

---

# On the Global Convergence of (Fast) Incremental Expectation Maximization Methods

---

**Belhal Karimi**

CMAP, École Polytechnique  
Palaiseau, France  
belhal.karimi@polytechnique.edu

**Hoi-To Wai**

The Chinese University of Hong Kong  
Shatin, Hong Kong  
htwai@se.cuhk.edu.hk

**Eric Moulines**

CMAP, École Polytechnique  
Palaiseau, France  
eric.moulines@polytechnique.edu

**Marc Lavielle**

INRIA Saclay  
Palaiseau, France  
marc.lavielle@inria.fr

## Abstract

The EM algorithm is one of the most popular algorithm for inference in latent data models. The original formulation of the EM algorithm does not scale to large data set, because the whole data set is required at each iteration of the algorithm. To alleviate this problem, [Neal and Hinton \[1998\]](#) have proposed an incremental version of the EM (iEM) in which at each iteration the conditional expectation of the latent data (E-step) is updated only for a mini-batch of observations. Another approach has been proposed by [Cappé and Moulines \[2009\]](#) in which the E-step is replaced by a stochastic approximation step, closely related to stochastic gradient. In this paper, we analyze incremental and stochastic version of the EM algorithm as well as the variance reduced-version of [\[Chen et al., 2018\]](#) in a common unifying framework. We also introduce a new version incremental version, inspired by the SAGA algorithm by [Defazio et al. \[2014\]](#). We establish non-asymptotic convergence bounds for global convergence. Numerical applications are presented in this article to illustrate our findings.

## 1 Introduction

Many problems in machine learning pertain to tackling an empirical risk minimization of the form

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

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 smooth convex regularization function and for each  $\theta \in \Theta$ ,  $g(y; \theta)$  is the (incomplete) likelihood of each individual observation. The objective function  $\bar{\mathcal{L}}(\theta)$  is possibly *non-convex* and is assumed to 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 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  $\{z_i\}_{i=1}^n$  are the (unobserved) latent variables. We consider the setting where the complete data likelihood belongs to the curved exponential family, *i.e.*,

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

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 the complete data sufficient statistics. Latent variable models are widely used in machine learning and

statistics; examples include mixture models for density estimation, clustering document, and topic modelling; see [McLachlan and Krishnan, 2007] and the references therein.

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

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

In the **M-step**, the surrogate function is minimized producing a new fit of the parameter  $\theta^k = \arg \max_{\theta \in \Theta} Q(\theta, \theta^{k-1})$ . The EM method has several appealing features – it is monotone where the likelihood do not decrease at each iteration, invariant with respect to the parameterization, numerically stable when the optimization set is well defined, etc. The EM method has been the subject of considerable interest since its formalization in [Dempster et al., 1977].

With the sheer size of data sets today, the bEM method is not applicable as the **E-step** (3) involves a full pass over the dataset of  $n$  observations. Several approaches based on stochastic optimization have been proposed to address this problem. Neal and Hinton [1998] proposed (but not analyzed) an incremental version of EM, referred to as the iEM method. Cappé and Moulines [2009] developed the online EM (sEM) method which uses a stochastic approximation procedure to track the sufficient statistics defined in (3). Recently, Chen et al. [2018] proposed a variance reduced sEM (sEM-VR) method which is inspired by the SVRG algorithm popular in stochastic convex optimization [Johnson and Zhang, 2013]. The applications of the above stochastic EM methods are numerous, especially with the iEM and sEM methods; e.g., [Thiesson et al., 2001] for inference with missing data, [Ng and McLachlan, 2003] for mixture models and unsupervised clustering, [Hinton et al., 2006] for inference of deep belief networks, [Hofmann, 1999] for probabilistic latent semantic analysis, [Wainwright et al., 2008, Blei et al., 2017] for variational inference of graphical models and [Ablin et al., 2018] for Independent Component Analysis.

This paper focuses on the theoretical aspect of stochastic EM methods by establishing novel *non-asymptotic* and *global* convergence rates for them. Our contributions are as follows.

- We offer two complementary views for the global convergence of EM methods – one focuses on the parameter space, and one on the sufficient statistics space. On one hand, the EM method can be studied as an *majorization-minimization* (MM) method in the parameter space. On the other hand, the EM method can be studied as a *scaled-gradient method* in the sufficient statistics space.
- Based on the two views described, we derive non-asymptotic convergence rate for stochastic EM methods. First, we show that the iEM method [Neal and Hinton, 1998] is a special instance of the MISO framework [Mairal, 2015], and takes  $\mathcal{O}(n/\epsilon)$  iterations to find an  $\epsilon$ -stationary point to the ML estimation problem. Second, the sEM-VR method [Chen et al., 2018] is an instance of variance reduced stochastic scaled-gradient method, which takes  $\mathcal{O}(n^{2/3}/\epsilon)$  iterations to find to an  $\epsilon$ -stationary point.
- Lastly, we develop a Fast Incremental EM (fiEM) method based on the SAGA algorithm [Defazio et al., 2014, Reddi et al., 2016b] for stochastic optimization. We show that the new method is again a scaled-gradient method with the same iteration complexity as sEM-VR. This new method offers trade-off between storage cost and computation complexity.

Importantly, our results capitalizes on the efficiency of stochastic EM methods applied on large datasets, and we support the above findings using numerical experiments.

**Prior Work** Since the empirical risk minimization problem (1) is typically *non-convex*, most prior work studying the convergence of EM methods considered either the *asymptotic* and/or *local* behaviors. For the classical study, the global convergence to a stationary point (either a local minimum or a saddle point) of the bEM method has been established by Wu et al. [1983] (by making the arguments

developed in [Dempster et al. \[1977\]](#) rigorous). The global convergence is a direct consequence of the EM method to be monotone. It is also known that in the neighborhood of a stationary point and under regularity conditions, the local rate of convergence of the bEM is linear and is given by the amount of *missing information* [[McLachlan and Krishnan, 2007](#), Chapters 3 and 4].

The convergence of the iEM method was first tackled by [Gunawardana and Byrne \[2005\]](#) exploiting the interpretation of the method as an alternating minimization procedure under the information geometric framework developed in [[Csiszár and Tusnády, 1984](#)]. Although the EM algorithm is presented as an alternation between the E-step and M-step, it is also possible to take a variational perspective on EM to view both steps as maximization steps. Nevertheless, [Gunawardana and Byrne \[2005\]](#) assume that the latent variables take only a finite number of values and the order in which the observations are processed remains the same from one pass to the other.

More recently, the *local but non-asymptotic convergence* of EM methods has been studied in several works. These results typically require the initializations to be within a neighborhood of an isolated stationary point and the (negated) log-likelihood function to be strongly convex locally. Such conditions are either difficult to verify in general or have been derived only for specific models; see for example [[Wang et al., 2015](#), [Xu et al., 2016](#), [Balakrishnan et al., 2017](#)] and the references therein. The local convergence of sEM-VR method has been studied in [[Chen et al., 2018](#), Theorem 1] but under a pathwise global stability condition. The authors' work [[Karimi et al., 2019](#)] provided the first global non-asymptotic analysis of the online (stochastic) EM method [[Cappé and Moulines, 2009](#)]. In comparison, the present work analyzes the variance reduced variants of EM method. Lastly, it is worthwhile to mention that [Zhu et al. \[2017\]](#) analyzed a variance reduced *gradient* EM method similar to [[Balakrishnan et al., 2017](#)].

## 2 Stochastic Optimization Techniques for EM methods

Let  $k \geq 0$  be the iteration number. The  $k$ th iteration of a generic stochastic EM method is composed of two sub-steps — firstly,

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

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,  $\mathcal{S}^{(k+1)}$  is a proxy for  $\bar{\mathbf{s}}(\hat{\boldsymbol{\theta}}^{(k)})$ , and  $\bar{\mathbf{s}}$  is defined in (3). Secondly, the M-step is given by

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

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

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

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

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

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

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

The stepsize is set to  $\gamma_{k+1} = 1$  for the iEM method;  $\gamma_{k+1} = \gamma$  is constant for the sEM-VR method. In the original version of the sEM method, the sequence of step  $\gamma_{k+1}$  is a diminishing step size. Moreover, for iEM we initialize with  $\mathcal{S}^{(0)} = \bar{\mathbf{s}}^{(0)}$ ; for sEM-VR, 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.

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

For iteration  $k \geq 0$ , the fiEM method draws *two* indices *independently* and uniformly 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  $t_{j_k}^k = \{k' : j_{k'} = j, k' < k\}$  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 the initialization  $\overline{\mathcal{S}}^{(0)} = \overline{\mathbf{s}}^{(0)}$ , we use a slightly different update rule from SAGA inspired by [Reddi et al., 2016b], as described by the following recursive updates

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

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

In the above, the update of  $\overline{\mathcal{S}}^{(k+1)}$  corresponds to an *unbiased estimate* of  $\overline{\mathbf{s}}^{(k)}$ , while the update for  $\overline{\mathcal{S}}^{(k+1)}$  maintains the structure that  $\overline{\mathcal{S}}^{(k)} = n^{-1} \sum_{i=1}^n \overline{\mathbf{s}}_i^{(t_i^k)}$  for any  $k \geq 0$ . The two updates of (10) are based on two different and independent indices  $i_k, j_k$  that are randomly drawn from  $\llbracket n \rrbracket$ . This is used for our fast convergence analysis in Section 3.

We summarize the iEM, sEM-VR, sEM, fiEM methods in Algorithm 1. The random termination number (11) is inspired by [Ghadimi and Lan, 2013] which enables one to show non-asymptotic convergence to stationary point for non-convex optimization. Due to their stochastic nature, the per-iteration complexity for all the stochastic EM methods are independent of  $n$ , unlike the bEM method. They are thus applicable to large datasets with  $n \gg 1$ .

---

**Algorithm 1** Stochastic EM methods.

---

- 1: **Input:** initializations  $\hat{\theta}^{(0)} \leftarrow 0, \hat{\mathbf{s}}^{(0)} \leftarrow \overline{\mathbf{s}}^{(0)}, K_{\max} \leftarrow \text{max. iteration number}$ .
- 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  $\overline{\mathcal{S}}^{(k+1)}$  using (8) or (7) or (9) or (10).
  - 6: Compute  $\hat{\mathbf{s}}^{(k+1)}$  via the sE-step (4).
  - 7: Compute  $\hat{\theta}^{(k+1)}$  via the M-step (5).
  - 8: **end for**
  - 9: **Return:**  $\hat{\theta}^{(K)}$ .
- 

## 2.1 Example: Gaussian Mixture Model

We discuss an example of learning a Gaussian Mixture Model (GMM) from a set of  $n$  observations  $\{y_i\}_{i=1}^n$ . We focus on a simplified setting where there are  $M$  components of unit variance and unknown means, the GMM is parameterized by  $\theta = (\{\omega_m\}_{m=1}^{M-1}, \{\mu_m\}_{m=1}^M) \in \Theta = \Delta^M \times \mathbb{R}^M$ , where  $\Delta^M \subseteq \mathbb{R}^{M-1}$  is the reduced  $M$ -dimensional probability simplex [see (29)]. 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  $\epsilon > 0$ . Furthermore, we use  $z_i \in \llbracket M \rrbracket$  as the latent label. The complete data log-likelihood is given by

$$\log f(z_i, y_i; \theta) = \sum_{m=1}^M \mathbb{1}_{\{m=z_i\}} [\log(\omega_m) - \mu_m^2/2] + \sum_{m=1}^M \mathbb{1}_{\{m=z_i\}} \mu_m y_i + \text{constant}, \quad (12)$$

where  $\mathbb{1}_{\{m=z_i\}} = 1$  if  $m = z_i$ ; otherwise  $\mathbb{1}_{\{m=z_i\}} = 0$ . The above can be rewritten in the same form as (2), particularly with  $S(y_i, z_i) \equiv (s_{i,1}^{(1)}, \dots, s_{i,M-1}^{(1)}, s_{i,1}^{(2)}, \dots, s_{i,M-1}^{(2)}, s_i^{(3)})$  and  $\phi(\theta) \equiv (\phi_1^{(1)}(\theta), \dots, \phi_{M-1}^{(1)}(\theta), \phi_1^{(2)}(\theta), \dots, \phi_{M-1}^{(2)}(\theta), \phi^{(3)}(\theta))$  such that

$$\begin{aligned} s_{i,m}^{(1)} &= \mathbb{1}_{\{z_i=m\}}, & \phi_m^{(1)}(\theta) &= \{\log(\omega_m) - \mu_m^2/2\} - \{\log(1 - \sum_{j=1}^{M-1} \omega_j) - \mu_M^2/2\}, \\ s_{i,m}^{(2)} &= \mathbb{1}_{\{z_i=m\}} y_i, & \phi_m^{(2)}(\theta) &= \mu_m, \quad s_i^{(3)} = y_i, \quad \phi^{(3)}(\theta) = \mu_M, \end{aligned} \quad (13)$$

and  $\psi(\theta) = -\{\log(1 - \sum_{m=1}^{M-1} \omega_m) - \mu_M^2/2\sigma^2\}$ . To evaluate the sE-step, the conditional expectation required by (6) can be computed in closed form, as they depend on  $\mathbb{E}_{\hat{\theta}^{(k)}}[\mathbb{1}_{\{z_i=m\}} | y = y_i]$  and  $\mathbb{E}_{\hat{\theta}^{(k)}}[y_i \mathbb{1}_{\{z_i=m\}} | y = y_i]$ . Moreover, the M-step (5) solves a strongly convex problem and can

be computed in closed form. Given a sufficient statistics  $\mathbf{s} \equiv (\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \mathbf{s}^{(3)})$ , the solution to (5) is:

$$\bar{\boldsymbol{\theta}}(\mathbf{s}) = \begin{pmatrix} (1 + \epsilon M)^{-1} (s_1^{(1)} + \epsilon, \dots, s_{M-1}^{(1)} + \epsilon)^\top \\ ((s_1^{(1)} + \delta)^{-1} s_1^{(2)}, \dots, (s_{M-1}^{(1)} + \delta)^{-1} s_{M-1}^{(2)})^\top \\ (1 - \sum_{m=1}^{M-1} s_m^{(1)} + \delta)^{-1} (s^{(3)} - \sum_{m=1}^{M-1} s_m^{(2)}) \end{pmatrix}. \quad (14)$$

The next section presents the main results of this paper for the convergence of stochastic EM methods. We shall use the above example on GMM to illustrate the required assumptions.

### 3 Global Convergence of Stochastic EM Methods

We establish non-asymptotic rates for the *global convergence* of the stochastic EM methods. We show that the iEM method is an instance of the incremental MM method; while sEM-VR, fiEM methods are instances of variance reduced *stochastic scaled gradient* methods. As we will see, the latter interpretation allows us to establish fast convergence rates of sEM-VR and fiEM methods. Detailed proofs for the theoretical results in this section are relegated to the appendix.

First, we list a few assumptions which will enable the convergence analysis performed later in this section. Define:

$$\mathcal{S} := \left\{ \sum_{i=1}^n \alpha_i \mathbf{s}_i : \mathbf{s}_i \in \text{conv} \{S(z, y_i) : z \in \mathcal{Z}\}, \alpha_i \in [-1, 1], i \in \llbracket 1, n \rrbracket \right\}, \quad (15)$$

where  $\text{conv}\{A\}$  denotes the closed convex hull of the set  $A$ . From (15), we observe that the iEM, sEM-VR, and fiEM methods generate  $\hat{\mathbf{s}}^{(k)} \in \mathcal{S}$  for any  $k \geq 0$ . Consider:

**H1.** *The sets  $\mathcal{Z}, \mathcal{S}$  are compact. There exists constants  $C_S, C_Z$  such that:*

$$C_S := \max_{\mathbf{s}, \mathbf{s}' \in \mathcal{S}} \|\mathbf{s} - \mathbf{s}'\| < \infty, \quad C_Z := \max_{i \in \llbracket 1, n \rrbracket} \int_{\mathcal{Z}} |S(z, y_i)| \mu(dz) < \infty. \quad (16)$$

**H1** depends on the latent data model used and can be satisfied by several practical models. For instance, the GMM in Section 2.1 satisfies (16) as the sufficient statistics are composed of indicator functions and observations. Other examples can also be found in Section 4. Denote by  $J_\kappa^\theta(\boldsymbol{\theta}')$  the Jacobian of the function  $\kappa : \boldsymbol{\theta} \mapsto \kappa(\boldsymbol{\theta})$  at  $\boldsymbol{\theta}' \in \Theta$ . Consider:

**H2.** *The function  $\phi$  is smooth and bounded on  $\text{int}(\Theta)$ , i.e., the interior of  $\Theta$ . For all  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \text{int}(\Theta)^2$ ,  $\|J_\phi^\theta(\boldsymbol{\theta}) - J_\phi^\theta(\boldsymbol{\theta}')\| \leq L_\phi \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|$  and  $\|J_\phi^\theta(\boldsymbol{\theta}')\| \leq C_\phi$ .*

**H3.** *The conditional distribution is smooth on  $\text{int}(\Theta)$ . For any  $i \in \llbracket 1, n \rrbracket$ ,  $z \in \mathcal{Z}$ ,  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \text{int}(\Theta)^2$ , we have  $|p(z|y_i; \boldsymbol{\theta}) - p(z|y_i; \boldsymbol{\theta}')| \leq L_p \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|$ .*

**H4.** *For any  $\mathbf{s} \in \mathcal{S}$ , the function  $\boldsymbol{\theta} \mapsto L(\mathbf{s}, \boldsymbol{\theta}) := R(\boldsymbol{\theta}) + \psi(\boldsymbol{\theta}) - \langle \mathbf{s} | \phi(\boldsymbol{\theta}) \rangle$  admits a unique global minimum  $\bar{\boldsymbol{\theta}}(\mathbf{s}) \in \text{int}(\Theta)$ . In addition,  $J_\phi^\theta(\bar{\boldsymbol{\theta}}(\mathbf{s}))$  is full rank and  $\bar{\boldsymbol{\theta}}(\mathbf{s})$  is  $L_\theta$ -Lipschitz.*

Under **H1**, the assumptions **H2** and **H3** are standard for the curved exponential family distribution and the conditional probability distributions, respectively; **H4** can be enforced by designing a strongly convex regularization function  $R(\boldsymbol{\theta})$  tailor made for  $\Theta$ . For instance, the penalization for GMM in Section 2.1 ensures  $\boldsymbol{\theta}^{(k)}$  is unique and lies in  $\text{int}(\Delta^M) \times \mathbb{R}^M$ , which can further imply the second statement in **H4**. We remark that for **H3**, it is possible to define the Lipschitz constant  $L_p$  independently for each data  $y_i$  to yield a refined characterization. We did not pursue such assumption to keep the notations simple.

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

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

**H5.** *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 a constant  $L_B$  such that for all  $\mathbf{s}, \mathbf{s}' \in \mathcal{S}^2$ , we have  $\|\mathbf{B}(\mathbf{s}) - \mathbf{B}(\mathbf{s}')\| \leq L_B \|\mathbf{s} - \mathbf{s}'\|$ .*

Again, **H5** is satisfied by practical models. For GMM in Section 2.1, it can be verified by deriving the closed form expression for  $\mathbf{B}(\mathbf{s})$  and using **H1**; also see the other example in Section 4. The derivation is, however, technical and will be relegated to the supplementary material.

Under H1, we have  $\|\hat{s}^{(k)}\| < \infty$  since  $S$  is compact. On the other hand, under H4, the EM methods generate  $\hat{\theta}^{(k)} \in \text{int}(\Theta)$  for any  $k \geq 0$ . Overall, these assumptions ensure that the EM methods operate in a ‘nice’ set throughout the optimization process.

### 3.1 Incremental EM method

We show that the iEM method is a special case of the MISO method [Mairal, 2015] utilizing the majorization minimization (MM) technique. The latter is a common technique for handling non-convex optimization. We begin by defining a surrogate function that majorizes  $\mathcal{L}_i$ :

$$Q_i(\theta; \theta') := - \int_{\mathcal{Z}} \{\log f(z_i, y_i; \theta) - \log p(z_i|y_i; \theta')\} p(z_i|y_i; \theta') \mu(dz_i). \quad (18)$$

The second term inside the bracket is a constant that does not depend on the first argument  $\theta$ . Since  $f(z_i, y_i; \theta) = p(z_i|y_i; \theta)g(y_i; \theta)$ , for all  $\theta' \in \Theta$ , we get  $Q_i(\theta; \theta') = -\log g(y_i; \theta') = \mathcal{L}_i(\theta')$ . For all  $\theta, \theta' \in \Theta$ , applying the Jensen inequality shows

$$Q_i(\theta, \theta') - \mathcal{L}_i(\theta) = \int \log \frac{p(z_i|y_i; \theta')}{p(z_i|y_i; \theta)} p(z_i|y_i; \theta') \mu(dz_i) \geq 0 \quad (19)$$

which is the Kullback-Leibler divergence between the conditional distribution of the latent data  $p(\cdot|y_i; \theta)$  and  $p(\cdot|y_i; \theta')$ . Hence, for all  $i \in \llbracket 1, n \rrbracket$ ,  $Q_i(\theta; \theta')$  is a majorizing surrogate to  $\mathcal{L}_i(\theta)$ , i.e., it satisfies for all  $\theta, \theta' \in \Theta$ ,  $Q_i(\theta; \theta') \geq \mathcal{L}_i(\theta)$  with equality when  $\theta = \theta'$ . For the special case of curved exponential family distribution, the M-step of the iEM method is expressed as

$$\begin{aligned} \hat{\theta}^{(k+1)} &\in \arg \min_{\theta \in \Theta} \left\{ R(\theta) + n^{-1} \sum_{i=1}^n Q_i(\theta; \hat{\theta}^{(\tau_i^{(k+1)})}) \right\} \\ &= \arg \min_{\theta \in \Theta} \left\{ R(\theta) + \psi(\theta) - \langle n^{-1} \sum_{i=1}^n \bar{s}_i^{(\tau_i^{(k+1)})} | \phi(\theta) \rangle \right\}. \end{aligned} \quad (20)$$

The iEM method can be interpreted through the MM technique — in the M-step,  $\hat{\theta}^{(k+1)}$  minimizes an upper bound of  $\bar{\mathcal{L}}(\theta)$ , while the sE-step updates the surrogate function in (20) which tightens the upper bound. Importantly, the error between the surrogate function and  $\mathcal{L}_i$  is a smooth function:

**Lemma 1.** *Assume H1, H2, H3, H4. Let  $e_i(\theta; \theta') := Q_i(\theta; \theta') - \mathcal{L}_i(\theta)$ . For any  $\theta, \bar{\theta}, \theta' \in \Theta^3$ , we have  $\|\nabla e_i(\theta; \theta') - \nabla e_i(\bar{\theta}; \theta')\| \leq L_e \|\theta - \bar{\theta}\|$ , where  $L_e := C_\phi C_Z L_p + C_S L_\phi$ .*

For non-convex optimization such as (1), it has been shown [Mairal, 2015, Proposition 3.1] that the incremental MM method converges asymptotically to a stationary solution of a problem. We strengthen their result by establishing a non-asymptotic rate, which is new to the literature.

**Theorem 1.** *Consider the iEM algorithm, i.e., Algorithm 1 with (7). Assume H1, H2, H3, H4. For any  $K_{\max} \geq 1$ , it holds that*

$$\mathbb{E}[\|\nabla \bar{\mathcal{L}}(\hat{\theta}^{(K)})\|^2] \leq n \frac{2L_e}{K_{\max}} \mathbb{E}[\bar{\mathcal{L}}(\hat{\theta}^{(0)}) - \bar{\mathcal{L}}(\hat{\theta}^{(K_{\max})})], \quad (21)$$

where  $L_e$  is defined in Lemma 1 and  $K$  is a uniform random variable on  $\llbracket 0, K_{\max} - 1 \rrbracket$  [cf. (11)] independent of the  $\{i_k\}_{k=0}^{K_{\max}}$ .

We remark that under suitable assumptions, our analysis in Theorem 1 also extends to several non-exponential family distribution models.

### 3.2 Stochastic EM as Scaled Gradient Methods

We interpret the sEM-VR and fiEM methods as *scaled gradient* methods on the sufficient statistics  $\hat{s}$ , tackling a *non-convex* optimization problem. The benefit of doing so is that we are able to demonstrate a faster convergence rate for these methods through motivating them as *variance reduced* optimization methods. The latter is shown to be more effective when handling large datasets [Reddi et al., 2016b,a, Allen-Zhu and Hazan, 2016] than traditional stochastic/deterministic optimization methods. To set our stage, we consider the minimization problem:

$$\min_{s \in S} V(s) := \bar{\mathcal{L}}(\bar{\theta}(s)) = R(\bar{\theta}(s)) + \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\bar{\theta}(s)), \quad (22)$$

where  $\bar{\theta}(s)$  is the unique map defined in the M-step (5). We first show that the stationary points of (22) has a one-to-one correspondence with the stationary points of (1):

**Lemma 2.** For any  $\mathbf{s} \in \mathcal{S}$ , it holds that

$$\nabla_{\mathbf{s}} V(\mathbf{s}) = \mathbf{J}_{\bar{\boldsymbol{\theta}}(\mathbf{s})}^{\mathbf{s}} \top \nabla_{\boldsymbol{\theta}} \bar{\mathcal{L}}(\bar{\boldsymbol{\theta}}(\mathbf{s})). \quad (23)$$

Assume H4. If  $\mathbf{s}^* \in \{\mathbf{s} \in \mathcal{S} : \nabla_{\mathbf{s}} V(\mathbf{s}) = 0\}$ , then  $\bar{\boldsymbol{\theta}}(\mathbf{s}^*) \in \{\boldsymbol{\theta} \in \Theta : \nabla_{\boldsymbol{\theta}} \bar{\mathcal{L}}(\boldsymbol{\theta}) = 0\}$ . Conversely, if  $\boldsymbol{\theta}^* \in \{\boldsymbol{\theta} \in \Theta : \nabla_{\boldsymbol{\theta}} \bar{\mathcal{L}}(\boldsymbol{\theta}) = 0\}$ , then  $\mathbf{s}^* = \bar{\mathbf{s}}(\boldsymbol{\theta}^*) \in \{\mathbf{s} \in \mathcal{S} : \nabla_{\mathbf{s}} V(\mathbf{s}) = 0\}$ .

The next lemmas show that the update direction,  $\hat{\mathbf{s}}^{(k)} - \boldsymbol{\mathcal{S}}^{(k+1)}$ , in the **sE-step** (4) of sEM-VR and fiEM methods is a *scaled gradient* of  $V(\mathbf{s})$ . We first observe the following conditional expectation:

$$\mathbb{E}[\hat{\mathbf{s}}^{(k)} - \boldsymbol{\mathcal{S}}^{(k+1)} | \mathcal{F}_k] = \hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)} = \hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}^{(k)})), \quad (24)$$

where  $\mathcal{F}_k$  is the  $\sigma$ -algebra generated by  $\{i_0, i_1, \dots, i_k\}$  (or  $\{i_0, j_0, \dots, i_k, j_k\}$  for fiEM).

The difference vector  $\mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s}))$  and the gradient vector  $\nabla_{\mathbf{s}} V(\mathbf{s})$  are correlated, as we observe:

**Lemma 3.** Assume H4, H5. For all  $\mathbf{s} \in \mathcal{S}$ ,

$$v_{\min}^{-1} \langle \nabla V(\mathbf{s}) | \mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s})) \rangle \geq \|\mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s}))\|^2 \geq v_{\max}^{-2} \|\nabla V(\mathbf{s})\|^2, \quad (25)$$

Combined with (24), the above lemma shows that the update direction in the **sE-step** (4) of sEM-VR and fiEM methods is a *stochastic scaled gradient* where  $\hat{\mathbf{s}}^{(k)}$  is updated with a stochastic direction whose mean is correlated with  $\nabla V(\mathbf{s})$ .

Furthermore, the expectation step's operator and the objective function in (22) are smooth functions:

**Lemma 4.** Assume H1, H3, H4, H5. For all  $\mathbf{s}, \mathbf{s}' \in \mathcal{S}$  and  $i \in \llbracket 1, n \rrbracket$ , we have

$$\|\bar{\mathbf{s}}_i(\bar{\boldsymbol{\theta}}(\mathbf{s})) - \bar{\mathbf{s}}_i(\bar{\boldsymbol{\theta}}(\mathbf{s}'))\| \leq L_{\mathbf{s}} \|\mathbf{s} - \mathbf{s}'\|, \quad \|\nabla V(\mathbf{s}) - \nabla V(\mathbf{s}')\| \leq L_V \|\mathbf{s} - \mathbf{s}'\|, \quad (26)$$

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

The following theorem establishes the (fast) non-asymptotic convergence rates of sEM-VR and fiEM methods, which are similar to [Reddi et al., 2016b,a, Allen-Zhu and Hazan, 2016]:

**Theorem 2.** Assume H1, H3, H4, H5. Denote  $\bar{L}_V = \max\{L_V, L_{\mathbf{s}}\}$  with the constants in Lemma 4.

- Consider the sEM-VR method, i.e., Algorithm 1 with (9). There exists a universal constant  $\mu \in (0, 1)$  (independent of  $n$ ) such that if we set the step size as  $\gamma = \frac{\mu v_{\min}}{L_V n^{2/3}}$  and the epoch length as  $m = \frac{n}{2\mu^2 v_{\min}^2 + \mu}$ , then for any  $K_{\max} \geq 1$  that is a multiple of  $m$ , it holds that

$$\mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(K)})\|^2] \leq n^{\frac{2}{3}} \frac{2\bar{L}_V}{\mu K_{\max}} \frac{v_{\max}^2}{v_{\min}^2} \mathbb{E}[V(\hat{\mathbf{s}}^{(0)}) - V(\hat{\mathbf{s}}^{(K_{\max})})]. \quad (27)$$

- Consider the fiEM method, i.e., Algorithm 1 with (10). Set  $\gamma = \frac{v_{\min}}{\alpha L_V n^{2/3}}$  such that  $\alpha = \max\{6, 1 + 4v_{\min}\}$ . For any  $K_{\max} \geq 1$ , it holds that

$$\mathbb{E}[\|\nabla V(\hat{\mathbf{s}}^{(K)})\|^2] \leq n^{\frac{2}{3}} \frac{\alpha^2 \bar{L}_V}{K_{\max}} \frac{v_{\max}^2}{v_{\min}^2} \mathbb{E}[V(\hat{\mathbf{s}}^{(0)}) - V(\hat{\mathbf{s}}^{(K_{\max})})]. \quad (28)$$

We recall that  $K$  in the above is a uniform and independent r.v. chosen from  $\llbracket K_{\max} - 1 \rrbracket$  [cf. (11)].

In the supplementary materials, we also provide a local convergence analysis for fiEM method which shows that the latter can achieve linear rate of convergence *locally* under a similar set of assumptions used in [Chen et al., 2018] for sEM-VR method.

**Comparing iEM, sEM-VR, and fiEM** Note that by (23) in Lemma 2, if  $\|\nabla_{\mathbf{s}} V(\hat{\mathbf{s}})\|^2 \leq \epsilon$ , then  $\|\nabla_{\boldsymbol{\theta}} \bar{\mathcal{L}}(\bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}))\|^2 = \mathcal{O}(\epsilon)$ , and vice versa, where the hidden constant is independent of  $n$ . In other words, the rates for iEM, sEM-VR, fiEM methods in Theorem 1 and 2 are comparable.

Importantly, the theorems show an intriguing comparison – to attain an  $\epsilon$ -stationary point with  $\|\nabla_{\boldsymbol{\theta}} \bar{\mathcal{L}}(\bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}))\|^2 \leq \epsilon$  or  $\|\nabla_{\mathbf{s}} V(\hat{\mathbf{s}})\|^2 \leq \epsilon$ , the iEM method requires  $\mathcal{O}(n/\epsilon)$  iterations (in expectation) while the sEM-VR, fiEM methods require only  $\mathcal{O}(n^{\frac{2}{3}}/\epsilon)$  iterations (in expectation). This comparison can be surprising since the iEM method is a monotone method as it guarantees decrease in the objective value; while the sEM-VR, fiEM methods are non-monotone. Nevertheless, it aligns with the recent analysis on stochastic variance reduction methods on non-convex problems. In the next section, we confirm the theory by observing a similar behavior numerically.

## 4 Numerical Examples

### 4.1 Gaussian Mixture Models

As described in Section 2.1, our goal is to fit a GMM model to a set of  $n$  observations  $\{y_i\}_{i=1}^n$  whose distribution is modeled as a Gaussian mixture of  $M$  components, each with a unit variance. Let  $z_i \in \llbracket M \rrbracket$  be the latent labels, the complete log-likelihood is given in (12), 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. The constraint set on  $\theta$  is given by

$$\Theta = \{\omega_m, m = 1, \dots, M-1 : \omega_m \geq 0, \sum_{m=1}^{M-1} \omega_m \leq 1\} \times \{\mu_m \in \mathbb{R}, m = 1, \dots, M\}. \quad (29)$$

In the following experiments of synthetic data, we generate samples from a GMM model with  $M = 2$  components with two mixtures with means  $\mu_1 = -\mu_2 = 0.5$ , see Appendix G.1 for details of the implementation and satisfaction of model assumptions for GMM inference. We aim at verifying the theoretical results in Theorem 1 and 2 of the dependence on sample size  $n$ .

**Fixed sample size** We use  $n = 10^4$  synthetic samples and run the bEM method until convergence (to double precision) to obtain the ML estimate  $\mu^*$ . We compare the bEM, sEM, iEM, sEM-VR and fiEM methods in terms of their precision measured by  $|\mu - \mu^*|^2$ . We set the stepsize of the sEM as  $\gamma_k = 3/(k+10)$ , and the stepsizes of the sEM-VR and the fiEM to a constant stepsize proportional to  $1/n^{2/3}$  and equal to  $\gamma = 0.003$ . The left plot of Figure 1 shows the convergence of the precision  $|\mu - \mu^*|^2$  for the different methods against the epoch(s) elapsed (one epoch equals  $n$  iterations). We observe that the sEM-VR and fiEM methods outperform the other methods, supporting our analytical results.

**Varying sample size** We compare the number of *iterations* required to reach a precision of  $10^{-3}$  as a function of the sample size from  $n = 10^3$  to  $n = 10^5$ . We average over 5 independent runs for each method using the same stepsizes as in the finite sample size case above. The right plot of Figure 1 confirms that our findings in Theorem 1 and 2 are sharp. It requires  $\mathcal{O}(n/\epsilon)$  (*resp.*  $\mathcal{O}(n^{2/3}/\epsilon)$ ) iterations to find a  $\epsilon$ -stationary point for the iEM (*resp.* sEM-VR and fiEM) method.

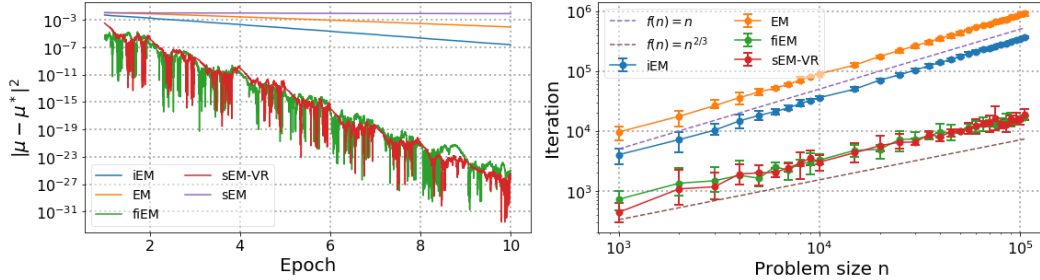


Figure 1: Performance of stochastic EM methods for fitting a GMM. (Left) Precision ( $|\mu^{(k)} - \mu^*|^2$ ) as a function of the epoch elapsed. (Right) Number of iterations to reach a precision of  $10^{-3}$ .

### 4.2 Probabilistic Latent Semantic Analysis

The second example considers probabilistic Latent Semantic Analysis (pLSA) whose aim is to classify documents into a number of topics. We are given a collection of documents  $\llbracket D \rrbracket$  with terms from a vocabulary  $\llbracket V \rrbracket$ . The data is summarized by a list of tokens  $\{y_i\}_{i=1}^n$  where each token is a pair of document and word  $y_i = (y_i^{(d)}, y_i^{(w)})$  which indicates that  $y_i^{(w)}$  appears in document  $y_i^{(d)}$ . The goal of pLSA is to classify the documents into  $K$  topics, which is modeled as a latent variable  $z_i \in \llbracket K \rrbracket$  associated with each token [Hofmann, 1999].

To apply stochastic EM methods for pLSA, we define  $\theta := (\theta^{(t,d)}, \theta^{(w,t)})$  as the parameter variable, where  $\theta^{(t,d)} = \{\theta_{d,k}^{(t,d)}\}_{\llbracket K-1 \rrbracket \times \llbracket D \rrbracket}$  and  $\theta^{(w,t)} = \{\theta_{k,v}^{(w,t)}\}_{\llbracket K \rrbracket \times \llbracket V-1 \rrbracket}$ . The constraint set  $\Theta$  is given as — for each  $d \in \llbracket D \rrbracket$ ,  $\theta_{d,\cdot}^{(t,d)} \in \Delta^K$  and for each  $k \in \llbracket K \rrbracket$ , we have  $\theta_{\cdot,k}^{(w,t)} \in \Delta^V$ , where  $\Delta^K, \Delta^V$



are the (reduced dimension)  $K, V$ -dimensional probability simplex; see (108) in the supplementary material for the precise definition. Furthermore, denote  $\theta_{d,K}^{(t|d)} = 1 - \sum_{k=1}^{K-1} \theta_{d,k}^{(t|d)}$  for each  $d \in \llbracket D \rrbracket$ , and  $\theta_{k,V}^{(w|t)} = 1 - \sum_{\ell=1}^{V-1} \theta_{k,\ell}^{(w|t)}$  for each  $k \in \llbracket K \rrbracket$ , the complete log likelihood for  $(y_i, z_i)$  is (up to an additive constant term):

$$\log f(z_i, y_i; \theta) = \sum_{k=1}^K \sum_{d=1}^D \log(\theta_{d,k}^{(t|d)}) \mathbb{1}_{\{k,d\}}(z_i, y_i^{(d)}) + \sum_{k=1}^K \sum_{v=1}^V \log(\theta_{k,v}^{(w|t)}) \mathbb{1}_{\{k,v\}}(z_i, y_i^{(w)}). \quad (30)$$

The penalization function is designed as

$$R(\theta^{(t|d)}, \theta^{(w|t)}) = -\log \text{Dir}(\theta^{(t|d)}; K, \alpha') - \log \text{Dir}(\theta^{(w|t)}; V, \beta'), \quad (31)$$

such that we ensure  $\bar{\theta}(s) \in \text{int}(\Theta)$ . We can apply the stochastic EM methods described in Section 2 on the pLSA problem. The implementation details are provided in Appendix G.2, therein we also verify the model assumptions required by our convergence analysis for pLSA.

**Experiment** We compare the stochastic EM methods on two FAO (UN Food and Agriculture Organization) datasets [Medelyan, 2009]. The first (*resp.* second) dataset consists of  $10^3$  (*resp.*  $10.5 \times 10^3$ ) documents and a vocabulary of size 300. The number of topics is set to  $K = 10$  and the stepsizes for the fiEM and sEM-VR are set to  $\gamma = 1/n^{2/3}$  while the stepsize for the sEM is set to  $\gamma_k = 1/(k + 10)$ . Figure 1 shows the evidence lower bound (ELBO) as a function of the number of epochs for the datasets. Again, the result shows that fiEM and sEM-VR methods achieve faster convergence than the competing EM methods, affirming our theoretical findings.

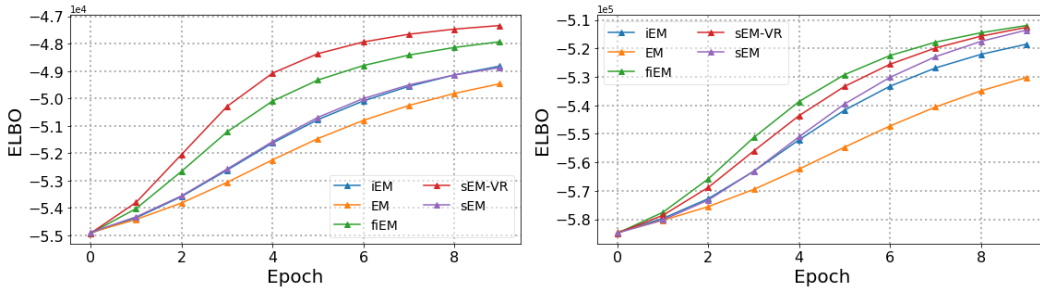


Figure 2: ELBO of the stochastic EM methods on FAO datasets as a function of number of epochs elapsed. (Left) small dataset with  $10^3$  documents. (Right) large dataset with  $10.5 \times 10^3$  documents.

## 5 Conclusion

This paper studies the global convergence for stochastic EM methods. Particularly, we focus on the inference of latent variable model with exponential family distribution and analyze the convergence of several stochastic EM methods. Our convergence results are *global* and *non-asymptotic*, and we offer two complimentary views on the existing stochastic EM methods — one interprets iEM method as an incremental MM method, and one interprets sEM-VR and fiEM methods as scaled gradient methods. The analysis shows that the sEM-VR and fiEM methods converge faster than the iEM method, and the result is confirmed via numerical experiments.

## Acknowledgement

BK and HTW contributed equally to this work. HTW’s work is supported by the CUHK Direct Grant #4055113.

## References

- P. Ablin, A. Gramfort, J.-F. Cardoso, and F. Bach. EM algorithms for ICA. *arXiv preprint arXiv:1805.10054*, 2018.
- Z. Allen-Zhu and E. Hazan. Variance reduction for faster non-convex optimization. In *International Conference on Machine Learning*, pages 699–707, 2016.
- S. Balakrishnan, M. J. Wainwright, and B. Yu. Statistical guarantees for the EM algorithm: From population to sample-based analysis. *Ann. Statist.*, 45(1):77–120, 02 2017. doi: 10.1214/16-AOS1435.
- D. M. Blei, A. Kucukelbir, and J. 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}.
- O. Cappé and E. 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.
- J. Chen, J. Zhu, Y. W. Teh, and T. Zhang. Stochastic expectation maximization with variance reduction. In *Advances in Neural Information Processing Systems*, pages 7978–7988, 2018.
- I. Csiszár and G. Tusnády. Information geometry and alternating minimization procedures. *Statist. Decisions*, suppl. 1:205–237, 1984. ISSN 0721-2631. Recent results in estimation theory and related topics.
- A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654, 2014.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38, 1977.
- S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- A. Gunawardana and W. Byrne. Convergence theorems for generalized alternating minimization procedures. *Journal of Machine Learning Research*, 6:2049–2073, 2005.
- G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- T. Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '99*, pages 50–57, New York, NY, USA, 1999. ACM. ISBN 1-58113-096-1. doi: 10.1145/312624.312649.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- B. Karimi, B. Miasojedow, E. Moulines, and H.-T. Wai. Non-asymptotic analysis of biased stochastic approximation schemes. In *Conference on Learning Theory*, 2019.
- J. Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.
- G. McLachlan and T. Krishnan. *The EM algorithm and extensions*, volume 382. John Wiley & Sons, 2007.

- O. Medelyan. *Human-competitive automatic topic indexing*. PhD thesis, The University of Waikato, 2009.
- R. M. Neal and G. E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, pages 355–368. Springer, 1998.
- S. Ng and G. McLachlan. On the choice of the number of blocks with the incremental EM algorithm for the fitting of normal mixtures. *Statistics and Computing*, 13(1):45–55, FEB 2003. ISSN 0960-3174. doi: {10.1023/A:1021987710829}.
- S. J. Reddi, A. Hefny, S. Sra, B. Póczos, and A. Smola. Stochastic variance reduction for nonconvex optimization. In *International conference on machine learning*, pages 314–323, 2016a.
- S. J. Reddi, S. Sra, B. Póczos, and A. Smola. Fast incremental method for nonconvex optimization. *arXiv preprint arXiv:1603.06159*, 2016b.
- B. Thiesson, C. Meek, and D. Heckerman. Accelerating EM for large databases. *Machine Learning*, 45(3):279–299, 2001. ISSN 0885-6125. doi: {10.1023/A:1017986506241}.
- M. J. Wainwright, M. I. Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- Z. Wang, Q. Gu, Y. Ning, and H. Liu. High dimensional em algorithm: Statistical optimization and asymptotic normality. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2521–2529. Curran Associates, Inc., 2015.
- C. J. Wu et al. On the convergence properties of the EM algorithm. *The Annals of statistics*, 11(1): 95–103, 1983.
- J. Xu, D. J. Hsu, and A. Maleki. Global analysis of expectation maximization for mixtures of two gaussians. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2676–2684. Curran Associates, Inc., 2016.
- R. Zhu, L. Wang, C. Zhai, and Q. Gu. High-dimensional variance-reduced stochastic gradient expectation-maximization algorithm. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 4180–4188. JMLR. org, 2017.