
Fast Two-Time-Scale Noisy EM Algorithms

Anonymous Author(s)

Affiliation

Address

email

Abstract

T.B.C

1 Introduction

We formulate the following empirical risk minimization problem as:

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

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 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 make the assumption of a complete model belonging to the curved exponential family, i.e.,

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

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.

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

Describe the main contributions of the paper as bullet points

Then small paragraph on the structure.

Some Notations

2 Two-Time-Scale Stochastic EM Algorithms

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

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

The **M-step** is given by

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

2.1 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 \tilde{s} of the deterministic quantity $\bar{s}(\theta)$:

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

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

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

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

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

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:

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

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

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

2.2 Incremental and Bi-Level Inexact EM Methods

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

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

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 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 the SAEM algorithm. Also if $\rho_k = 1$, $\gamma_k = 1$ and $\mathcal{S}^{(k)} = \tilde{S}^{(k)}$, then we recover the Monte Carlo EM algorithm.

We now introduce three variants of the SAEM update depending on different definitions of the proxy $\mathcal{S}^{(k)}$ 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 $\tau_i^k = \max\{k' : i_{k'} = i, k' < k\}$ be the iteration index where $i \in \llbracket 1, n \rrbracket$ is last drawn prior to iteration k . For iteration $k \geq 0$, the fiSAEM method draws *two* indices *independently* and uniformly 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 $t_j^k = \{k' : j_{k'} = j, k' < k\}$ 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 $\bar{\mathcal{S}}^{(0)} = \bar{s}^{(0)}$, 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.

2.3 Two-Time-Scale Noisy EM methods

We now introduce the general method derived using the two variance reduction techniques described above. Algorithm 1 leverages both levels (7) and (8) in order to output a vector of fitted parameters $\hat{\theta}^{(K)}$ where K is some randomly chosen termination point.

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 $\hat{s}^{(k+1)}$.

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 $\hat{S}_{i_k}^{(k)}$ 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 $\hat{S}^{(k+1)}$ 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:** $\hat{\theta}^{(K)}$.
-

3 Global and Finite Time Analysis of the Scheme

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

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

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

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

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

75 **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
 76 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.

77 Similar to [Karimi et al., 2019], we denote by $H_L^\theta(\mathbf{s}, \theta)$ the Hessian (w.r.t to θ for a given value of
 78 $\mathbf{s})$ of the function $\theta \mapsto L(\mathbf{s}, \theta) = R(\theta) + \psi(\theta) - \langle \mathbf{s} | \phi(\theta) \rangle$, and define

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

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

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

86 Indeed the expectations (3) are never available and requires Monte Carlo approximation. Thus, at
 87 iteration $k + 1$, we introduce the errors when approximating the quantity $\bar{s}_i(\hat{\theta}(\hat{\mathbf{s}}^{(k-1)}))$. For all
 88 $i \in \llbracket 1, n \rrbracket$, $r > 0$ and $\vartheta \in \Theta$, define:

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

89 For instance, we consider that the MC approximation is unbiased if for all $i \in \llbracket 1, n \rrbracket$ and $m \in$
 90 $\llbracket 1, M \rrbracket$, the samples $z_{i,m} \sim p(z_i|y_i; \theta)$ are i.i.d. under the posterior distribution, i.e., $\mathbb{E}[\eta_{i,\vartheta}^{(r)} | \mathcal{F}_r] = 0$
 91 where \mathcal{F}_r is the filtration up to iteration r .

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

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

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

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

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

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

98 3.1 Global Convergence of Incremental Noisy EM Algorithms

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

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

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

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

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

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

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

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

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

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

113 3.2 Global Convergence of Two-Time-Scale Noisy EM Algorithms

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

116 **TO COMPLETE**

117 4 Numerical Examples

118 4.1 Gaussian Mixture Models

119 Given n observations $\{y_i\}_{i=1}^n$, we want to fit a Gaussian Mixture Model (GMM) whose distribution
 120 is modeled as a Gaussian mixture of M components, each with a unit variance. Let $z_i \in \llbracket M \rrbracket$ be
 121 the latent labels of each component, the complete log-likelihood is defined as:

$$\log f(z_i, y_i; \theta) = \sum_{m=1}^M \mathbb{1}_{\{m\}}(z_i) [\log(\omega_m) - \mu_m^2/2] + \sum_{m=1}^M \mathbb{1}_{\{m\}}(z_i) \mu_m y_i + \text{constant} . \quad (23)$$

122 where $\theta := (\omega, \mu)$ with $\omega = \{\omega_m\}_{m=1}^{M-1}$ are the mixing weights with the convention $\omega_M =$
 123 $1 - \sum_{m=1}^{M-1} \omega_m$ and $\mu = \{\mu_m\}_{m=1}^M$ are the means. We use the penalization $R(\theta) = \frac{\delta}{2} \sum_{m=1}^M \mu_m^2 -$
 124 $\log \text{Dir}(\omega; M, \epsilon)$ where $\delta > 0$ and $\text{Dir}(\cdot; M, \epsilon)$ is the M dimensional symmetric Dirichlet distribu-
 125 tion with concentration parameter $\epsilon > 0$. The constraint set on θ is given by

$$\Theta = \{\omega_m, m = 1, \dots, M-1 : \omega_m \geq 0, \sum_{m=1}^{M-1} \omega_m \leq 1\} \times \{\mu_m \in \mathbb{R}, m = 1, \dots, M\}. \quad (24)$$

126 In the following experiments of synthetic data, we generate samples from a GMM model with $M =$
 127 2 components with two mixtures with means $\mu_1 = -\mu_2 = 0.5$.

128 We use $n = 10^4$ synthetic samples and run the bEM method until convergence (to double preci-
 129 sion) to obtain the ML estimate μ^* . We compare the bEM, SAEM, iSAEM, vrSAEM and fiSAEM
 130 methods in terms of their precision measured by $|\mu - \mu^*|^2$. The left plot of Figure ?? shows the
 131 convergence of the precision $|\mu - \mu^*|^2$ for the different methods against the epoch(s) elapsed (one
 132 epoch equals n iterations). We observe that the vrSAEM and fiSAEM methods outperform the other
 133 methods, supporting our analytical results.

134 4.2 Deformable Template Model for Image Analysis

135 We now run our different methods using an example taken from [?]. Let $(y_i, i \in \llbracket 1, n \rrbracket)$ be observed
 136 images. Let $u \in \mathcal{U} \subset \mathbb{R}^2$ denote the pixel index on the image and $x_u \in \mathcal{D} \subset \mathbb{R}^2$ its location.

137 The model used in this experiment suggests that each image y_i is a deformation of a template, noted
 138 $I : \mathcal{D} \rightarrow \mathbb{R}$, common to all images of the dataset:

$$y_i(u) = I(x_u - \Phi_i(x_u)) + \varepsilon_i(u) \quad (25)$$

139 where $\phi_i : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is a deformation function, and $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ is an observation error.

140 The template model, given $(p_k, k \in \llbracket 1, k_p \rrbracket)$ landmarks on the template, a fixed known kernel \mathbf{K}_p
 141 and a vector of parameters $\beta \in \mathbb{R}^{k_p}$ is defined as follows:

$$I_\xi = \mathbf{K}_p \beta, \quad \text{where} \quad (\mathbf{K}_p \beta)(x) = \sum_{k=1}^{k_p} \mathbf{K}_p(x, p_k) \beta_k \quad (26)$$

142 Besides, we parameterize the deformation model given some landmarks $(g_k, k \in \llbracket 1, k_g \rrbracket)$ and a
 143 fixed kernel $\mathbf{K}_{\mathbf{g}}$ as:

$$\Phi_i(x) = (\mathbf{K}_{\mathbf{g}} z_i)(x) = \sum_{k=1}^{k_s} \mathbf{K}_{\mathbf{g}}(x, g_k) \left(z_i^{(1)}(k), z_i^{(2)}(k) \right) \quad (27)$$

144 where $z_i \sim (0, \Gamma)$ and $z_i \in (\mathbb{R}^{k_g})^2$.

145 **5 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.

159 A Proof of Theorem 1

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

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

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

$$v_{\min}^{-1} \langle \nabla V(\mathbf{s}) | \mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s})) \rangle \geq \|\mathbf{s} - \bar{\mathbf{s}}(\bar{\boldsymbol{\theta}}(\mathbf{s}))\|^2 \geq v_{\max}^{-2} \|\nabla V(\mathbf{s})\|^2, \quad (28)$$

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

$$V(\hat{\mathbf{s}}^{(k+1)}) \leq V(\hat{\mathbf{s}}^{(k)}) - \gamma_{k+1} \langle \hat{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle + \frac{\gamma_{k+1}^2 L_V}{2} \|\hat{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)}\|^2 \quad (29)$$

166 Taking the expectation on both sides and using the growth condition (28), we obtain:

$$\begin{aligned} \mathbb{E}[V(\hat{\mathbf{s}}^{(k+1)})] &\leq \mathbb{E}[V(\hat{\mathbf{s}}^{(k)})] - \gamma_{k+1} v_{\min} \mathbb{E} \left[\left\| \bar{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(k)} \right\|^2 \right] + \mathbb{E} \left[\frac{\gamma_{k+1}^2 L_V}{2} \|\hat{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)}\|^2 \right] \\ &\quad - \gamma_{k+1} \mathbb{E} \left[\langle \bar{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right] \end{aligned} \quad (30)$$

167 We then establish an auxiliary Lemma yielding an upper-bound on the quantity

168 $\mathbb{E} \left[\langle \bar{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right]$ where:

$$\bar{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)} = \bar{\mathbf{s}}^{(k)} - \left(\tilde{\mathbf{S}}^{(k)} + \frac{1}{n} (\tilde{\mathbf{S}}_{i_k}^{(k)} - \tilde{\mathbf{S}}_{i_k}^{(\tau_{i_k}^k)}) \right) \quad (31)$$

169

Lemma 3.

$$\mathbb{E} \left[\langle \bar{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right] \leq \quad (32)$$

170 Using Lemma 2:

$$\begin{aligned} \mathbb{E}[V(\hat{\mathbf{s}}^{(k+1)})] &\leq \mathbb{E}[V(\hat{\mathbf{s}}^{(k)})] - \gamma_{k+1} \left(v_{\min} - \frac{\gamma_{k+1} L_V}{2} \right) \mathbb{E} \left[\left\| \bar{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(k)} \right\|^2 \right] \\ &\quad + \gamma_{k+1}^2 L_V L_s^2 \left(1 - \frac{1}{n} \right)^2 n^{-1} \sum_{i=1}^n \mathbb{E} \left[\left\| \hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(\tau_i^k)} \right\|^2 \right] + \frac{\gamma_{k+1}^2 L_V C}{M_k} \\ &\quad - \gamma_{k+1} \mathbb{E} \left[\langle \bar{\mathbf{s}}^{(k)} - \tilde{\mathbf{S}}^{(k+1)} | \nabla V(\hat{\mathbf{s}}^{(k)}) \rangle \right] \end{aligned} \quad (33)$$

171 Besides,

$$n^{-1} \sum_{i=1}^n \mathbb{E} \left[\left\| \hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(\tau_i^{k+1})} \right\|^2 \right] = n^{-1} \sum_{i=1}^n \left(\frac{1}{n} \mathbb{E} [\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2] + \frac{n-1}{n} \mathbb{E} [\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(\tau_i^k)}\|^2] \right) \quad (34)$$

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

$$\begin{aligned} &\mathbb{E} [\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2] \\ &= \mathbb{E} \left[\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2 + \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 + 2 \langle \hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)} | \hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)} \rangle \right] \\ &= \mathbb{E} \left[\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2 + \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 - 2 \gamma_{k+1} \langle \hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)} | \hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)} \rangle \right] \\ &\leq \mathbb{E} \left[\|\hat{\mathbf{s}}^{(k+1)} - \hat{\mathbf{s}}^{(k)}\|^2 + \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 + \frac{\gamma_{k+1}}{\beta_{k+1}} \|\hat{\mathbf{s}}^{(k)} - \bar{\mathbf{s}}^{(k)}\|^2 + \gamma_{k+1} \beta_{k+1} \|\hat{\mathbf{s}}^{(k)} - \hat{\mathbf{s}}^{(t_i^k)}\|^2 \right] \end{aligned} \quad (35)$$

173

□