

---

# Fast Two-Time-Scale Noisy EM Algorithms

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 T.B.C

## 2 1 Introduction

3 We formulate the following empirical risk minimization as:

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

4 where  $\{y_i\}_{i=1}^n$  are the observations,  $\Theta$  is a convex subset of  $\mathbb{R}^d$  for the parameters,  $R : \Theta \rightarrow \mathbb{R}$  is a  
5 smooth convex regularization function and for each  $\theta \in \Theta$ ,  $g(y; \theta)$  is the (incomplete) likelihood of  
6 each individual observation. The objective function  $\bar{\mathcal{L}}(\theta)$  is possibly *non-convex* and is assumed to  
7 be lower bounded  $\bar{\mathcal{L}}(\theta) > -\infty$  for all  $\theta \in \Theta$ .

8 In the latent variable model,  $g(y_i; \theta)$ , is the marginal of the complete data likelihood defined as  
9  $f(z_i, y_i; \theta)$ , 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 (unobserved) latent vari-  
10 ables. We make the assumption of a complete model belonging to the curved exponential family,  
11 i.e.,

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

12 where  $\psi(\theta)$ ,  $h(z_i, y_i)$  are scalar functions,  $\phi(\theta) \in \mathbb{R}^k$  is a vector function, and  $S(z_i, y_i) \in \mathbb{R}^k$  is  
13 the complete data sufficient statistics.

14 **Prior Work** Cite Kuhn [Kuhn et al., 2019] (for ISAEM) and incremental EM like papers. As well  
15 as Optim papers (Variance reduction, SAGA etc.)

## 16 2 Expectation Maximization Algorithm

17 Full batch EM is a two steps procedure. The E-step amounts to computing the conditional expecta-  
18 tion of the complete data sufficient statistics,

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

19 The M-step is given by

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

### 20 3 Monte Carlo Integration and Stochastic Approximation

21 For complex and possibly nonlinear models, the expectation under the posterior distribution defined  
 22 in (3) is not tractable. In that case, the first solution involves computing a Monte Carlo integration  
 23 of that latter term. For all  $i \in \llbracket 1, n \rrbracket$ , draw for  $m \in \llbracket 1, M \rrbracket$ , samples  $z_{i,m} \sim p(z_i|y_i; \theta)$  and compute  
 24 the MC integration  $\tilde{s}$  of the deterministic quantity  $\bar{s}(\theta)$ :

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

25 and compute  $\hat{\theta} = \bar{\theta}(\hat{s})$ .

26 This algorithm bypasses the intractable expectation issue but is rather computationally expensive in  
 27 order to reach point wise convergence ( $M$  needs to be large).

28 As a result, an alternative to that stochastic algorithm is to use a Robbins-Monro (RM) type of  
 29 update. We denote

$$\tilde{S}^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \tilde{S}_i^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \frac{1}{M} \sum_{m=1}^M S(z_{i,m}^{(k)}, y_i) \quad (6)$$

30 where  $z_{i,m}^{(k)} \sim p(z_i|y_i; \theta^{(k)})$ . At iteration  $k$ , the sufficient statistics  $\hat{s}^{(k+1)}$  is approximated as follows:

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

31 where  $\{\gamma_k\}_{k=1}^\infty \in [0, 1]$  is a sequence of decreasing step sizes to ensure asymptotic convergence.  
 32 This is called the Stochastic Approximation of the EM (SAEM), see [Delyon et al., 1999] and allows  
 33 a smooth convergence to the target parameter. It represents the *first level* of our algorithm (needed  
 34 to temper the variance and noise implied by MC integration).

35 In the next section, we derive variants of this algorithm to adapt of the sheer size of data of today's  
 36 applications.

### 37 4 Incremental and Bi-Level Inexact EM Methods

38 Strategies to scale to large datasets include classical incremental and variance reduced variants. We  
 39 will explicit a general update that will cover those variants and that represents the *second level* of our  
 40 algorithm, namely the incremental update of the noisy statistics  $\hat{S}^{(k)}$  inside the RM type of update.

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

41 Note  $\{\rho_k\}_{k=1}^\infty \in [0, 1]$  is a sequence of step sizes,  $\mathcal{S}^{(k)}$  is a proxy for  $\tilde{S}^{(k)}$ , If the stepsize is equal  
 42 to one and the proxy  $\mathcal{S}^{(k)} = \hat{S}^{(k)}$ , i.e., computed in a full batch manner as in (6), then we recover  
 43 the SAEM algorithm. Also if  $\rho_k = 1$ ,  $\gamma_k = 1$  and  $\mathcal{S}^{(k)} = \tilde{S}^{(k)}$ , then we recover the Monte Carlo  
 44 EM algorithm.

45 We now introduce three variants of the SAEM update depending on different definitions of the proxy  
 46  $\mathcal{S}^{(k)}$  and the choice of the stepsize  $\rho_k$ . Let  $i_k \in \llbracket 1, n \rrbracket$  be a random index drawn at iteration  $k$  and  
 47  $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$  be the iteration index where  $i \in \llbracket 1, n \rrbracket$  is last drawn prior to  
 48 iteration  $k$ . For iteration  $k \geq 0$ , the fiSAEM method draws *two* indices *independently* and uniformly  
 49 as  $i_k, j_k \in \llbracket 1, n \rrbracket$ . In addition to  $\tau_i^k$  which was defined *w.r.t.*  $i_k$ , we define  $t_j^k = \{k' : j_{k'} = j, k' <$   
 50  $k\}$  to be the iteration index where the sample  $j \in \llbracket 1, n \rrbracket$  is last drawn as  $j_k$  prior to iteration  $k$ . With  
 51 the initialization  $\bar{\mathcal{S}}^{(0)} = \bar{s}^{(0)}$ , we use a slightly different update rule from SAGA inspired by [Reddi

et al., 2016]. Then, we obtain:

$$(iSAEM [Karimi, 2019, Kuhn et al., 2019]) \quad \mathcal{S}^{(k+1)} = \mathcal{S}^{(k)} + \frac{1}{n} (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(\tau_{i_k}^k)}) \quad (9)$$

$$(vrSAEM This paper) \quad \mathcal{S}^{(k+1)} = \tilde{S}^{(\ell(k))} + (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(\ell(k))}) \quad (10)$$

$$(fiSAEM This paper) \quad \mathcal{S}^{(k+1)} = \bar{\mathcal{S}}^{(k)} + (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(t_{i_k}^k)}) \quad (11)$$

$$\bar{\mathcal{S}}^{(k+1)} = \bar{\mathcal{S}}^{(k)} + n^{-1} (\tilde{S}_{j_k}^{(k)} - \tilde{S}_{j_k}^{(t_{j_k}^k)}). \quad (12)$$

The stepsize is set to  $\rho_{k+1} = 1$  for the iSAEM method;  $\rho_{k+1} = \gamma$  is constant for the vrSAEM and fiSAEM methods. Moreover, for iSAEM we initialize with  $\mathcal{S}^{(0)} = \tilde{S}^{(0)}$ ; for vrSAEM we set an epoch size of  $m$  and define  $\ell(k) := m \lfloor k/m \rfloor$  as the first iteration number in the epoch that iteration  $k$  is in.

---

**Algorithm 1** Two-Time-Scale Noisy EM methods.

---

- 1: **Input:** initializations  $\hat{\theta}^{(0)} \leftarrow 0, \hat{s}^{(0)} \leftarrow \hat{S}^{(0)}, K_{\max} \leftarrow \text{max. iteration number.}$
- 2: Set the terminating iteration number,  $K \in \{0, \dots, K_{\max} - 1\}$ , as a discrete r.v. with:

$$P(K = k) = \frac{\gamma_k}{\sum_{\ell=0}^{K_{\max}-1} \gamma_{\ell}}. \quad (13)$$

- 3: **for**  $k = 0, 1, 2, \dots, K$  **do**
- 4:   Draw index  $i_k \in \llbracket 1, n \rrbracket$  uniformly (and  $j_k \in \llbracket 1, n \rrbracket$  for fiSAEM).
- 5:   Compute  $\hat{S}_i^{(k)}$  using the MC-step (5), for the drawn indices.
- 6:   Compute the surrogate sufficient statistics  $\mathcal{S}^{(k+1)}$  using (9) or (10) or (11).
- 7:   Compute  $\tilde{S}^{(k+1)}$  and  $\hat{s}^{(k+1)}$  using respectively (8) and (7):

$$\begin{aligned} \tilde{S}^{(k+1)} &= \tilde{S}^{(k)} + \rho_{k+1} (\mathcal{S}^{(k+1)} - \tilde{S}^{(k)}) \\ \hat{s}^{(k+1)} &= \hat{s}^{(k)} + \gamma_{k+1} (\tilde{S}^{(k+1)} - \hat{s}^{(k)}) \end{aligned} \quad (14)$$

- 8:   Compute  $\hat{\theta}^{(k+1)}$  via the M-step (4).
  - 9: **end for**
  - 10: **Return:**  $\hat{\theta}^{(K)}$ .
- 

The updates in (14) is said to have two timescales as the step sizes satisfy  $\lim_{k \rightarrow \infty} \gamma_k / \rho_k < 1$  such that  $\tilde{S}^{(k+1)}$  is updated at a faster timescale than  $\hat{s}^{(k+1)}$ .

## 5 Finite Time Analysis

First, we consider the following minimization problem on the statistics space:

$$\min_{\mathbf{s} \in \mathcal{S}} V(\mathbf{s}) := \bar{\mathcal{L}}(\bar{\theta}(\mathbf{s})) = R(\bar{\theta}(\mathbf{s})) + \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\bar{\theta}(\mathbf{s})) \quad (15)$$

It has been shown that this minimization problem is equivalent to the optimization problem (1), see [Karimi et al., 2019, Lemma2]

**H1.**  $\Theta$  is an open set of  $\mathbb{R}^d$  and the sets  $Z, S$  are measurable open sets such that:

$$S \supset \left\{ n^{-1} \sum_{i=1}^n u_i, u_i \in \text{conv}(\bar{\mathbf{s}}_i(\theta)) \right\} \quad (16)$$

where  $\bar{\mathbf{s}}_i(\theta)$  is defined in (3).

**H2.** The conditional distribution is smooth on  $\text{int}(\Theta)$ . For any  $i \in \llbracket 1, n \rrbracket$ ,  $z \in Z$ ,  $\theta, \theta' \in \text{int}(\Theta)^2$ , we have  $|p(z|y_i; \theta) - p(z|y_i; \theta')| \leq L_p \|\theta - \theta'\|$ .

We also recall from the introduction that we consider curved exponential family models. besides:

68 **H3.** For any  $s \in S$ , the function  $\theta \mapsto L(s, \theta) := R(\theta) + \psi(\theta) - \langle s | \phi(\theta) \rangle$  admits a unique global  
69 minimum  $\bar{\theta}(s) \in \text{int}(\Theta)$ . In addition,  $J_\phi^\theta(\bar{\theta}(s))$  is full rank and  $\bar{\theta}(s)$  is  $L_\theta$ -Lipschitz.

70 Similar to [Karimi et al., 2019], we denote by  $H_L^\theta(s, \theta)$  the Hessian (w.r.t to  $\theta$  for a given value of  
71  $s$ ) of the function  $\theta \mapsto L(s, \theta) = R(\theta) + \psi(\theta) - \langle s | \phi(\theta) \rangle$ , and define

$$B(s) := J_\phi^\theta(\bar{\theta}(s)) \left( H_L^\theta(s, \bar{\theta}(s)) \right)^{-1} J_\phi^\theta(\bar{\theta}(s))^\top. \quad (17)$$

72 **H4.** It holds that  $v_{\max} := \sup_{s \in S} \|B(s)\| < \infty$  and  $0 < v_{\min} := \inf_{s \in S} \lambda_{\min}(B(s))$ . There exists  
73 a constant  $L_B$  such that for all  $s, s' \in S^2$ , we have  $\|B(s) - B(s')\| \leq L_B \|s - s'\|$ .

74 We now formulate the main difference with the work done in [Karimi et al., 2019]. The class of  
75 algorithms we develop in this paper are two time-scale where the first stage corresponds to the  
76 variance reduction trick used in [Karimi et al., 2019] in order to accelerate incremental methods and  
77 kill the variance induced by the index sampling. The second stage is the Robbins-Monro type of  
78 update that aims to kill the variance induced by the MC approximations

79 Indeed the expectations (3) are never available and requires Monte Carlo approximation. Thus, at  
80 iteration  $k + 1$ , we introduce the errors when approximating the quantity  $\bar{s}_i(\hat{\theta}(\hat{s}^{(k-1)}))$ . For all  
81  $i \in \llbracket 1, n \rrbracket$ ,  $r > 0$  and  $\vartheta \in \Theta$ , define:

$$\eta_{i,\vartheta}^{(r)} := \tilde{S}_i^{(r)} - \bar{s}_i(\vartheta) \quad (18)$$

82 For instance, we consider that the MC approximation is unbiased if for all  $i \in \llbracket 1, n \rrbracket$  and  $m \in$   
83  $\llbracket 1, M \rrbracket$ , the samples  $z_{i,m} \sim p(z_i | y_i; \theta)$  are i.i.d. under the posterior distribution, i.e.,  $\mathbb{E}[\eta_{i,\vartheta}^{(r)} | \mathcal{F}_r] = 0$   
84 where  $\mathcal{F}_r$  is the filtration up to iteration  $r$ .

85 The following results are derived under the assumption of control of the fluctuations implied by the  
86 approximation stated as follows:

87 **H5.** There exist a positive sequence of MC batch size  $\{M_k\}_{k>0}$  and constants  $(C, C_\eta)$  such that for  
88 all  $k > 0$ ,  $i \in \llbracket 1, n \rrbracket$  and  $\vartheta \in \Theta$ :

$$\mathbb{E} \left[ \left\| \eta_{i,\vartheta}^{(r)} \right\|^2 \right] \leq \frac{C_\eta}{M_r} \quad \text{and} \quad \mathbb{E} \left[ \left\| \mathbb{E}[\eta_{i,\vartheta}^{(r)} | \mathcal{F}_r] \right\|^2 \right] \leq \frac{C}{M_r} \quad (19)$$

89 **Lemma 1.** [Karimi et al., 2019] Assume H2, H3, H4. For all  $s, s' \in S$  and  $i \in \llbracket 1, n \rrbracket$ , we have

$$\|\bar{s}_i(\bar{\theta}(s)) - \bar{s}_i(\bar{\theta}(s'))\| \leq L_s \|s - s'\|, \quad \|\nabla V(s) - \nabla V(s')\| \leq L_V \|s - s'\|, \quad (20)$$

90 where  $L_s := C_Z L_p L_\theta$  and  $L_V := v_{\max}(1 + L_s) + L_B C_S$ .

## 91 5.1 Global Convergence of Incremental Noisy EM Algorithms

92 Following the asymptotic analysis of update (9), we present a finite-time analysis of the incremental  
93 variant of the Stochastic Approximation of the EM algorithm.

94 The first intermediate result is the computation of the quantity  $\hat{S}^{(k+1)} - \hat{s}^{(k)}$ , which corresponds to  
95 the dirft term of (7) and reads as follows:

96 **Lemma 2.** Assume H1. The update (9) is equivalent to the following update on the resulting statis-  
97 tics

$$\hat{s}^{(k+1)} = \hat{s}^{(k)} + \gamma_{k+1} \left( n^{-1} \sum_{i=1}^n \tilde{S}_i^{(\tau_i^k)} - \hat{s}^{(k)} \right) \quad (21)$$

98 where  $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$ . Also:

$$\mathbb{E} \left[ \left\| \hat{S}^{(k+1)} - \hat{s}^{(k)} \right\|^2 \right] \leq \mathbb{E} \left[ \left\| \bar{s}^{(k)} - \hat{s}^{(k)} \right\|^2 \right] + 2L_s^2 \left( 1 - \frac{1}{n} \right)^2 n^{-1} \sum_{i=1}^n \mathbb{E} \left[ \left\| \hat{s}^{(k)} - \hat{s}^{(\tau_i^k)} \right\|^2 \right] + \frac{2C}{M_k} \quad (22)$$

99 where  $\bar{s}^{(k)}$  is defined by (3).

100 The following main result for the iSAEM algorithm is derived under a control of the Monte Carlo  
 101 fluctuations as described by assumption H 5. Typically, the controls exhibited below are of interest  
 102 when the number of MC samples  $M_k$  increase with the iteration index  $f$ .

103 **Theorem 1.** *Let  $K_{\max}$  be a positive integer. Let  $\{\gamma_k, k \in \mathbb{N}\}$  be a sequence of positive step sizes  
 104 and consider the iSAEM sequence  $\{\hat{\mathbf{s}}^{(k)}, k \in \mathbb{N}\}$  obtained with  $\rho_{k+1} = 1$  for any  $k$ .*

105 *Assume that  $\hat{\mathbf{s}}^{(k)} \in \mathcal{S}$  for any  $k \leq K_{\max}$ .*

106 **Proof** Under some regularity conditions of the Lyapunov function  $V$ , cf. Lemma 20, and the fol-  
 107 lowing growth condition for all  $\mathbf{s} \in \mathcal{S}$ ,

$$v_{\min}^{-1} \langle \nabla V(\mathbf{s}) | \mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s})) \rangle \geq \|\mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s}))\|^2 \geq v_{\max}^{-2} \|\nabla V(\mathbf{s})\|^2, \quad (23)$$

108 proven in [Karimi et al., 2019, Lemma 3], we can write:

$$V(\hat{\mathbf{s}}^{(k+1)}) \leq V(\hat{\mathbf{s}}^{(k)}) - \gamma_{k+1} \langle \hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle + \frac{\gamma_{k+1}^2 L_V}{2} \|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)}\|^2 \quad (24)$$

109 Taking the expectation on both sides and using the growth condition (23), we obtain:

$$\begin{aligned} \mathbb{E}[V(\hat{\mathbf{s}}^{(k+1)})] &\leq \mathbb{E}[V(\hat{\mathbf{s}}^{(k)})] - \gamma_{k+1} v_{\min} \mathbb{E} \left[ \left\| \bar{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(k)} \right\|^2 \right] + \mathbb{E} \left[ \frac{\gamma_{k+1}^2 L_V}{2} \|\hat{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)}\|^2 \right] \\ &\quad - \gamma_{k+1} \mathbb{E} \left[ \langle \bar{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right] \end{aligned} \quad (25)$$

110 We then establish an auxiliary Lemma yielding an upper-bound on the quantity  
 111  $\mathbb{E} \left[ \langle \bar{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right]$  where:

$$\bar{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)} = \bar{\mathbf{s}}^{(k)} - \left( \tilde{S}^{(k)} + \frac{1}{n} (\tilde{S}_{i_k}^{(k)} - \tilde{S}_{i_k}^{(\tau_{i_k}^k)}) \right) \quad (26)$$

112

**Lemma 3.**

$$\mathbb{E} \left[ \langle \bar{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right] \leq \quad (27)$$

113 Using Lemma 2:

$$\begin{aligned} \mathbb{E}[V(\hat{\mathbf{s}}^{(k+1)})] &\leq \mathbb{E}[V(\hat{\mathbf{s}}^{(k)})] - \gamma_{k+1} \left( v_{\min} - \frac{\gamma_{k+1} L_V}{2} \right) \mathbb{E} \left[ \left\| \bar{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(k)} \right\|^2 \right] \\ &\quad + \gamma_{k+1}^2 L_V L_{\mathbf{s}}^2 \left( 1 - \frac{1}{n} \right)^2 n^{-1} \sum_{i=1}^n \mathbb{E} \left[ \left\| \hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(\tau_i^k)} \right\|^2 \right] + \frac{\gamma_{k+1}^2 L_V C}{M_k} \\ &\quad - \gamma_{k+1} \mathbb{E} \left[ \langle \bar{\mathbf{s}}^{(k)} - \tilde{S}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right] \end{aligned} \quad (28)$$

114 Besides,

$$n^{-1} \sum_{i=1}^n \mathbb{E} \left[ \left\| \hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(\tau_i^{k+1})} \right\|^2 \right] = n^{-1} \sum_{i=1}^n \left( \frac{1}{n} \mathbb{E} [\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2] + \frac{n-1}{n} \mathbb{E} [\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(\tau_i^k)}\|^2] \right) \quad (29)$$

115 yielding for any numbers  $\beta_k > 0$ ,

$$\begin{aligned} &\mathbb{E} [\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2] \\ &= \mathbb{E} \left[ \|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2 + \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 + 2 \langle \hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)} | \hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)} \rangle \right] \\ &= \mathbb{E} \left[ \|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2 + \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 - 2 \gamma_{k+1} \langle \hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)} | \hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)} \rangle \right] \\ &\leq \mathbb{E} \left[ \|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2 + \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 + \frac{\gamma_{k+1}}{\beta_{k+1}} \|\hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)}\|^2 + \gamma_{k+1} \beta_{k+1} \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 \right] \end{aligned} \quad (30)$$

116

□

## 117 **5.2 Global Convergence of Two-Time-Scale Noisy EM Algorithms**

118 We now proceed by giving our main result regarding the global convergence of the fiSAEM algo-  
119 rithm.

## 120 **6 Numerical Examples**

### 121 **6.1 Gaussian Mixture Models**

122 Graphs obtained and relevant

### 123 **6.2 Deep Latent Variable Models using noisy EM**

124 See if makes sense to use EM instead of Variational Inference

### 125 **6.3 Deformable Template Model for Image Analysis**

126 See Kuhn et.al. paper.

## 127 **7 Conclusion**

## References

- B. Delyon, M. Lavielle, and E. 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>.
- B. Karimi. *Non-Convex Optimization for Latent Data Models: Algorithms, Analysis and Applications*. PhD thesis, 2019.
- B. Karimi, H.-T. Wai, É. Moulines, and M. Lavielle. On the global convergence of (fast) incremental expectation maximization methods. In *Advances in Neural Information Processing Systems*, pages 2833–2843, 2019.
- E. Kuhn, C. Matias, and T. Rebafka. Properties of the stochastic approximation em algorithm with mini-batch sampling. *arXiv preprint arXiv:1907.09164*, 2019.
- S. J. Reddi, S. Sra, B. Póczos, and A. Smola. Fast incremental method for nonconvex optimization. *arXiv preprint arXiv:1603.06159*, 2016.