|s_{k-1} - s^*\| > \rho)$. Note that r'_k verifies (AVE2). Under (AVE1), we have

$$(97) \quad [E(\|\mu_n^{(0)}\|^2)]^{1/2} \leq C'(H^*, \beta) C^{1/2}(\rho) \left(\sum_{k=1}^n \gamma_k^2 \exp(-2\beta\bar{\gamma}(k,n)) \right)^{1/2}$$

and thus $\mu_n^{(0)} = O_2(\gamma_n^{1/2})$, by application of Lemma 4. An application of the Abel's transform shows that

$$(98) \quad \mu_n^{(0)} = \psi_*(n,1)E(n,1) + \sum_{k=2}^n \gamma_k \psi_*(n,k) H^* E(n,k),$$

where $E(n,k) \triangleq \sum_{j=k}^n \gamma_j e_j(\rho)$ for $k \leq n$, $E(n,k) = 0$ otherwise. Let δ be such that $\alpha < \delta < 1$. Under (AVE1), it holds that, w.p.1, $\sup_{(n,k) \in \mathbb{N} \times \mathbb{N}} \|E(n,k)\| < \infty$

and $\lim_{n \rightarrow \infty} \sup_{n - [n^\delta] \leq k \leq n} \|E(n, k)\| = 0$. Under (AVE3(α)), we have

$$\sum_{k=1}^{n - [n^\delta]} \gamma_k \|\psi_*(n, k)\| \leq C'(H^*, \beta)(n - [n^\delta]) \exp(-\beta[n^\delta] \gamma_{n - [n^\delta]})$$

and thus, w.p.1,

$$\begin{aligned} & \left\| \sum_{k=1}^{n - [n^\delta]} \gamma_k \psi_*(n, k) H^* E(n, k) \right\| \rightarrow 0 \quad \text{as } n \rightarrow \infty, \\ & \left\| \sum_{k=n - [n^\delta] + 1}^n \gamma_k \psi_*(n, k) H^* E(n, k) \right\| \\ & \leq \sup_{n - [n^\delta] \leq k \leq n} \|E(n, k)\| \|H^*\|_2 \sum_{k=1}^n \gamma_k \exp(-\beta \bar{\gamma}(k, n)) \rightarrow 0 \quad \text{as } n \rightarrow \infty, \end{aligned}$$

showing that $\mu_n^{(0)} - o(1)$ w.p.1. Similarly, we have

$$(99) \quad \|\mu_n^{(1)}\| \leq C'(H^*, \beta) \sum_{k=1}^n \gamma_k^{3/2} \exp(-\beta \bar{\gamma}(k, n)) \gamma_k^{-1/2} \|r'_k\|$$

and, by application of Lemma 4, $\mu_n^{(1)} \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0) = O(\gamma_n^{1/2})$ w.p.1. Equations (97) and (99) together imply that $\mu_n \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0) = O_2(\gamma_n^{1/2}) + O(\gamma_n^{1/2})$ w.p.1.

Define $\psi(n, k) = (I + \gamma_n H_n) \cdots (I + \gamma_{k+1} H_{k+1})$ for $n > k$ [$\psi(n, n) = I$, $\psi(n, k) = 0$, for $n < k$]. Since $\mu_n \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0) o(1)$ w.p.1, we have $\|H_n - H(s^*)\|_2 \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0) o(1)$ w.p.1. By Lemma 5, there exist $\beta > 0$ and a w.p.1 finite r.v. $C(H, \beta) < \infty$, such that, for all $0 \leq k \leq n$,

$$(100) \quad \|\psi(n, k)\|_2 \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0\right) \leq C(H, \beta) \exp(-\beta \bar{\gamma}(k, n)) \quad \text{w.p.1.}$$

Notice that

$$\rho_n = \sum_{k=1}^n \gamma_k \psi(n, k) \varepsilon'_k + \psi(n, 0) s_0$$

and $\varepsilon'_k \mathbb{1}(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0) = O(\|\mu_{k-1}\|^2)$ w.p.1; by applying (100) there exists a r.v. $K(\beta, \rho)$ w.p.1 finite such that, for all $n \geq 1$,

$$(101) \quad \begin{aligned} & \|\rho_n\| \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0\right) \\ & \leq K(\beta, \rho) \sum_{k=1}^n \gamma_k \exp(-\beta \bar{\gamma}(k, n)) \|\mu_{k-1}\|^2 \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0\right) \\ & \quad + K(\beta, \rho) \exp(-\beta \bar{\gamma}(0, n)) \|s_0 - s^*\| \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0\right). \end{aligned}$$

By application of Lemma 4, we have

$$(102) \quad E\left(\sum_{k=1}^n \gamma_k \exp(-\beta \bar{\gamma}(k, n)) \|\mu_{k-1}^{(0)}\|^2\right) = O(\gamma_n),$$

$$(103) \quad \sum_{k=1}^n \gamma_k \exp(-\beta \bar{\gamma}(k, n)) \|\mu_{k-1}^{(1)}\|^2 \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0\right) =_{\text{w.p.1}} O(\gamma_n),$$

$$(104) \quad \exp(-\beta \bar{\gamma}(0, n)) \|s_0 - s_*\| \mathbb{1}\left(\lim_{n \rightarrow \infty} \|s_n - s_*\| = 0\right) =_{\text{w.p.1}} O(\gamma_n)$$

which concludes the proof. \square

REMARK 5. This lemma serves for Theorem 4 the same purpose as Assumptions (A3.1)–(A3.8) for Theorem 3.1 in Kushner and Yin [(1997), pages 336–338]. Note that the relation (94) does not imply (A3.1) [restated with our notations, (A3.1) requires that $(s_n - s^*) \mathbb{1}(\|s_n - s^*\| \leq \rho) = O_2(\gamma_n^{1/2})$]. The assumptions used in the previous lemma are weaker than those required to show (A3.1).

PROOF OF THEOREM 4. Write $s_n - s^* = \mu_n + \rho_n$ where μ_n and ρ_n are defined in (92) and (93), respectively. Note that $\bar{s}_n - s^* = \bar{\mu}_n + \bar{\rho}_n$, where $\bar{\mu}_n$ and $\bar{\rho}_n$ denote Cesaro’s mean of μ_n and ρ_n ,

$$\bar{\mu}_n \triangleq n^{-1} \sum_{k=1}^n \mu_k \quad \text{and} \quad \bar{\rho}_n \triangleq n^{-1} \sum_{k=1}^n \rho_k.$$

The proof is in two steps: (i) we show that $\sqrt{n} \bar{\rho}_n = o_P(1)$ and (ii) we then show that $\sqrt{n}(\bar{\mu}_n - H(s^*)^{-1} \bar{\nu}_n) = o_P(1)$, where $\bar{\nu}_n$ is Cesaro’s mean of ν_n , $\bar{\nu}_n \triangleq n^{-1} \sum_{k=1}^n \nu_k$.

(i) Applying Lemma 6, (94), there exist a nonnegative w.p.1 bounded r.v. X and a sequence of nonnegative r.v.’s $\{Y_n\}_{n \in \mathbb{N}}$, such that $Y_n = O_1(1)$ and $\rho_k = \gamma_k X Y_k$ w.p.1. Property (i) follows from

$$\sqrt{n} \bar{\rho}_n \leq_{\text{w.p.1}} X \left(n^{-1/2} \sum_{i=1}^n \gamma_i Y_i \right) = O(1) O_1(n^{1/2-\alpha}).$$

(ii) Using (92), we may write, for $n \geq 1$, $\mu_{n-1} = H(s^*)^{-1}(\gamma_n^{-1}(\mu_n - \mu_{n-1}) - \nu_n)$. Thus,

$$\begin{aligned} \bar{\mu}_n &= H(s^*)^{-1} \left(-n^{-1} \sum_{k=1}^n \nu_{k+1} + n^{-1} \sum_{k=1}^n \gamma_{k+1}^{-1} (\mu_{k+1} - \mu_k) \right) \\ &= H(s^*)^{-1} \left(-n^{-1} \sum_{k=1}^n \nu_{k+1} - n^{-1} \gamma_2^{-1} \mu_1 + n^{-1} \gamma_{n+1}^{-1} \mu_{n+1} \right. \\ &\quad \left. + n^{-1} \sum_{k=2}^n (\gamma_k^{-1} - \gamma_{k+1}^{-1}) \mu_k \right). \end{aligned}$$

The proof is concluded by showing that $n^{-1/2} \sum_{k=2}^n (\gamma_k^{-1} - \gamma_{k+1}^{-1}) \mu_k = o_P(1)$. Using Lemma 6, (94), there exists a nonnegative r.v. X and a sequence of non-negative r.v.'s $\{Y_n\}$ such that $Y_n = O_2(1)$ and, for all $k \geq 1$, $\|\mu_k\| = \gamma_k^{1/2} X Y_k$ w.p.1. Under (AVE3(α)), $\gamma_k^{1/2} |\gamma_k^{-1} - \gamma_{k+1}^{-1}| \leq C_\gamma k^{\alpha/2-1}$, for some constant $C_\gamma < \infty$. Thus

$$n^{-1/2} \sum_{k=2}^n (\gamma_k^{-1} - \gamma_{k+1}^{-1}) \|\mu_k\| \leq_{\text{w.p.1}} C_\gamma X n^{-1/2} \sum_{k=2}^n k^{\alpha/2-1} Y_k$$

and the proof is concluded by noting that $n^{-1/2} \sum_{k=2}^n k^{\alpha/2-1} Y_k = O_1(n^{(\alpha-1)/2}) = o_P(1)$. \square

APPENDIX C

Miscellaneous results.

LEMMA 7. Assume (M1)–(M4) and (MAX1)–(MAX2). Then, for all $(s, \theta) \in \mathcal{S} \times \Theta$, $\sum_{i=0}^\infty \gamma_{i+n+1} \|h(s; \hat{\theta}^{(i)}(s; \theta)) - h(s)\| < \infty$ and

$$(105) \quad \|v_n(s; \theta)\| \leq C(s; \theta) [1 - \rho(s)]^{-1} \sup_{\{\theta' \in \Theta: \|\theta' - \hat{\theta}(s)\| \leq C(s; \theta)\}} \|\partial_\theta \bar{s}(\theta)\|_2 \gamma_{n+1},$$

where $v_n(s; \theta) = \lim_{m \rightarrow \infty} v_{n,m}(s; \theta)$. Moreover, the set of functions $\{\gamma_n^{-1} v_n(s; \theta)\}_{n \geq 0}$ is equicontinuous, that is, for any compact subsets $\mathcal{X}_\mathcal{S} \subset \mathcal{S}$ and $\mathcal{X}_\Theta \subset \Theta$ and for any $\varepsilon > 0$, there is $\delta > 0$, such that for all $n \geq 0$,

$$(106) \quad \sup_{\{s, s' \in \mathcal{X}_\mathcal{S}, \|s-s'\| \leq \delta\}} \sup_{\{\theta, \theta' \in \mathcal{X}_\Theta, \|\theta-\theta'\| \leq \delta\}} [\gamma_n^{-1} \|v_n(s; \theta) - v_n(s'; \theta')\|] \leq \varepsilon.$$

The proof of this lemma is straightforward and is omitted.

PROOF OF THEOREM 8. As mentioned above, we need to prove that (i) the remainder term $R_n = o(1)$ w.p.1. Under the compactness assumption (A(i)–(ii)), there exists w.p.1 a compact set $\mathcal{X}_\mathcal{S} \subset \mathcal{S}$, such that the sequence $\{S_n\}$ is in $\mathcal{X}_\mathcal{S}$ for all n (note that this compact set depends on the trajectory). Assumptions (MAX1) and (MAX2) then imply that there exists w.p.1 a compact set $\mathcal{X}_\Theta \subset \Theta$ such that $\{\theta_n\}_{n \in \mathbb{N}} \subset \mathcal{X}_\Theta$. Under (SAEM1)–(SAEM3), $\sum_{k=0}^n \gamma_k E_k$ converges w.p.1 (see Theorem 6.2). This implies that, w.p.1,

$$(107) \quad S_n - S_{n-1} = o(1),$$

$$(108) \quad \theta_n - \hat{\theta}(S_{n-1}; \theta_{n-1}) = \hat{\theta}(S_n; \theta_{n-1}) - \hat{\theta}(S_{n-1}; \theta_{n-1}) = o(1).$$

Using Lemma 7, we thus have under (MAX2) that $v_n(s_n; \theta_n) = O(\gamma_n)$ w.p.1, which implies, under (M3) and (MAX1) $h(S_n) - h(\hat{S}_n) = O(\gamma_n)$ w.p.1 [$h(s)$ is continuously differentiable on \mathcal{S}]. The proof of (i) then follows from the equicontinuity of the set of functions $\{\gamma_n^{-1} v_n(s; \theta)\}_{n \in \mathbb{N}}$ (Lemma 7). Equiva-

tions (107) and (108) together with (106) indeed imply that

$$\gamma_n^{-1}(v_n(S_n; \theta_n) - v_n(S_{n-1}; \hat{\theta}(S_{n-1}; \theta_{n-1}))) = o(1),$$

which concludes the proof. \square

REFERENCES

- ANDRADOTTIR, S. (1995). A stochastic approximation algorithm with varying bounds. *Oper. Res.* 1037–1048.
- BRANDIERE, O. and DUFLO, M. (1996). Les algorithmes stochastiques contournent-ils les pièges ? *Ann. Inst. H. Poincaré* **32** 395–427.
- BROCKER, T. (1975). *Differentiable Germs and Catastrophes*. Cambridge Univ. Press.
- CELEUX, G. and DIEBOLT, J. (1988). A probabilistic teacher algorithm for iterative maximum likelihood estimation. In *Classification and Related Methods of Data Analysis* (H. H. Bock, ed.) 617–623. North-Holland, Amsterdam.
- CELEUX, G. and DIEBOLT, J. (1992). A stochastic approximation type EM algorithm for the mixture problem. *Stochastics Stochastics Rep.* **41** 127–146.
- CHEN, H. F., GUO, L. and GAO, A. J. (1988). Convergence and robustness of the Robbins–Monro algorithm truncated at randomly varying bounds. *Stoch. Process Appl.* **27** 217–231.
- CHICKIN, D. O. and POZNYAK, A. S. (1984). On the asymptotical properties of the stochastic approximation procedure under dependent noise. *Automat. Remote Control* **44**.
- CHICKIN, D. O. (1988). On the convergence of the stochastic approximation procedure under dependent noise. *Automat. Remote Control* **48**.
- DE JONG, P. and SHEPHARD, N. (1995). The simulation smoother for time series model. *Biometrika* **82** 339–350.
- DELYON, B. (1996). General results on stochastic approximation. *IEEE Trans. Automat. Control*. To appear.
- DELYON, B. and JUDITSKI, A. (1992). Stochastic approximation with averaging of trajectories. *Stochastics Stochastics Rep.* **39** 107–118.
- DOUKHAN, P. (1994). *Mixing: Properties and Examples. Lecture Notes in Statist.* Springer, Berlin.
- DIEBOLT, J. and CELEUX, G. (1996). Asymptotic properties of a stochastic EM algorithm for estimating mixture proportions. *Stochastic Models* **9** 599–613.
- DIEBOLT, J. and IP, E. H. S. (1996). A stochastic EM algorithm for approximating the maximum likelihood estimate. In *Markov Chain Monte Carlo in Practice* (W. R. Gilks, S. T. Richardson, D. J. Spiegelhalter, eds.). Chapman and Hall, London.
- DUFLO, M. (1997). *Random Iterative Models*. Springer, Berlin.
- DEMPSTER, A., LAIRD, N. and RUBIN, D. (1977). Maximum-likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. Ser. B* **39** 1–38.
- FORT, J. C. and PAGÈS, G. (1996). Convergence of stochastic algorithms: from Kushner–Clark theorem to the Lyapounov functional method. *Adv. in Appl. Probab.* **28** 1072–1094.
- GEYER, C. J. and THOMPSON, E. A. (1992). Constrained Monte-Carlo maximum likelihood for dependent data. *J. Roy. Statist. Soc. Ser. B* **54** 657–699.
- GEYER, C. J. (1994). On the convergence of Monte-Carlo maximum likelihoods calculations. *J. Roy. Statist. Soc. Ser. B* **56** 261–274.
- GEYER, C. J. (1996). Likelihood inference for spatial point processes. In *Current Trends in Stochastic Geometry and its Applications* (W. S. Kendall, O. E. Barndorff-Nielsen and M. C. van Lieshout, eds.) Chapman and Hall, London. To appear.
- HALL, P. and HEYDE, C. C. (1980). *Limit Theory and Its Applications*. Academic Press, New York.
- HORN, R. and JOHNSON, C. (1985). *Matrix Analysis*. Cambridge Univ. Press.
- IBRAGIMOV, I. and HAS'MINSKI, R. (1981). *Statistical Estimation: Asymptotic Theory*. Springer, New York.
- KUSHNER, H. and CLARK, D. (1978). *Stochastic Approximation for Constrained and Unconstrained Systems*. Springer, New York.

- KUSHNER, H. and YIN, G. (1997). *Stochastic Approximation Algorithms and Applications*. Springer, Berlin.
- LANGE, K. (1995). A gradient algorithm locally equivalent to the EM algorithm. *J. Roy. Statist. Soc. Ser. B* **75** 425–437.
- LIU, C. and RUBIN, D. (1994). The ECME algorithm: a simple extension of EM and ECM with faster monotone convergence. *Biometrika* **81** 633–648.
- LITTLE, R. J. and RUBIN, D. B. (1987). *Statistical Analysis with Missing Data*. Wiley, New York.
- LOUIS, T. A. (1982). Finding the observed information matrix when using the EM algorithm. *J. Roy. Statist. Soc. Ser. B* **44** 226–233.
- MACDONALD, I. L. and ZUCCHINI, W. (1997). *Hidden Markov and Other Models for Discrete-valued Time-series*. Chapman and Hall, London.
- MEILIJSON, I. (1989). A fast improvement to the EM algorithm on its own terms. *J. Roy. Statist. Soc. Ser. B* **51** 127–138.
- MENG, X. and RUBIN, D. (1993). Maximum likelihood via the ECM algorithm: a general framework. *Biometrika* **80** 267–278.
- MENG, X. (1994). On the rate of convergence of the ECM algorithm. *Ann. Statist.* **22** 326–339.
- MURRAY, G. (1977). Discussion on P. Dempster et al., Maximum-likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc.* 27–28.
- POLYAK, B. T. (1990). New stochastic approximation type procedures. *Automatica i Telemekh.* 98–107. (English translation in *Automat. Remote Control* **51**.)
- POLYAK, B. T. and JUDITSKI, A. B. (1992). Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.* **30** 838–855.
- RAO, C. (1965). *Linear Statistical Inference and Its Applications*. Wiley, New York.
- SHEPHARD, N. (1994). Partial non-Gaussian state space. *Biometrika* **81** 115–131.
- TANNER, M. (1993). *Tools for Statistical Inference: Methods for Exploration of Posterior Distributions and Likelihood Functions*. Springer, Berlin.
- TITTERINGTON, D. M., SMITH, A. F. M. and MAKOV, U. E. (1985). *Statistical Analysis of Finite Mixture Distribution*. Wiley, New York.
- WEI, G. and TANNER, M. (1990). A Monte-Carlo implementation of the EM algorithm and the Poor's Man's data augmentation algorithm. *J. Amer. Statist. Assoc.* **85** 699–704.
- WU, C. (1983). On the convergence property of the EM algorithm. *Ann. Statist.* **11** 95–103.
- YOUNES, L. (1989). Parameter estimation for imperfectly observed Gibbsian fields. *Probab. Theory and Related Fields* **82** 625–645.
- YOUNES, L. (1992). Parameter estimation for imperfectly observed Gibbs fields and some comments on Chalmond's EM Gibbsian algorithm. *Stochastic Models, Statistical Methods and Algorithms in Image Analysis. Lecture Notes in Statistics* **74**. Springer, Berlin.

M. LAVIELLE
 UFR DE MATHÉMATIQUES ET INFORMATIQUE
 UNIVERSITÉ PARIS V
 AND
 CNRS URA 743
 UNIVERSITÉ PARIS-SUD
 E-MAIL: lavielle@stats.matups.fr

B. DELYON
 IRISA/INRIA
 CAMPUS UNIVERSITAIRE DE BEAULIEU
 35042 RENNES CEDEX
 FRANCE
 E-MAIL: delyon@irisa.fr

E. MOULINES
 ECOLE NATIONALE SUPÉRIEURE DES TÉLÉCOMMUNICATIONS
 CNRS-URA 820
 46, RUE BARRAULT
 75634 PARIS CEDEX 13
 FRANCE
 E-MAIL: moulines@sig.enst.fr