
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:

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

4 where $\{y_i\}_{i=1}^n$ are the observations, Θ 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.*,

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

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 (?) (for ISAEM) and incremental EM like papers. As well as Optim papers
14 (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,

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

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 (??) 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

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

25 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:

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

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 (?) and allows a smooth
28 convergence to the target parameter. It represents the *first level* of our algorithm (needed to temper
29 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.

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

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

47 Then, we obtain:

$$(iSAEM \text{ (??)}) \quad \mathcal{S}^{(k)} = \mathcal{S}^{(k)} + \frac{1}{n} (\hat{S}_{i_k}^{(k)} - \hat{S}_{i_k}^{(\tau_{i_k}^k)}) \quad (8)$$

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

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

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

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

52 5 Finite Time Analysis

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

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

54 It has been shown that this minimization problem is equivalent to the optimization problem (??), see
 55 (?, Lemma2)

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

$$\mathcal{S} \supset \left\{ n^{-1} \sum_{i=1}^n u_i, u_i \in \text{conv}(\bar{\mathbf{s}}_i(\theta)) \right\} \quad (14)$$

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

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

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

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

63 Similar to (?), we denote by $H_L^{\theta}(\mathbf{s}, \theta)$ the Hessian (w.r.t to θ for a given value of \mathbf{s}) of the function
 64 $\theta \mapsto L(\mathbf{s}, \theta) = R(\theta) + \psi(\theta) - \langle \mathbf{s} | \phi(\theta) \rangle$, and define

$$B(\mathbf{s}) := J_{\phi}^{\theta}(\bar{\theta}(\mathbf{s})) \left(H_L^{\theta}(\mathbf{s}, \bar{\theta}(\mathbf{s})) \right)^{-1} J_{\phi}^{\theta}(\bar{\theta}(\mathbf{s}))^{\top}. \quad (15)$$

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

67 We now formulate the main difference with the work done in (?). The class of algorithms we
 68 develop in this paper are two time-scale where the first stage corresponds to the variance reduction
 69 trick used in (?) in order to accelerate incremental methods and kill the variance induced by the
 70 index sampling. The second stage is the Robbins-Monro type of update that aims to kill the variance
 71 induced by the MC approximations

72 Indeed the expectations (??) are never available and requires Monte Carlo approximation. Thus,
 73 at iteration $k + 1$, we introduce the errors when approximating the quantity $\bar{\mathbf{s}}_i(\hat{\boldsymbol{\theta}}(\hat{\mathbf{s}}^{(k-1)}))$. For all
 74 $i \in \llbracket 1, n \rrbracket$, $r > 0$ and $\vartheta \in \Theta$, define:

$$\eta_{i,\vartheta}^{(r)} := \hat{S}_i^{(r)} - \bar{\mathbf{s}}_i(\vartheta) \quad (16)$$

75 For instance, we consider that the MC approximation is unbiased if for all $i \in \llbracket 1, n \rrbracket$ and $m \in$
 76 $\llbracket 1, M \rrbracket$, the samples $z_{i,m} \sim p(z_i | y_i; \theta)$ are i.i.d. under the posterior distribution, i.e., $\mathbb{E}[\eta_{i,\vartheta}^{(r)} | \mathcal{F}_r] = 0$
 77 where \mathcal{F}_r is the filtration up to iteration r .

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

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

$$\mathbb{E} \left[\left\| \eta_{i,\vartheta}^{(r)} \right\|^2 \right] \leq \frac{C_\eta}{M_r} \quad \text{and} \quad \mathbb{E} \left[\mathbb{E}[\eta_{i,\vartheta}^{(r)} | \mathcal{F}_r] \right] \leq \frac{C}{M_r} \quad (17)$$

82 **Lemma 1.** (?) Assume $H??$, $H??$, $H??$. For all $\mathbf{s}, \mathbf{s}' \in \mathcal{S}$ and $i \in \llbracket 1, n \rrbracket$, we have

$$\|\bar{\mathbf{s}}_i(\bar{\boldsymbol{\theta}}(\mathbf{s})) - \bar{\mathbf{s}}_i(\bar{\boldsymbol{\theta}}(\mathbf{s}'))\| \leq L_s \|\mathbf{s} - \mathbf{s}'\|, \quad \|\nabla V(\mathbf{s}) - \nabla V(\mathbf{s}')\| \leq L_V \|\mathbf{s} - \mathbf{s}'\|, \quad (18)$$

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

84 5.1 Global Convergence of Incremental Noisy EM Algorithms

85 Following the asymptotic analysis of update (??), we present a finite-time analysis of the incremental
 86 variant of the Stochastic Approximation of the EM algorithm.

87 The first intermediate result is the computation of the quantity $\hat{S}^{(k+1)} - \hat{\mathbf{s}}^{(k)}$, which corresponds to
 88 the drift term of (??) and reads as follows:

89 **Lemma 2.** Assume $H??$. The update (??) is equivalent to the following update on the resulting
 90 statistics

$$\hat{\mathbf{s}}^{(k+1)} = \hat{\mathbf{s}}^{(k)} + \gamma_{k+1} \left(n^{-1} \sum_{i=1}^n \hat{S}_i^{(\tau_i^k)} - \hat{\mathbf{s}}^{(k)} \right) \quad (19)$$

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

$$\mathbb{E} \left[\hat{S}^{(k+1)} - \hat{\mathbf{s}}^{(k)} | \mathcal{F}_k \right] = (\bar{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(k)}) + \left(1 - \frac{1}{n} \right) \left(n^{-1} \sum_{i=1}^n \hat{S}_i^{(\tau_i^k)} - \bar{\mathbf{s}}^{(k)} \right) + \frac{C}{M_k} \quad (20)$$

92 where $\bar{\mathbf{s}}^{(k)}$ is defined by (??).

93 **Theorem 1.** *ff*

94 5.2 Global Convergence of Two-Time-Scale Noisy EM Algorithms

95 6 Numerical Examples

96 6.1 Gaussian Mixture Models

97 Graphs obtained and relevant

98 **6.2 Deep Latent Variable Models using noisy EM**

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

100 **6.3 Deformable Template Model for Image Analysis**

101 See Kuhn et.al. paper.

102 **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.

