

---

# Fast Bi-Level and Incremental 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  $\theta \in \Theta$ . In the latent variable model,  $g(y_i; \theta)$ , is the marginal  
8 of the complete data likelihood defined 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  
9  $\{z_i\}_{i=1}^n$  are the (unobserved) latent variables. We make the assumption of a complete model be-  
10 longing to the curved exponential family, *i.e.*,

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

11 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  
12 the complete data sufficient statistics.

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

## 15 2 Expectation Maximization Algorithm

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

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

18 The **M-step** is given by

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

### 19 3 Monte Carlo Integration and Stochastic Approximation

For complex and possibly nonlinear 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 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 the MC integration  $\hat{s}$  of the deterministic quantity  $\bar{s}(\theta)$ :

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

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

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

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

$$\hat{S}^{(k)} = \frac{1}{n} \sum_{i=1}^n \hat{S}_i^{(k)} = \frac{1}{n} \sum_{i=1}^n \frac{1}{M} \sum_{m=1}^M S(z_{i,m}^{(k)}, y_i) \quad (5)$$

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

$$\text{SA-step : } \hat{s}^{(k)} = \hat{s}^{(k-1)} + \gamma_k (\hat{S}^{(k)} - \hat{s}^{(k-1)}) \quad (6)$$

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

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

### 32 4 Incremental and Bi-Level Inexact EM Methods

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

$$\text{Inexact-step : } \hat{S}^{(k)} = \hat{S}^{(k-1)} + \rho_{k+1} (\mathcal{S}^{(k)} - \hat{S}^{(k-1)}), \quad (7)$$

36 Note  $\{\rho_k\}_{k=1}^\infty \in [0, 1]$  is a sequence of step sizes,  $\mathcal{S}^{(k+1)}$  is a proxy for  $\hat{S}^{(k)}$ . If the stepsize is equal  
37 to one and the proxy  $\mathcal{S}^{(k+1)} = \hat{S}^{(k)}$ , i.e., computed in a full batch manner as in (5), then we recover  
38 the SAEM algorithm. Also if  $\rho_k = 1$ ,  $\gamma_k = 1$  and  $\mathcal{S}^{(k+1)} = \hat{S}^{(k)}$ , then we recover the Monte Carlo  
39 EM algorithm.

40 We now introduce three variants of the SAEM update depending on different definitions of the proxy  
41  $\mathcal{S}^{(k)}$  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  
42  $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$  be the iteration index where  $i \in \llbracket 1, n \rrbracket$  is last drawn prior to  
43 iteration  $k$ . For iteration  $k \geq 0$ , the fiSAEM method draws *two* indices *independently* and uniformly  
44 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  $t_j^k = \{k' : j_{k'} = j, k' <$   
45  $k\}$  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  
46 the initialization  $\bar{\mathcal{S}}^{(0)} = \bar{s}^{(0)}$ , we use a slightly different update rule from SAGA inspired by [Reddi

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

$$(iSAEM [Kuhn et al., 2019, Karimi, 2019]) \quad \mathcal{S}^{(k+1)} = \mathcal{S}^{(k)} + \frac{1}{n} (\hat{S}_{i_k}^{(k)} - \hat{S}_{i_k}^{(\tau_{i_k}^k)}) \quad (8)$$

$$(vrSAEM This paper) \quad \mathcal{S}^{(k+1)} = \hat{S}^{(\ell(k))} + (\hat{S}_{i_k}^{(k)} - \hat{S}_{i_k}^{(\ell(k))}) \quad (9)$$

$$(fiSAEM This paper) \quad \mathcal{S}^{(k+1)} = \overline{\mathcal{S}}^{(k)} + (\hat{S}_{i_k}^{(k)} - \hat{S}_{i_k}^{(t_{i_k}^k)}) \quad (10)$$

$$\overline{\mathcal{S}}^{(k+1)} = \overline{\mathcal{S}}^{(k)} + n^{-1} (\hat{S}_{j_k}^{(k)} - \hat{S}_{j_k}^{(t_{j_k}^k)}). \quad (11)$$

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

---

**Algorithm 1** Bi-Level Stochastic Approximation 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:

$$P(K = k) = \frac{\gamma_k}{\sum_{\ell=0}^{K_{\max}-1} \gamma_\ell}. \quad (12)$$

- 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 the surrogate sufficient statistics  $\mathcal{S}^{(k+1)}$  using (8) or (9) or (10).
  - 6:   Compute  $\hat{S}^{(k+1)}$  via the Inexact-step (7).
  - 7:   Compute  $\hat{s}^{(k+1)}$  via the SA-step (6).
  - 8:   Compute  $\hat{\theta}^{(k+1)}$  via the M-step (4).
  - 9: **end for**
  - 10: **Return:**  $\hat{\theta}^{(K)}$ .
- 

## 52 5 Finite Time Analysis

53 Finite analysis of iSAEM vrSAEM and fiSAEM .

54 Analysis in the curved exponential family assumption.

55 Suboptimality condition would be:  $\mathbb{E}[\|\nabla V(\hat{s}^{(K)})\|^2]$  where

$$\min_{\mathbf{s} \in \mathcal{S}} V(\mathbf{s}) := \overline{\mathcal{L}}(\overline{\boldsymbol{\theta}}(\mathbf{s})) = \mathbf{R}(\overline{\boldsymbol{\theta}}(\mathbf{s})) + \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\overline{\boldsymbol{\theta}}(\mathbf{s})) \quad (13)$$

56 is the Lyapunov function minimized here.

## 57 6 Numerical Examples

### 58 6.1 Gaussian Mixture Models

59 Graphs obtained and relevant

### 60 6.2 Logistic Regression with Missing values OR random effects

61 To Be Done

## 62 7 Conclusion

63 **References**

- 64 B. Delyon, M. Lavielle, and E. Moulines. Convergence of a stochastic approximation version of  
65 the em algorithm. *Ann. Statist.*, 27(1):94–128, 03 1999. doi: 10.1214/aos/1018031103. URL  
66 <https://doi.org/10.1214/aos/1018031103>.
- 67 B. Karimi. *Non-Convex Optimization for Latent Data Models: Algorithms, Analysis and Applica-*  
68 *tions*. PhD thesis, 2019.
- 69 E. Kuhn, C. Matias, and T. Rebafka. Properties of the stochastic approximation em algorithm with  
70 mini-batch sampling. *arXiv preprint arXiv:1907.09164*, 2019.
- 71 S. J. Reddi, S. Sra, B. Póczos, and A. Smola. Fast incremental method for nonconvex optimization.  
72 *arXiv preprint arXiv:1603.06159*, 2016.

