
Fast Bi-Level and Incremental Stochastic Approximation of the EM Algorithm

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.,

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

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 (for ISAEM) and incremental EM like papers .As well as Optim papers
14 (Variance reduction, SAGA etc.)

15 2 Expectation Maximization Methods

16 The basic "batch" EM (bEM) method iteratively computes a sequence of estimates $\{\theta^k, k \in \mathbb{N}\}$
17 with an initial parameter θ^0 . Each iteration of bEM is composed of two steps. In the **E-step**, a
18 surrogate function is computed as $\theta \mapsto Q(\theta, \theta^{k-1}) = \sum_{i=1}^n Q_i(\theta, \theta^{k-1})$ where $Q_i(\theta, \theta') :=$
19 $-\int_{\mathcal{Z}} \log f(z_i, y_i; \theta) p(z_i | y_i; \theta') \mu(dz_i)$ such that $p(z_i | y_i; \theta) := f(z_i, y_i; \theta) / g(y_i, \theta)$ is the condi-
20 tional probability density of the latent variables z_i given the observations y_i . When $f(z_i, y_i; \theta)$ is a
21 curved exponential family model, the **E-step** amounts to computing the conditional expectation of
22 the complete data sufficient statistics,

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

23 In the **M-step**, the surrogate function is minimized producing a new fit of the parameter $\theta^k =$
24 $\arg \max_{\theta \in \Theta} Q(\theta, \theta^{k-1})$.

25 3 Monte Carlo Integration and Stochastic Approximation

26 4 Incremental and Bi-Level Inexact EM MEthods

27 We first describe the stochastic EM methods to be analyzed under a unified framework. The k th
28 iteration of a generic stochastic EM method is composed of two sub-steps —

$$\text{sE-step : } \hat{\mathbf{s}}^{(k+1)} = \hat{\mathbf{s}}^{(k)} - \gamma_{k+1}(\hat{\mathbf{s}}^{(k)} - \mathcal{S}^{(k+1)}), \quad (4)$$

29 which is a stochastic version of the E-step in (3). Note $\{\gamma_k\}_{k=1}^\infty \in [0, 1]$ is a sequence of step sizes,
30 $\mathcal{S}^{(k+1)}$ is a proxy for $\bar{\mathbf{s}}(\hat{\boldsymbol{\theta}}^{(k)})$, and $\bar{\mathbf{s}}$ is defined in (3). The M-step is given by

$$\text{M-step: } \hat{\boldsymbol{\theta}}^{(k+1)} = \bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}^{(k+1)}) := \arg \min_{\boldsymbol{\theta} \in \Theta} \{ \mathbf{R}(\boldsymbol{\theta}) + \psi(\boldsymbol{\theta}) - \langle \hat{\mathbf{s}}^{(k+1)} | \phi(\boldsymbol{\theta}) \rangle \}, \quad (5)$$

31 which is controlled by the sufficient statistics determined by the sE-step. The stochastic EM meth-
32 ods differ in the way that $\mathcal{S}^{(k+1)}$ is computed. Existing methods employ stochastic approximation
33 or variance reduction without the need to fully compute $\bar{\mathbf{s}}(\hat{\boldsymbol{\theta}}^{(k)})$. To simplify notations, we define

$$\bar{\mathbf{s}}_i^{(k)} := \bar{\mathbf{s}}_i(\hat{\boldsymbol{\theta}}^{(k)}) = \int_{\mathcal{Z}} S(z_i, y_i) p(z_i | y_i; \hat{\boldsymbol{\theta}}^{(k)}) \mu(dz_i) \quad \text{and} \quad \bar{\mathbf{s}}^{(\ell)} := \bar{\mathbf{s}}(\hat{\boldsymbol{\theta}}^{(\ell)}) = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{s}}_i^{(\ell)}. \quad (6)$$

34 Note that if $\mathcal{S}^{(k+1)} = \bar{\mathbf{s}}^{(k)}$ and $\gamma_{k+1} = 1$, eq. (4) reduces to the E-step in the classical bEM method.
35 To describe the stochastic EM methods, let $i_k \in \llbracket 1, n \rrbracket$ be a random index drawn at iteration k and
36 $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$ be the iteration index where $i \in \llbracket 1, n \rrbracket$ is last drawn prior to
37 iteration k , we have:

$$(iEM \text{ [Neal and Hinton, 1998]}) \quad \mathcal{S}^{(k+1)} = \mathcal{S}^{(k)} + \frac{1}{n} (\bar{\mathbf{s}}_{i_k}^{(k)} - \bar{\mathbf{s}}_{i_k}^{(\tau_{i_k}^k)}) \quad (7)$$

$$(sEM \text{ [Cappé and Moulines, 2009]}) \quad \mathcal{S}^{(k+1)} = \bar{\mathbf{s}}_{i_k}^{(k)} \quad (8)$$

$$(sEM\text{-VR [Chen et al., 2018]}) \quad \mathcal{S}^{(k+1)} = \bar{\mathbf{s}}^{(\ell(k))} + (\bar{\mathbf{s}}_{i_k}^{(k)} - \bar{\mathbf{s}}_{i_k}^{(\ell(k))}) \quad (9)$$

38 The stepsize is set to $\gamma_{k+1} = 1$ for the iEM method; $\gamma_{k+1} = \gamma$ is constant for the sEM-VR method.
39 In the original version of the sEM method, the sequence of step γ_{k+1} is a diminishing step size.
40 Moreover, for iEM we initialize with $\mathcal{S}^{(0)} = \bar{\mathbf{s}}^{(0)}$; for sEM-VR, we set an epoch size of m and
41 define $\ell(k) := m \lfloor k/m \rfloor$ as the first iteration number in the epoch that iteration k is in.

42 **fiEM** Our analysis framework can handle a new, yet natural application of a popular variance
43 reduction technique to the EM method. The new method, called fiEM, is developed from the SAGA
44 method [Defazio et al., 2014] in a similar vein as in sEM-VR.

45 For iteration $k \geq 0$, the fiEM method draws *two* indices *independently* and uniformly as $i_k, j_k \in$
46 $\llbracket 1, n \rrbracket$. In addition to τ_i^k which was defined w.r.t. i_k , we define $t_j^k = \{k' : j_{k'} = j, k' < k\}$ to be
47 the iteration index where the sample $j \in \llbracket 1, n \rrbracket$ is last drawn as j_k prior to iteration k . With the
48 initialization $\bar{\mathcal{S}}^{(0)} = \bar{\mathbf{s}}^{(0)}$, we use a slightly different update rule from SAGA inspired by [Reddi
49 et al., 2016], as described by the following recursive updates

$$\mathcal{S}^{(k+1)} = \bar{\mathcal{S}}^{(k)} + (\bar{\mathbf{s}}_{i_k}^{(k)} - \bar{\mathbf{s}}_{i_k}^{(t_{i_k}^k)}), \quad \bar{\mathcal{S}}^{(k+1)} = \bar{\mathcal{S}}^{(k)} + n^{-1} (\bar{\mathbf{s}}_{j_k}^{(k)} - \bar{\mathbf{s}}_{j_k}^{(t_{j_k}^k)}). \quad (10)$$

Algorithm 1 Stochastic EM methods.

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

$$P(K = k) = \frac{\gamma_k}{\sum_{\ell=0}^{K_{\max}-1} \gamma_\ell}. \quad (11)$$

- 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 fiEM).
 - 5: Compute the surrogate sufficient statistics $\mathcal{S}^{(k+1)}$ using (8) or (7) or (9) or (10).
 - 6: Compute $\hat{\mathbf{s}}^{(k+1)}$ via the sE-step (4).
 - 7: Compute $\hat{\boldsymbol{\theta}}^{(k+1)}$ via the M-step (5).
 - 8: **end for**
 - 9: **Return:** $\hat{\boldsymbol{\theta}}^{(K)}$.
-

50 where we set a constant step size as $\gamma_{k+1} = \gamma$.

51 In the above, the update of $\mathcal{S}^{(k+1)}$ corresponds to an *unbiased estimate* of $\bar{\mathbf{s}}^{(k)}$, while the update for
52 $\bar{\mathcal{S}}^{(k+1)}$ maintains the structure that $\bar{\mathcal{S}}^{(k)} = n^{-1} \sum_{i=1}^n \bar{\mathbf{s}}_i^{(t_i^k)}$ for any $k \geq 0$. The two updates of (10)
53 are based on two different and independent indices i_k, j_k that are randomly drawn from $\llbracket n \rrbracket$. This is
54 used for our fast convergence analysis in Section 5.

55 **5 Finite Time Analysis**

56 **6 Numerical Examples**

57 **7 Conclusion**

58 **References**

- 59 O. Cappé and E. Moulines. On-line expectation–maximization algorithm for latent data models.
60 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613, 2009.
- 61 J. Chen, J. Zhu, Y. W. Teh, and T. Zhang. Stochastic expectation maximization with variance reduc-
62 tion. In *Advances in Neural Information Processing Systems*, pages 7978–7988, 2018.
- 63 A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support
64 for non-strongly convex composite objectives. In *Advances in neural information processing*
65 *systems*, pages 1646–1654, 2014.
- 66 A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the
67 EM algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38,
68 1977.
- 69 S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic pro-
70 gramming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 71 G. McLachlan and T. Krishnan. *The EM algorithm and extensions*, volume 382. John Wiley & Sons,
72 2007.
- 73 R. M. Neal and G. E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and
74 other variants. In *Learning in graphical models*, pages 355–368. Springer, 1998.
- 75 S. J. Reddi, S. Sra, B. Póczos, and A. Smola. Fast incremental method for nonconvex optimization.
76 *arXiv preprint arXiv:1603.06159*, 2016.

