

---

# Fast Two-Timescale Stochastic EM Algorithms

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The Expectation-Maximization (EM) algorithm is a popular choice for learning  
2       latent variable models. Variants of the EM have been initially introduced by [28],  
3       using incremental updates to scale to large datasets, and by [33, 12], using Monte  
4       Carlo (MC) approximations to bypass the intractable conditional expectation of  
5       the latent data for most nonconvex models. In this paper, we propose a general  
6       class of methods called Two-Timescale EM Methods based on a two-stage ap-  
7       proach of stochastic updates to tackle an essential nonconvex optimization task  
8       for latent variable models. We motivate the choice of a double dynamic by invoking  
9       the variance reduction virtue of each stage of the method on both sources of  
10      noise: the index sampling for the incremental update and the MC approximation.  
11      We establish finite-time and global convergence bounds for nonconvex objective  
12      functions. Numerical applications are also presented to illustrate our findings.

## 1 Introduction

14   Learning latent variable models is critical for modern machine learning problems, see (e.g.,) [26]  
15   for references. We formulate the training of such model as an empirical risk minimization problem:

$$\min_{\theta \in \Theta} \bar{L}(\theta) := L(\theta) + r(\theta) \quad \text{with} \quad L(\theta) = \frac{1}{n} \sum_{i=1}^n L_i(\theta) := \frac{1}{n} \sum_{i=1}^n \{ -\log g(y_i; \theta) \}, \quad (1)$$

16   where  $\{y_i\}_{i=1}^n$  are observations,  $\Theta \subset \mathbb{R}^d$  is the parameters set and  $r : \Theta \rightarrow \mathbb{R}$  is a smooth regular-  
17   izer. The objective function  $\bar{L}(\theta)$  is possibly *nonconvex* and is assumed to be lower bounded. In the  
18   latent variable model, the likelihood  $g(y_i; \theta)$ , is the marginal of the complete data likelihood defined  
19   as  $f(z_i, y_i; \theta)$ , i.e.,  $g(y_i; \theta) = \int_{\mathcal{Z}} f(z_i, y_i; \theta) \mu(dz_i)$ , where  $\{z_i\}_{i=1}^n$  are the latent variables. In this  
20   paper, we assume that the complete model belongs to the curved exponential family [14]:

$$f(z_i, y_i; \theta) = h(z_i, y_i) \exp \left( \langle S(z_i, y_i) | \phi(\theta) \rangle - \psi(\theta) \right), \quad (2)$$

21   where  $\psi(\theta)$ ,  $h(z_i, y_i)$  are scalar functions,  $\phi(\theta) \in \mathbb{R}^k$  is a vector function, and  $\{S(z_i, y_i) \in \mathbb{R}^k\}_{i=1}^n$   
22   is the vector of sufficient statistics. Batch EM [13, 34], the method of reference for (1), is comprised  
23   of two steps. The **E-step** computes the conditional expectation of the sufficient statistics of (2):

$$\text{E-step: } \bar{s}(\theta) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta) \quad \text{where} \quad \bar{s}_i(\theta) = \int_{\mathcal{Z}} S(z_i, y_i) p(z_i | y_i; \theta) \mu(dz_i), \quad (3)$$

24   and the **M-step** is given by

$$\text{M-step: } \hat{\theta} = \bar{\theta}(\bar{s}(\theta)) := \arg \min_{\vartheta \in \Theta} \{ r(\vartheta) + \psi(\vartheta) - \langle \bar{s}(\theta) | \phi(\vartheta) \rangle \}. \quad (4)$$

25   Two caveats of this method are the following: (a) with the explosion of data, the first step of the EM  
26   is computationally inefficient as it requires, at each iteration, a full pass over the dataset; and (b) the  
27   complexity of modern models makes the expectation in (3) intractable. So far, and to the best of our  
28   knowledge, both challenges have been addressed separately, as detailed in the sequel.

**Prior Work:** Inspired by stochastic optimization procedures, [28] and [8] develop respectively an incremental and an online variant of the E-step in models where the expectation is computable, and were then extensively used and studied in [30, 23, 7]. Some improvements of those methods have been provided and analyzed, globally and in finite-time, in [20] where variance reduction techniques taken from the optimization literature have been efficiently applied to scale the EM algorithm to large datasets. Regarding the computation of the expectation under the posterior distribution, the Monte Carlo EM (MCEM) has been introduced in the seminal paper [33] where an MC approximation for this expectation is computed. A variant of that algorithm is the Stochastic Approximation of the EM (SAEM) in [12] leveraging the power of Robbins-Monro update [32] to ensure pointwise convergence of the vector of estimated parameters using a decreasing stepsize rather than increasing the number of MC samples. The MCEM and the SAEM have been successfully applied in mixed effects models [25, 16, 4] or to do inference for joint modeling of time to event data coming from clinical trials in [10], unsupervised clustering in [29], variational inference of graphical models in [5] among other applications. Recently, an incremental variant of the SAEM was proposed in [22] showing positive empirical results but its analysis is limited to asymptotic consideration.

**Contributions:** This paper *introduces* and *analyzes* a new class of methods which purpose is to update two proxies for the target expected quantities in a two-timescale manner. Those approximated quantities are then used to optimize the objective function (1) for modern examples and settings using the M-step of the EM algorithm. The main contributions of the paper are:

- We propose a two-timescale method based on (i) Stochastic Approximation (SA), to alleviate the problem of computing MC approximations, and on (ii) Incremental updates, to scale to large datasets. We describe in details the edges of each level of our method based on variance reduction arguments. Such class of algorithms has two advantages. First, it naturally leverages variance reduction and Robbins-Monro type of updates to tackle large-scale and highly nonlinear learning tasks. Then, it gives a simple formulation as a *scaled-gradient method* which makes the global analysis and the implementation accessible.
- We also establish global (independent of the initialization) and finite-time (true at each iteration) upper bounds on a classical sub-optimality condition in the nonconvex literature [18, 15], *i.e.*, the second order moment of the gradient of the objective function. We discuss the double dynamic of those bounds due to the two-timescale property of our algorithm update and we theoretically stress the advantages of introducing variance reduction in a *Stochastic Approximation* [32] scheme.

In Section 2 we formalize both incremental and Monte Carlo variants of the EM. Then, we introduce our two-timescale class of EM algorithms for which we derive several global statistical guarantees in Section 3 for possibly *nonconvex* functions. Section 4 is devoted to numerical illustrations. The supplementary material of this paper includes proofs of our theoretical results.

## 2 Two-Timescale Stochastic EM Algorithms

We recall and formalize in this section the different methods found in the literature that aim at solving the intractable expectation and the large-scale problem. We then provide the general framework of our method that efficiently tackles the optimization problem (1).

### 2.1 Monte Carlo Integration and Stochastic Approximation

As mentioned in the Introduction, for complex and possibly nonconvex models, the expectation under the posterior distribution defined in (3) is not tractable. In that case, the first solution involves computing a Monte Carlo integration of that latter. For all  $i \in [n]$ , where  $[n] := \{1, \dots, n\}$ , draw  $\{z_{i,m} \sim p(z_i|y_i; \theta)\}_{m=1}^M$  samples and compute the MC integration  $\tilde{s}$  of  $\bar{s}(\theta)$  defined by (3):

$$\text{MC-step : } \tilde{s} := \frac{1}{n} \sum_{i=1}^n \frac{1}{M} \sum_{m=1}^M S(z_{i,m}, y_i). \quad (5)$$

Then update the parameter  $\hat{\theta} = \bar{\theta}(\tilde{s})$ . This algorithm bypasses the intractable expectation issue but is rather computationally expensive in order to reach point wise convergence ( $M$  needs to be large). An alternative to that stochastic algorithm is to use a Robbins-Monro (RM) type of update. We

denote, at iteration  $k$ , the number of samples  $M_k$  and the following MC approximation  $\tilde{S}^{(k+1)}$ :

$$\tilde{S}^{(k+1)} := \frac{1}{n} \sum_{i=1}^n \tilde{S}_i^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \frac{1}{M_k} \sum_{m=1}^{M_k} S(z_{i,m}^{(k)}, y_i) \quad \text{where} \quad z_{i,m}^{(k)} \sim p(z_i | y_i; \theta^{(k)}). \quad (6)$$

Then, the RM update of the sufficient statistics  $\hat{s}^{(k+1)}$  reads:

$$\text{SA-step : } \hat{s}^{(k+1)} = \hat{s}^{(k)} + \gamma_{k+1} (\tilde{S}^{(k+1)} - \hat{s}^{(k)}), \quad (7)$$

where  $\{\gamma_k\}_{k \geq 1} \in (0, 1)$  is a sequence of decreasing stepsizes to ensure asymptotic convergence. This is called the Stochastic Approximation of the EM (SAEM) and has been shown to converge to a maximum likelihood of the observations under very general conditions [12]. In simple scenarios, the samples  $\{z_{i,m}\}_{m=1}^M$  are conditionally independent and identically distributed with distribution  $p(z_i, \theta)$ . Nevertheless, in most cases, since the loss function between the observed data  $y_i$  and the latent variable  $z_i$  can be nonconvex, sampling exactly from this distribution is not an option and the MC batch is sampled by Markov Chain Monte Carlo (MCMC) algorithm [27, 6]. It has been proved in [21] that (7) converges almost surely when coupled with an MCMC sampling procedure.

**Role of the stepsize  $\gamma_k$ :** The sequence of decreasing positive integers  $\{\gamma_k\}_{k \geq 1}$  controls the convergence of the algorithm. It is inefficient to start with small values for the stepsize  $\gamma_k$  and large values for the number of simulations  $M_k$ . Rather, it is recommended that one decreases  $\gamma_k$ , as in  $\gamma_k = 1/k^\alpha$ , with  $\alpha \in (0, 1)$ , and keeps a constant and small number  $M_k$  bypassing the computationally involved sampling step in (5). In practice,  $\gamma_k$  is set equal to 1 during the first few iterations to let the iterates explore the parameter space without memory and converge quickly to a neighborhood of the target estimate. The Stochastic Approximation is performed during the remaining iterations ensuring the almost sure convergence of the vector of estimates.

This Robbins-Monro type of update constitutes the *first level* of our algorithm, needed to temper the variance and noise introduced by the Monte Carlo integration. In the next section, we derive variants of this algorithm to adapt to the sheer size of data of today's applications and formalize the *second level* of our class of two-timescale EM methods.

## 2.2 Incremental and Two-Stage Stochastic EM Methods

Efficient strategies to scale to large datasets include incremental [28] and variance reduced [11, 19] methods. We will explicit a general update that covers those latter variants and that represents the *second level* of our algorithm, i.e., the incremental update of the noisy statistics  $\tilde{S}^{(k+1)}$  in (7):

$$\text{Incremental-step : } \tilde{S}^{(k+1)} = \tilde{S}^{(k)} + \rho_{k+1} (\mathcal{S}^{(k+1)} - \tilde{S}^{(k)}). \quad (8)$$

Note that  $\{\rho_k\}_{k \geq 1} \in (0, 1)$  is a sequence of stepsizes,  $\mathcal{S}^{(k)}$  is a proxy for  $\tilde{S}^{(k)}$ . If the stepsize is equal to one and the proxy  $\mathcal{S}^{(k)} = \tilde{S}^{(k)}$ , i.e., computed in a full batch manner as in (6), then we recover the SAEM algorithm. Also if  $\rho_k = 1$ ,  $\gamma_k = 1$  and  $\mathcal{S}^{(k)} = \tilde{S}^{(k)}$ , then we recover the MCEM [33]. For all methods, we define a random index drawn at iteration  $k$ , noted  $i_k \in [n]$ , and  $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$  as the iteration index where  $i \in [n]$  is last drawn prior to iteration  $k$ . The proposed fitTEM method draws two indices *independently* and uniformly as  $i_k, j_k \in [n]$ . Thus, we define  $t_j^k = \{k' : j_{k'} = j, k' < k\}$  to be the iteration index where the sample  $j \in [n]$  is last drawn as  $j_k$  prior to iteration  $k$  in addition to  $\tau_i^k$  which was defined w.r.t.  $i_k$ . Recall  $\tilde{S}_{i_k}^{(k)} = \frac{1}{M_k} \sum_{m=1}^{M_k} S(z_{i_k,m}^{(k)}, y_{i_k})$

**Table 1** Proxies for the Incremental-step (8)

1: iSAEM	$\mathcal{S}^{(k+1)} = \mathcal{S}^{(k)} + n^{-1} (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(\tau_{i_k}^k)})$
2: vrTTEM	$\mathcal{S}^{(k+1)} = \tilde{S}^{(\ell(k))} + (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(\ell(k))})$
3: fitTEM	$\mathcal{S}^{(k+1)} = \overline{\mathcal{S}}^{(k)} + (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(t_{i_k}^k)})$
	$\overline{\mathcal{S}}^{(k+1)} = \overline{\mathcal{S}}^{(k)} + n^{-1} (\tilde{S}_{j_k}^{(k)} - \tilde{S}_{j_k}^{(t_{j_k}^k)})$

where  $z_{i_k,m}^{(k)}$  are samples drawn from  $p(z_{i_k} | y_{i_k}; \theta^{(k)})$ . The stepsize in (8) is set to  $\rho_{k+1} = 1$  for the iSAEM method and we initialize with  $\mathcal{S}^{(0)} = \tilde{S}^{(0)}$ ;  $\rho_{k+1} = \rho$  is constant for the vrTTEM and fitTEM methods. Note that we initialize as follows  $\overline{\mathcal{S}}^{(0)} = \tilde{S}^{(0)}$  for the fitTEM which can be seen as a slightly modified version of SAGA inspired by [31]. For vrTTEM we set an epoch size of  $m$  and we define  $\ell(k) := m \lfloor k/m \rfloor$  as the first iteration number in the epoch that iteration  $k$  is in.

122 **Two-Timescale Stochastic EM methods:** We now introduce the general method derived using the  
 123 two variance reduction techniques described above. Algorithm 1 leverages both levels (7) and (8) in  
 124 order to output a vector of fitted parameters  $\hat{\theta}^{(K_m)}$  where  $K_m$  is the total number of iterations.

---

**Algorithm 1** Two-Timescale Stochastic EM methods.

---

- 1: **Input:**  $\hat{\theta}^{(0)} \leftarrow 0, \hat{s}^{(0)} \leftarrow \tilde{S}^{(0)}, \{\gamma_k\}_{k>0}, \{\rho_k\}_{k>0}$  and  $K_m \in \mathbb{N}^*$ .
- 2: **for**  $k = 0, 1, 2, \dots, K_m - 1$  **do**
- 3:   Draw index  $i_k \in [n]$  uniformly (and  $j_k \in [n]$  for fitTEM).
- 4:   Compute  $\tilde{S}_{i_k}^{(k)}$  using the MC-step (5), for the drawn indices.
- 5:   Compute the surrogate sufficient statistics  $\mathcal{S}^{(k+1)}$  using Lines 1, 2 or 3 in Table 1.
- 6:   Compute  $\tilde{S}^{(k+1)}$  and  $\hat{s}^{(k+1)}$  using respectively (8) and (7):

$$\begin{aligned}\tilde{S}^{(k+1)} &= \tilde{S}^{(k)} + \rho_{k+1}(\mathcal{S}^{(k+1)} - \tilde{S}^{(k)}) \\ \hat{s}^{(k+1)} &= \hat{s}^{(k)} + \gamma_{k+1}(\tilde{S}^{(k+1)} - \hat{s}^{(k)})\end{aligned}\tag{9}$$

- 7:   Compute  $\hat{\theta}^{(k+1)} = \bar{\theta}(\hat{s}^{(k+1)})$  via the M-step.
  - 8: **end for**
- 

125 The update in (9) is said to have a two-timescale property as the stepsizes satisfy  $\lim_{k \rightarrow \infty} \gamma_k / \rho_k < 1$   
 126 such that  $\tilde{S}^{(k+1)}$  is updated at a faster time-scale, determined by  $\rho_{k+1}$ , than  $\hat{s}^{(k+1)}$ , determined by  
 127  $\gamma_{k+1}$ . The next section introduces the main results of this paper and establishes global and finite-  
 128 time bounds for the three different updates of our scheme.

### 129 3 Finite Time Analysis of the Two-Timescale Scheme

130 Following [8], it can be shown that stationary points of the objective function (1) corresponds to the  
 131 stationary points of the following *nonconvex* Lyapunov function:

$$\min_{\mathbf{s} \in \mathcal{S}} V(\mathbf{s}) := \bar{L}(\bar{\theta}(\mathbf{s})) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\bar{\theta}(\mathbf{s})) + \mathbf{r}(\bar{\theta}(\mathbf{s})), \tag{10}$$

132 that we propose to study in this article.

#### 133 3.1 Assumptions and Intermediate Lemmas

134 Several important assumptions required to derive convergence guarantees read as follows:

135 **A1.** *The sets  $\mathcal{Z}, \mathcal{S}$  are compact. There exist constants  $C_S, C_Z$  such that:*

$$C_S := \max_{\mathbf{s}, \mathbf{s}' \in \mathcal{S}} \|\mathbf{s} - \mathbf{s}'\| < \infty, \quad C_Z := \max_{i \in [n]} \int_{\mathcal{Z}} |S(z, y_i)| \mu(dz) < \infty. \tag{11}$$

136 **A2.** *For any  $i \in [n]$ ,  $z \in \mathcal{Z}$ ,  $\theta, \theta' \in \text{int}(\Theta)^2$ , we have  $|p(z|y_i; \theta) - p(z|y_i; \theta')| \leq L_p \|\theta - \theta'\|$   
 137 where  $\text{int}(\Theta)$  denotes the interior of  $\Theta$ .*

138 We also recall that we consider curved exponential family models assuming the following:

139 **A3.** *For any  $\mathbf{s} \in \mathcal{S}$ , the function  $\theta \mapsto L(\mathbf{s}, \theta) := \mathbf{r}(\theta) + \psi(\theta) - \langle \mathbf{s} | \phi(\theta) \rangle$  admits a unique global  
 140 minimum  $\bar{\theta}(\mathbf{s}) \in \text{int}(\Theta)$ . In addition,  $J_{\phi}^{\theta}(\bar{\theta}(\mathbf{s}))$  is full rank,  $L_p$ -Lipschitz and  $\bar{\theta}(\mathbf{s})$  is  $L_t$ -Lipschitz.*

141 We denote by  $H_L^{\theta}(\mathbf{s}, \theta)$  the Hessian (w.r.t to  $\theta$  for a given value of  $\mathbf{s}$ ) of the function  $\theta \mapsto L(\mathbf{s}, \theta) =$   
 142  $\mathbf{r}(\theta) + \psi(\theta) - \langle \mathbf{s} | \phi(\theta) \rangle$ , and define  $B(\mathbf{s}) := J_{\phi}^{\theta}(\bar{\theta}(\mathbf{s})) \left( H_L^{\theta}(\mathbf{s}, \bar{\theta}(\mathbf{s})) \right)^{-1} J_{\phi}^{\theta}(\bar{\theta}(\mathbf{s}))^{\top}$ .

143 **A4.** *It holds that  $v_{\max} := \sup_{\mathbf{s} \in \mathcal{S}} \|B(\mathbf{s})\| < \infty$  and  $0 < v_{\min} := \inf_{\mathbf{s} \in \mathcal{S}} \lambda_{\min}(B(\mathbf{s}))$ . There exists  
 144 a constant  $L_b$  such that for all  $\mathbf{s}, \mathbf{s}' \in \mathcal{S}^2$ , we have  $\|B(\mathbf{s}) - B(\mathbf{s}')\| \leq L_b \|\mathbf{s} - \mathbf{s}'\|$ .*

145 The class of algorithms we develop in this paper is composed of two levels where the second stage  
 146 corresponds to the variance reduction trick used in [20] in order to accelerate incremental methods  
 147 and reduce the variance introduced by the index sampling. The first stage is the Robbins-Monro  
 148 update that aims at reducing the Monte Carlo noise of  $\tilde{S}^{(k+1)}$  at iteration  $k$  denoted as follows:

$$\eta_i^{(k)} := \tilde{S}_i^{(k)} - \bar{s}_i(\vartheta^{(k)}) \quad \text{for all } i \in [n] \quad \text{and } k > 0. \tag{12}$$

For instance, we consider that the MC approximation is unbiased if for all  $i \in [n]$  and  $m \in [M]$ , the samples  $z_{i,m} \sim p(z_i|y_i; \theta)$  are i.i.d. under the posterior distribution, i.e.,  $\mathbb{E}[\eta_i^{(k)}|\mathcal{F}_k] = 0$  where  $\mathcal{F}_k$  is the filtration up to iteration  $k$ . The following results are derived under the assumption that the fluctuations implied by the approximation are bounded:

**A5.** For all  $k > 0$ ,  $i \in [n]$ , it holds:  $\mathbb{E}[\|\eta_i^{(k)}\|^2] \leq \infty$  and  $\mathbb{E}[\|\mathbb{E}[\eta_i^{(k)}|\mathcal{F}_k]\|^2] \leq \infty$ .

Note that typically, the controls exhibited above are vanishing when the number of MC samples  $M_k$  increases with  $k$ . We now state two important results on the Lyapunov function; its smoothness:

**Lemma 1.** [20] Assume A1-A4. For all  $\mathbf{s}, \mathbf{s}' \in \mathcal{S}$  and  $i \in [n]$ , we have

$$\|\bar{\mathbf{s}}_i(\bar{\boldsymbol{\theta}}(\mathbf{s})) - \bar{\mathbf{s}}_i(\bar{\boldsymbol{\theta}}(\mathbf{s}'))\| \leq L_s \|\mathbf{s} - \mathbf{s}'\|, \quad \|\nabla V(\mathbf{s}) - \nabla V(\mathbf{s}')\| \leq L_V \|\mathbf{s} - \mathbf{s}'\|, \quad (13)$$

where  $L_s := C_Z L_p L_t$  and  $L_V := v_{\max}(1 + L_s) + L_b C_S$ .

We also establish a growth condition on the gradient of  $V$  related to the mean field of the algorithm:

**Lemma 2.** Assume A3 and A4. For all  $\mathbf{s} \in \mathcal{S}$ ,

$$v_{\min}^{-1} \langle \nabla V(\mathbf{s}) | \mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s})) \rangle \geq \|\mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s}))\|^2 \geq v_{\max}^{-2} \|\nabla V(\mathbf{s})\|^2. \quad (14)$$

We present in the following sections a finite-time and global (independent of the initialization) analysis of both the incremental and two-timescale variants our method.

### 3.2 Global Convergence of Incremental Stochastic EM Algorithms

The following result for the iSAEM algorithm is derived under the control of the Monte Carlo fluctuations as described by Assumption A5 and is built upon an intermediary Lemma, characterizing the quantity of interest  $(\hat{S}^{(k+1)} - \hat{\mathbf{s}}^{(k)})$ :

**Lemma 3.** Assume A1. The iSAEM update (1) is equivalent to the following update on the statistics  $\hat{\mathbf{s}}^{(k+1)} = \hat{\mathbf{s}}^{(k)} + \gamma_{k+1} (\sum_{i=1}^n \tilde{S}_i^{(\tau_i^k)} - \hat{\mathbf{s}}^{(k)})$ . Also:

$$\mathbb{E}[\tilde{S}^{(k+1)} - \hat{\mathbf{s}}^{(k)}] = \mathbb{E}[\bar{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(k)}] + \left(1 - \frac{1}{n}\right) \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \tilde{S}_i^{(\tau_i^k)} - \bar{\mathbf{s}}^{(k)} \right] + \frac{1}{n} \mathbb{E}[\eta_{i_k}^{(k+1)}],$$

where  $\bar{\mathbf{s}}^{(k)}$  is defined by (3) and  $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$ .

Then, the following non-asymptotic convergence rate can be derived for the iSAEM algorithm:

**Theorem 1.** Assume A1-A5. Consider the iSAEM sequence  $\{\hat{\mathbf{s}}^{(k)}\}_{k \geq 0} \in \mathcal{S}$  obtained with  $\rho_{k+1} = 1$  for any  $k \leq K_m$  where  $K_m$  is a positive integer. Let  $\{\gamma_k = 1/(k^a \alpha c_1 \bar{L})\}_{k \geq 0}$ , where  $a \in (0, 1)$ , be a sequence of stepsizes,  $c_1 = v_{\min}^{-1}$ ,  $\alpha = \max\{8, 1 + 6v_{\min}\}$ ,  $\bar{L} = \max\{L_s, L_V\}$ ,  $\beta = c_1 \bar{L}/n$ . Then:

$$v_{\max}^{-2} \sum_{k=0}^{K_m} \tilde{\alpha}_k \mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(k)})\|^2] \leq \mathbb{E}[V(\hat{\mathbf{s}}^{(0)}) - V(\hat{\mathbf{s}}^{(K_m)})] + \sum_{k=0}^{K_m-1} \tilde{\Gamma}_k \mathbb{E}[\|\eta_{i_k}^{(k)}\|^2].$$

Note that, in Theorem 1, the convergence bound is composed of an initialization term  $V(\hat{\mathbf{s}}^{(0)}) - V(\hat{\mathbf{s}}^{(K_m)})$  and suffers from the Monte Carlo noise introduced by the posterior sampling step, see the second term on the RHS of the inequality. We observe, in the next section, that when variance reduction is applied ( $\rho_k < 1$ ), a second phase of convergence will be included in our bounds.

### 3.3 Global Convergence of Two-Timescale Stochastic EM Algorithms

We now deal with the analysis of Algorithm 1 when variance reduction is applied i.e.,  $\rho < 1$ . Two important intermediate Lemmas are needed in order to establish finite-time bounds for the vrTTEM and the fitTEM methods. We first derive an identity for the drift term of the vrTTEM :

**Lemma 4.** Consider the vrTTEM update (2) with  $\rho_k = \rho$ , it holds for all  $k > 0$

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)}\|^2] &\leq 2\rho^2 \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)}\|^2] + 2\rho^2 L_s^2 \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(\ell(k))}\|^2] \\ &\quad + 2(1 - \rho)^2 \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2] + 2\rho^2 \mathbb{E}[\|\eta_{i_k}^{(k+1)}\|^2], \end{aligned}$$

where we recall that  $\ell(k)$  is the first iteration number in the epoch that iteration  $k$  is in.

183 The second one derives an identity for the quantity  $\mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)}\|^2]$  using the fitTEM update:

184 **Lemma 5.** Consider the fitTEM update (3) with  $\rho_k = \rho$ . It holds for all  $k > 0$  that

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)}\|^2] &\leq 2\rho^2 \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)}\|^2] + 2\rho^2 \frac{L_s^2}{n} \sum_{i=1}^n \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2] \\ &\quad + 2(1 - \rho)^2 \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2] + 2\rho^2 \mathbb{E}[\|\eta_{i_k}^{(k+1)}\|^2], \end{aligned}$$

185 where  $L_s$  is the smoothness constant defined in Lemma 1.

186 Let  $K$  be an independent discrete r.v. drawn from  $\{1, \dots, K_m\}$  with distribution  $\{\gamma_{k+1}/P_m\}_{k=0}^{K_m-1}$ ,  
187 then, for any  $K_m > 0$ , the convergence criterion used in our study reads

$$\mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(K)})\|^2] = \frac{1}{P_m} \sum_{k=0}^{K_m-1} \gamma_{k+1} \mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(k)})\|^2],$$

188 where  $P_m = \sum_{\ell=0}^{K_m-1} \gamma_\ell$  and the expectation is over the stochasticity of the algorithm. Denote  
189  $\Delta V = V(\hat{\mathbf{s}}^{(0)}) - V(\hat{\mathbf{s}}^{(K_m)})$ . We now state the main result regarding the vrTTEM method:

190 **Theorem 2.** Assume A1-A5. Consider the vrTTEM sequence  $\{\hat{\mathbf{s}}^{(k)}\}_{k>0} \in \mathcal{S}$  for any  $k \leq K_m$  where  
191  $K_m$  is a positive integer. Let  $\{\gamma_{k+1} = 1/(k^a \bar{L})\}_{k>0}$ , where  $a \in (0, 1)$ , be a sequence of stepsizes,  
192  $\bar{L} = \max\{L_s, L_V\}$ ,  $\rho = \mu/(c_1 \bar{L} n^{2/3})$ ,  $m = nc_1^2/(2\mu^2 + \mu c_1^2)$  and a constant  $\mu \in (0, 1)$ . Then:

$$\mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(K)})\|^2] \leq \frac{2n^{2/3} \bar{L}}{\mu P_m v_{\min}^2 v_{\max}^2} \left( \mathbb{E}[\Delta V] + \sum_{k=0}^{K_m-1} \tilde{\eta}^{(k+1)} + \chi^{(k+1)} \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2] \right).$$

193 Furthermore, the fitTEM method has the following convergence rate:

194 **Theorem 3.** Assume A1-A5. Consider the fitTEM sequence  $\{\hat{\mathbf{s}}^{(k)}\}_{k>0} \in \mathcal{S}$  for any  $k \leq K_m$  where  
195  $K_m$  be a positive integer. Let  $\{\gamma_{k+1} = 1/(k^a \alpha c_1 \bar{L})\}_{k>0}$ , where  $a \in (0, 1)$ , be a sequence of  
196 positive stepsizes,  $\alpha = \max\{2, 1 + 2v_{\min}\}$ ,  $\bar{L} = \max\{L_s, L_V\}$ ,  $\beta = 1/(\alpha n)$ ,  $\rho = 1/(\alpha c_1 \bar{L} n^{2/3})$   
197 and  $c_1(k\alpha - 1) \geq c_1(\alpha - 1) \geq 2$ ,  $\alpha \geq 2$ . Then:

$$\mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(K)})\|^2] \leq \frac{4\alpha \bar{L} n^{2/3}}{P_m v_{\min}^2 v_{\max}^2} \left( \mathbb{E}[\Delta V] + \sum_{k=0}^{K_m-1} \Xi^{(k+1)} + \Gamma^{(k+1)} \mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2] \right).$$

198 Note that in those two bounds, the quantities  $\tilde{\eta}^{(k+1)}$  and  $\Xi^{(k+1)}$  depend only on the Monte Carlo  
199 noises  $\mathbb{E}[\|\eta_{i_k}^{(k)}\|^2]$ ,  $\mathbb{E}[\|\mathbb{E}[\eta_i^{(r)} | \mathcal{F}_r]\|^2]$ , bounded under Assumption A5, and some constants.

200 **Remarks:** Theorem 2 and Theorem 3 exhibit in their convergence bounds *two different phases*. The  
201 upper bounds display a *bias term* due to the initial conditions, i.e., the term  $\Delta V$ , and a *double*  
202 *dynamic* burden exemplified by the term  $\mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2]$ . Indeed, the following remarks are  
203 worth doing on this quantity: (i) This term is the price we pay for the two-timescale dynamic and  
204 corresponds to the gap between the two *asynchronous* updates (one on  $\hat{\mathbf{s}}^{(k)}$  and the other on  $\tilde{S}^{(k)}$ ).  
205 (ii) It is readily understood that if  $\rho = 1$ , i.e., there is no variance reduction, then for any  $k > 0$

$$\mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2] = \mathbb{E}[\|\mathbf{S}^{(k+1)} - \tilde{S}^{(k+1)}\|^2] = 0 \quad \text{with} \quad \hat{\mathbf{s}}^{(0)} = \tilde{S}^{(0)} = 0,$$

206 which strengthen the fact that this quantity characterizes the impact of the variance reduction tech-  
207 nique introduced in our class of methods. The following Lemma characterizes this gap:

208 **Lemma 6.** Considering a decreasing stepsize  $\gamma_k \in (0, 1)$  and a constant  $\rho \in (0, 1)$ , we have

$$\mathbb{E}[\|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k)}\|^2] \leq \frac{\rho}{1 - \rho} \sum_{\ell=0}^k (1 - \gamma_\ell)^2 (\mathbf{S}^{(\ell)} - \tilde{S}^{(\ell)}),$$

209 where  $\mathbf{S}^{(k)}$  is defined either by Line 2 (vrTTEM) or Line 3 (fitTEM).



## 4 Numerical Examples

This section presents several numerical applications for our proposed class of Algorithms 1.

### 4.1 Gaussian Mixture Models

We begin by a simple and illustrative example. The authors acknowledge that the following model can be trained using deterministic EM-type of algorithms but propose to apply stochastic methods, including theirs, in order to compare their performances. Given  $n$  observations  $\{y_i\}_{i=1}^n$ , we want to fit a Gaussian Mixture Model (GMM) whose distribution is modeled as a mixture of  $M$  Gaussian components, each with a unit variance. Let  $z_i \in [M]$  be the latent labels of each component, the complete log-likelihood is defined as follows:

$$\log f(z_i, y_i; \theta) = \sum_{m=1}^M \mathbb{1}_{\{m\}}(z_i) [\log(\omega_m) - \mu_m^2/2] + \sum_{m=1}^M \mathbb{1}_{\{m\}}(z_i) \mu_m y_i + \text{constant}.$$

where  $\theta := (\omega, \mu)$  with  $\omega = \{\omega_m\}_{m=1}^{M-1}$  are the mixing weights with the convention  $\omega_M = 1 - \sum_{m=1}^{M-1} \omega_m$  and  $\mu = \{\mu_m\}_{m=1}^M$  are the means. We use the penalization  $r(\theta) = \frac{\delta}{2} \sum_{m=1}^M \mu_m^2 - \log \text{Dir}(\omega; M, \epsilon)$  where  $\delta > 0$  and  $\text{Dir}(\cdot; M, \epsilon)$  is the  $M$  dimensional symmetric Dirichlet distribution with concentration parameter  $\epsilon > 0$ . The constraint set is given by  $\Theta = \{\omega_m, m = 1, \dots, M-1 : \omega_m \geq 0, \sum_{m=1}^{M-1} \omega_m \leq 1\} \times \{\mu_m \in \mathbb{R}, m = 1, \dots, M\}$ . In the following experiments on synthetic data, we generate 50 synthetic datasets of size  $n = 10^5$  from a GMM model with  $M = 2$  components of means  $\mu_1 = -\mu_2 = 0.5$ . We run the EM method until convergence (to double precision) to obtain the ML estimate  $\mu^*$  averaged on 50 datasets. We compare the EM, iEM (incremental EM), SAEM, iSAEM, vrTTEM and fitTTEM methods in terms of their precision measured by  $|\mu - \mu^*|^2$ . We set the stepsize of the SA-step for all method as  $\gamma_k = 1/k^\alpha$  with  $\alpha = 0.5$ , and the stepsize  $\rho_k$  for the vrTTEM and the fitTTEM to a constant stepsize equal to  $1/n^{2/3}$ . The number of MC samples is fixed to  $M = 10$ . Figure 1 shows the precision  $|\mu - \mu^*|^2$  for the different methods through the epoch(s) (one epoch equals  $n$  iterations). The vrTTEM and fitTTEM methods outperform the other stochastic methods, supporting the benefits of our scheme.

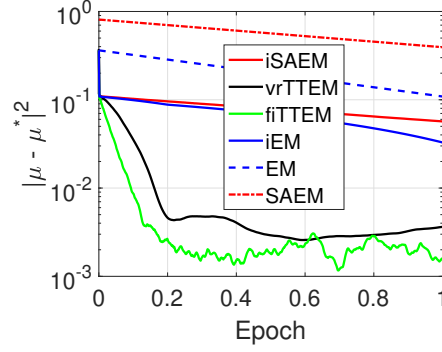


Figure 1: Precision  $|\mu^{(k)} - \mu^*|^2$  per epoch

### 4.2 Deformable Template Model for Image Analysis

Let  $(y_i, i \in [n])$  be observed gray level images defined on a grid of pixels. Let  $u \in \mathcal{U} \subset \mathbb{R}^2$  denote the pixel index on the image and  $x_u \in \mathcal{D} \subset \mathbb{R}^2$  its location. The model used in this experiment suggests that each image  $y_i$  is a deformation of a template, noted  $I : \mathcal{D} \rightarrow \mathbb{R}$ , common to all images of the dataset:

$$y_i(u) = I(x_u - \Phi_i(x_u, z_i)) + \varepsilon_i(u) \quad (15)$$

where  $\phi_i : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is a deformation function,  $z_i$  some latent variable parameterizing this deformation and  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  is an observation error. The template model, given  $\{p_k\}_{k=1}^{k_p}$  landmarks on the template, a fixed known kernel  $\mathbf{K}_p$  and a vector of parameters  $\beta \in \mathbb{R}^{k_p}$  is defined as follows:

$$I_\xi = \mathbf{K}_p \beta, \quad \text{where} \quad (\mathbf{K}_p \beta)(x) = \sum_{k=1}^{k_p} \mathbf{K}_p(x, p_k) \beta_k.$$

Given a set of landmarks  $\{g_k\}_{k=1}^{k_g}$  and a fixed kernel  $\mathbf{K}_g$ , we parameterize the deformation  $\Phi_i$  as:

$$\Phi_i = \mathbf{K}_g z_i \quad \text{where} \quad (\mathbf{K}_g z_i)(x) = \sum_{k=1}^{k_g} \mathbf{K}_g(x, g_k) \left( z_i^{(1)}(k), z_i^{(2)}(k) \right),$$

where we put a Gaussian prior on the latent variables,  $z_i \sim \mathcal{N}(0, \Gamma)$  and  $z_i \in (\mathbb{R}^{k_g})^2$ . The vector of parameters we estimate is thus  $\theta = (\beta, \Gamma, \sigma)$ . The complete model (15) belongs to the curved exponential family, see [2], which vector of sufficient statistics for all  $i \in [n]$  is defined by  $S(y_i, z_i) = (\mathbf{K}_{p, z_i}^\top y_i, \mathbf{K}_{p, z_i}^\top \mathbf{K}_{p, z_i}, z_i^\top z_i)$  where we denote  $\mathbf{K}_{p, z_i} = \mathbf{K}_{p, z_i}(x_u - \phi_i(x_u, z_i), p_j)$ . Then, the two-timescale M-step (4) yields the following parameter updates

253  $\bar{\theta}(\hat{s}) = (\beta(\hat{s}) = \hat{s}_2^{-1}(z)\hat{s}_1(z), \Gamma(\hat{s}) = \hat{s}_3(z)/n, \sigma(\hat{s}) = \beta(\hat{s})^\top \hat{s}_2(z)\beta(\hat{s}) - 2\beta(\hat{s})\hat{s}_1(z))$  where  
 254  $\hat{s} = (\hat{s}_1(z), \hat{s}_2(z), \hat{s}_3(z))$  is the vector of statistics obtained via update (9) in Algorithm 1.

255 **Numerical Experiment:** We apply model (15) and our Algorithm 1 to a collection of handwritten  
 256 digits, called the US postal database [17], featuring  $n = 1000$ ,  $(16 \times 16)$ -pixel images for each  
 257 class of digits from 0 to 9. The main challenge with this dataset stems from the geometric dispersion  
 258 within each class of digit as shown Figure 2 for digit 5. We thus ought to use our deformable  
 259 template model (15) in order to account for both sources of variability: the intrinsic template to each  
 260 class of digit and the small and local deformations in each observed image.



Figure 2: Training set of the USPS database (20 images for digit 5)

261 Figure 3 shows the resulting synthetic images for digit 5 through several epochs, for the batch  
 262 method, the online SAEM, the incremental SAEM and the various two-timescale methods. For  
 263 all methods, the initialization of the template (16) is the mean of the gray level images. In our  
 264 experiments, we have chosen Gaussian kernels for both,  $K_p$  and  $K_g$ , defined on  $\mathbb{R}^2$  and centered  
 265 on the landmark points  $\{p_k\}_{k=1}^{k_p}$  and  $\{g_k\}_{k=1}^{k_g}$  with standard respective standard deviations of 0.12  
 266 and 0.3. We set  $k_p = 15$  and  $k_g = 6$  equidistributed landmarks points on the grid for the training  
 267 procedure. Those hyperparameters are inspired by relevant studies [1, 3]. In particular, the choice  
 268 of the geometric covariance, indexed by  $g$ , in such study is critical since it has a direct impact on  
 269 the *sharpness* of the templates. As for the photometric hyperparameter, indexed by  $p$ , both the  
 270 template and the geometry are impacted, in the sense that with a large photometric variance, the  
 271 kernel centered on one landmark *spreads out* to many of its neighbors.

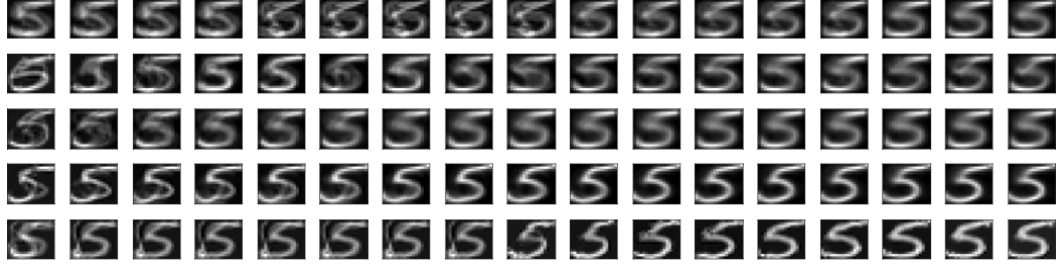


Figure 3: (USPS Digits) Estimation of the template. From top to bottom: batch, online, iSAEM, vrT-TEM and fitTEM through 7 epochs. Note that Batch method templates are replicated in-between epochs for a fair comparison with incremental variants.

272 As the iterations proceed, the templates become sharper. Figure 3 displays the virtue of the vrTTEM  
 273 and fitTEM methods that obtain a more *contrasted* and *accurate* template estimate. The incremental  
 274 and online version are looking much better on the very first epochs compared to the batch method,  
 275 which is intuitive given the high computational cost of the latter. After a few epochs, the batch  
 276 SAEM estimates similar template as the incremental an online methods due to their high variance.  
 277 Our variance reduced and fast incremental variants are effective in the long run and sharpen the final  
 278 template estimates contrasting between the background and the regions of interest in the image.

## 279 5 Conclusion

280 This paper introduces a new class of two-timescale EM methods for learning latent variable models.  
 281 In particular, the models dealt with in this paper belong to the curved exponential family and are  
 282 possibly nonconvex. The nonconvexity of the problem is tackled using a Robbins-Monro type of  
 283 update, which represents the *first level* of our class of methods. The scalability with the number  
 284 of samples is performed through a variance reduced and incremental update, the *second* and last  
 285 level of our newly introduced scheme. The various algorithms are interpreted as scaled gradient  
 286 methods, in the space of the sufficient statistics, and our convergence results are *global*, in the sense  
 287 of independence of the initial values, and *non-asymptotic*, *i.e.*, true for any random termination  
 288 number. Numerical examples illustrate the benefits of our scheme on synthetic and real tasks.



## 6 Broader Impact

Our work aims at improving training procedures for latent data models. Latent data models are particularly interesting in several impactful domains such as sociology, economy or pharmacology. Indeed, in those latter domains, a special instance of latent data models, namely mixed-effect models, can be employed. It considers a latent structure in order to take into account the variability between subjects, which can be individuals, companies or patients. In that case, our class of algorithms becomes useful as demonstrated in the additional experiment in the supplementary material. It is worth noting that other types of latent variable models could be used for impactful research, such as the missing data framework when the considered observations are sensible, yet missing.

## References

- [1] Stéphanie Allasonnière and Estelle Kuhn. Stochastic algorithm for parameter estimation for dense deformable template mixture model. *arXiv preprint arXiv:0802.1521*, 2008.
- [2] Stéphanie Allasonnière, Yali Amit, and Alain Trounev. Towards a coherent statistical framework for dense deformable template estimation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(1):3–29, 2007.
- [3] Stéphanie Allasonnière, Estelle Kuhn, Alain Trounev, et al. Construction of bayesian deformable models via a stochastic approximation algorithm: a convergence study. *Bernoulli*, 16(3):641–678, 2010.
- [4] Charlotte Baey, Samis Trevezas, and Paul-Henry Cournède. A non linear mixed effects model of plant growth and estimation via stochastic variants of the em algorithm. *Communications in Statistics-Theory and Methods*, 45(6):1643–1669, 2016.
- [5] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational Inference: A Review for Statisticians. *Journal of the American statistical Association*, 112(518):859–877, JUN 2017. ISSN 0162-1459. doi: {10.1080/01621459.2017.1285773}.
- [6] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of markov chain monte carlo*. CRC press, 2011.
- [7] Olivier Cappé. Online EM algorithm for hidden markov models. *Journal of Computational and Graphical Statistics*, 20(3):728–749, 2011.
- [8] Olivier Cappé and Eric Moulines. On-line expectation–maximization algorithm for latent data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3): 593–613, 2009.
- [9] Bradley P Carlin and Siddhartha Chib. Bayesian model choice via markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(3):473–484, 1995.
- [10] Arindom Chakraborty and Kalyan Das. Inferences for joint modelling of repeated ordinal scores and time to event data. *Computational and mathematical methods in medicine*, 11(3): 281–295, 2010.
- [11] Jianfei Chen, Jun Zhu, Yee Whye Teh, and Tong Zhang. Stochastic expectation maximization with variance reduction. In *Advances in Neural Information Processing Systems*, pages 7978–7988, 2018.
- [12] Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the em algorithm. *Ann. Statist.*, 27(1):94–128, 03 1999. doi: 10.1214/aos/1018031103. URL <https://doi.org/10.1214/aos/1018031103>.

- 332 [13] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete  
333 data via the EM algorithm. *Journal of the royal statistical society. Series B (methodological)*,  
334 pages 1–38, 1977.
- 335 [14] Bradley Efron et al. Defining the curvature of a statistical problem (with applications to second  
336 order efficiency). *The Annals of Statistics*, 3(6):1189–1242, 1975.
- 337 [15] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex  
338 stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 339 [16] James P Hughes. Mixed effects models with censored data with application to hiv rna levels.  
340 *Biometrics*, 55(2):625–629, 1999.
- 341 [17] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on*  
342 *pattern analysis and machine intelligence*, 16(5):550–554, 1994.
- 343 [18] Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *arXiv*  
344 *preprint arXiv:1712.07897*, 2017.
- 345 [19] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive vari-  
346 ance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- 347 [20] Belhal Karimi, Hoi-To Wai, Éric Moulines, and Marc Lavielle. On the global convergence  
348 of (fast) incremental expectation maximization methods. In *Advances in Neural Information*  
349 *Processing Systems*, pages 2833–2843, 2019.
- 350 [21] Estelle Kuhn and Marc Lavielle. Coupling a stochastic approximation version of em with an  
351 mcmc procedure. *ESAIM: Probability and Statistics*, 8:115–131, 2004.
- 352 [22] Estelle Kuhn, Catherine Matias, and Tabea Rebafka. Properties of the stochastic approximation  
353 em algorithm with mini-batch sampling. *arXiv preprint arXiv:1907.09164*, 2019.
- 354 [23] Percy Liang and Dan Klein. Online em for unsupervised models. In *Proceedings of human*  
355 *language technologies: The 2009 annual conference of the North American chapter of the*  
356 *association for computational linguistics*, pages 611–619, 2009.
- 357 [24] Florian Maire, Eric Moulines, and Sidonie Lefebvre. Online em for functional data, 2016.  
358 URL <http://arxiv.org/abs/1604.00570>. cite arxiv:1604.00570v1.pdf.
- 359 [25] Charles E McCulloch. Maximum likelihood algorithms for generalized linear mixed models.  
360 *Journal of the American statistical Association*, 92(437):162–170, 1997.
- 361 [26] Geoffrey McLachlan and Thriyambakam Krishnan. *The EM algorithm and extensions*, volume  
362 382. John Wiley & Sons, 2007.
- 363 [27] Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science  
364 & Business Media, 2012.
- 365 [28] Radford M Neal and Geoffrey E Hinton. A view of the EM algorithm that justifies incremental,  
366 sparse, and other variants. In *Learning in graphical models*, pages 355–368. Springer, 1998.
- 367 [29] SK Ng and GJ McLachlan. On the choice of the number of blocks with the incremental EM  
368 algorithm for the fitting of normal mixtures. *Statistics and Computing*, 13(1):45–55, FEB  
369 2003. ISSN 0960-3174. doi: {10.1023/A:1021987710829}.
- 370 [30] Hien D Nguyen, Florence Forbes, and Geoffrey J McLachlan. Mini-batch learning of expo-  
371 nential family finite mixture models. *Statistics and Computing*, pages 1–18, 2020.

- 372 [31] Sashank J Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola. Fast incremental method for  
373 nonconvex optimization. *arXiv preprint arXiv:1603.06159*, 2016.
- 374 [32] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of math-*  
375 *ematical statistics*, pages 400–407, 1951.
- 376 [33] Greg CG Wei and Martin A Tanner. A monte carlo implementation of the em algorithm and  
377 the poor man’s data augmentation algorithms. *Journal of the American statistical Association*,  
378 85(411):699–704, 1990.
- 379 [34] CF Jeff Wu. On the convergence properties of the em algorithm. *The Annals of statistics*,  
380 pages 95–103, 1983.