
Distributed and Private Stochastic EM Methods via Quantized and Compressed MCMC

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To be completed

2 1 Introduction

3 We consider the distributed minimization of the following negated log incomplete data likelihood

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

4 where n denotes the number of workers, $\{y_i\}_{i=1}^n$ are observations, $\theta \in \mathbb{R}^d$ is the parameters set and
5 $R : \theta \rightarrow \mathbb{R}$ is a smooth regularizer.

6 The objective $L(\theta)$ is possibly nonconvex and is assumed to be lower bounded. In the latent data
7 model, the likelihood $g(y_i; \theta)$, is the marginal distribution of the complete data likelihood, noted
8 $f(z_i, y_i; \theta)$, such that

$$g(y_i; \theta) = \int_{\mathcal{Z}} f(z_i, y_i; \theta) \mu(dz_i), \quad (2)$$

9 where $\{z_i\}_{i=1}^n$ are the vectors of latent variables associated to the observations $\{y_i\}_{i=1}^n$.

10 We also consider a special case of that problem since the complete likelihood pertains to the curved
11 exponential family:

$$f(z_i, y_i; \theta) = h(z_i, y_i) \exp(\langle S(z_i, y_i) | \phi(\theta) \rangle - \psi(\theta)), \quad (3)$$

12 where $\psi(\theta)$, $h(z_i, y_i)$ are scalar functions, $\phi(\theta) \in \mathbb{R}^k$ is a vector function, and $\{S(z_i, y_i) \in \mathbb{R}^k\}_{i=1}^n$
13 is the vector of sufficient statistics. We refer the readers to [Efron et al., 1975] for details on this
14 subclass of problems which is of high interest given the broad range of problems that fall under
15 this assumption. In the centralized settings, *i.e.*, when all data points are stored in a central server,
16 a reference tool for learning such a model is called the EM algorithm [Dempster et al., 1977, Wu,
17 1983]. Comprised of two steps, the E-step computes an aggregated sum of expectations as follows:

$$\bar{s}(\theta) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta) \quad \text{where} \quad \bar{s}_i(\theta) := \int_{\mathcal{Z}} S(z_i, y_i) p(z_i | y_i; \theta) dz_i, \quad (4)$$

18 and the M-step is given by

$$\bar{\theta}(\bar{s}(\theta)) := \arg \min_{\vartheta \in \Theta} \{ r(\vartheta) + \psi(\vartheta) - \langle \bar{s}(\theta) | \phi(\vartheta) \rangle \}. \quad (5)$$

19 **1.1 Our motivations**

20 **Expectations are not tractable:** Sampling for those approximations are costly.

21 **Need for distributed computing:** MovieLens, Large n, compute time, decentralized infrastructure

22 **Need for privacy and communication efficiency:** Sensible data (hospital, user data...) that can
23 not be moved. Low bandwidth devices (compute should be light).

24 **1.2 Our contributions**

25 2 Related Work

26 EM algorithms:

27 In the case when the computation of the expectation under the posterior distribution is impossible,
28 the Monte Carlo EM (MCEM) has been introduced in [Wei and Tanner \[1990\]](#) where a Monte Carlo
29 (MC) approximation for this expectation is computed. A variant of that algorithm is the Stochastic
30 Approximation of the EM (SAEM) in [Delyon et al. \[1999\]](#) leveraging the power of Robbins-Monro
31 update [Robbins and Monro \[1951\]](#) to ensure pointwise convergence of the vector of estimated pa-
32 rameters using a decreasing stepsize rather than increasing the number of MC samples. The MCEM
33 and the SAEM have been successfully applied in mixed effects models [McCulloch \[1997\]](#), [Hughes](#)
34 [\[1999\]](#), [Baey et al. \[2016\]](#) or to do inference for joint modeling of time to event data coming from
35 clinical trials in [Chakraborty and Das \[2010\]](#), unsupervised clustering in [Ng and McLachlan \[2003\]](#),
36 variational inference of graphical models in [Blei et al. \[2017\]](#) among other applications. An incre-
37 mental variant of the SAEM was proposed in [Kuhn et al. \[2019\]](#) showing positive empirical results
38 but its analysis is limited to asymptotic consideration. Two-timescale methods of the SAEM have
39 been proposed in [\[Karimi and Li, 2020\]](#) to accelerate the convergence. Gradient-based methods
40 have been developed and analyzed in [Zhu et al. \[2017\]](#) but they remain out of the scope of this paper
41 as they tackle the high-dimensionality issue.

42 Distributed methods:

43 [\[Morral et al., 2012\]](#) [\[Srivastava et al., 2019\]](#)

44 Traditional decentralized optimization methods include well-know algorithms such as
45 ADMM [\[Boyd et al., 2011\]](#), Dual Averaging [\[Duchi et al., 2011\]](#), Distributed Subgradient
46 Descent [\[Nedic and Ozdaglar, 2009\]](#). More recent algorithms include Extra [\[Shi et al., 2015\]](#),
47 Next [\[Di Lorenzo and Scutari, 2016\]](#), Prox-PDA [\[?\]](#), GNSD [\[?\]](#), and Choco-SGD [\[?\]](#). While
48 these algorithms are commonly used in applications other than deep learning, recent algorithmic
49 advances in the machine learning community have shown that decentralized optimization can also
50 be useful for training deep models such as neural networks. ? demonstrate that a stochastic version
51 of Decentralized Subgradient Descent can outperform parameter server-based algorithms when
52 the communication cost is high. No existing work, to our knowledge, has seriously considered
53 integrating *EM methods* in the setting of decentralized learning. One noteworthy work [\[Morral](#)
54 [et al., 2012\]](#) proposes a decentralized version of the Online EM [\[Cappé and Moulines, 2009\]](#) and it
55 is proven to satisfy some non-standard convergence.

56 MCMC and Quantization:

57 [\[Chopin and Ducrocq, 2021\]](#)

58 [\[Vono et al., 2021\]](#)

59 Federated Learning methods:

60 3 On the Decentralization of the EM algorithm

61 3.1 Distributed SAEM

62 We first consider the plain distributed version of the sEM which does not tackle any privacy or
63 communication bottlenecks. We precise that we perform periodic locals models averaging. It goes
64 as follows:

Algorithm 1 Distributed SAEM with Periodic Locals Models Averaging

- 1: **Input:** Compression operator $\mathcal{C}(\cdot)$, number of rounds R , initial parameter θ_0 .
- 2: **for** $r = 1$ to R **do**
- 3: **for** parallel for device $i \in D^r$ **do**
- 4: Set $\hat{\theta}_i^{(r)} = \hat{\theta}^{(r)}$. {Initialize each worker with current global model}
- 5: Draw M samples $z_{i,m}^{(r+1)}$ under model $\hat{\theta}_i^{(r)}$ via MCMC: {Local MCMC step}
- 6: Compute the local statistics $\tilde{S}_i^{(r+1)} = S(z_{i,m}^{(r+1)})$. {Local statistics}
- 7: Worker computes **local model**: {(Local) M-Step using local statistics}

$$\hat{\theta}_i^{(r+1)} = \bar{\theta}(\tilde{S}_i^{(r+1)})$$

- 8: Worker sends local model $\hat{\theta}_i^{(r+1)}$ to server.
- 9: **end for**
- 10: Server computes **global model** by periodic averaging {Local model averaging}

$$\hat{\theta}^{(r+1)} := \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i^{(r+1)}$$

- 11: **end for**
-

65 3.2 Federated SAEM with Quantization and Compression

66 While Algorithm 2 is a distributed variant of the SAEM, it is neither (a) private nor (b)
67 communication-efficient.

68 **Privacy:** Indeed, we remark that broadcasting the vector of statistics are a potential breach to the
69 data observations as their expression is related y and the latent data z . With a simple knowledge of
70 the model used, the data could be retrieved if one extracts those statistics.

71 **Communication bottlenecks:** Also regarding (b), the broadcast of n vector of statistics $S(y_i, z_i)$
72 can be cumbersome when the size of the latent space and the parameter space of the model are huge.

73 For computational purposes and privacy enhanced matter, I have chosen to study and develop the
74 second algorithms that I proposed in my last week's report. In that algorithm, one does not compute
75 a periodic averaging of the local models (this would requires performing as many M-steps as there
76 are workers). Rather, workers compute local statistics and send them to the central server for a
77 periodic averaging of those vectors and the latter computes one M-step to update the global model.

Algorithm 2 FL-SAEM with Periodic Statistics Averaging

- 1: **Input:** number of rounds R , initial parameter θ_0 , number of MC samples $\{M_r\}_{r>0}$.
- 2: Init: $\theta_0 \in \Theta \subseteq \mathbb{R}^d$, as the global model and $\bar{\theta}_0 = \frac{1}{n} \sum_{i=1}^n \theta_0$.
- 3: **for** $r = 1$ to R **do**
- 4: **for** parallel for device $i \in D^r$ **do**
- 5: Set $\hat{\theta}_i^{(0,r)} = \hat{\theta}^{(r)}$.
- 6: Draw M_r samples $z_{i,m}^{(r)}$ under model $\hat{\theta}_i^{(r)}$
- 7: Compute the surrogate sufficient statistics $\tilde{S}_i^{(r+1)}$
- 8: Workers send local statistics $\tilde{S}_i^{(k+1)}$ to server.
- 9: **end for**
- 10: Server computes **global model using the aggregated statistics:**

$$\hat{\theta}^{(r+1)} = \bar{\theta}(\tilde{S}^{(r+1)})$$

where $\tilde{S}^{(r+1)} = (\tilde{S}_i^{(r+1)}, i \in D_r)$ and send global model back to the devices.

- 11: **end for**
-

78 3.3 Embedded methods to comply with Federated settings

79 **Line 6 – Quantization:** The first step is to quantize the gradient in the Stochastic Langevin Dynam-
80 ics step used in our sampling scheme Line 6 of Algorithm 2. Inspired by [Alistarh et al., 2017], we
81 use an extension of the QSGD algorithm for our latent samples. Define the quantization operator as
82 follows:

$$\mathcal{C}_j^{(\ell)}(g, \xi_j) = \|v\| \cdot \text{sign}(g_j) \cdot (\lfloor \ell |g_j| / \|v\| \rfloor + \mathbf{1}\{\xi_j \leq \ell |g_j| / \|v\| - \lfloor \ell |g_j| / \|v\| \rfloor\}) / \ell \quad (6)$$

83 where ℓ is the level of quantization and $j \in [d]$ denotes the dimension of the gradient.

84 Hence, for the sampling step, Line 6, we use the modified SGLD below, to be compliant with the
85 privacy of our method.

Algorithm 3 Langevin Dynamics with Quantization for worker i

- 1: **Input:** Current local model $\hat{\theta}_i^{(r)}$ for worker $i \in \llbracket 1, n \rrbracket$.
- 2: Draw M samples $\{z_i^{(r,m)}\}_{m=1}^M$ from the posterior distribution $p(z_i|y_i; \hat{\theta}_i^{(k)})$ via Langevin diffusion with a quantized gradient:
- 3: **for** $k = 1$ to K **do**
- 4: Compute the quantized gradient of $\nabla \log p(z_i|y_i; \hat{\theta}_i^{(k)})$:

$$g_i(k, m) = \mathcal{C}_j^{(\ell)}\left(\nabla_j f_{\theta_t}(z_i^{(k-1,m)}), \xi_j^{(k)}\right) \quad (7)$$

where $\xi_j^{(k)}$ is a realization of a uniform random variable.

- 5: Sample the latent data using the following chain:

$$z_i^{(k,m)} = z_i^{(k-1,m)} + \frac{\gamma_k}{2} g_i(k, m) + \sqrt{\gamma_k} B_k, \quad (8)$$

where B_t denotes the Brownian motion and $m \in [M]$ denotes the MC sample.

- 6: **end for**
 - 7: Assign $\{z_i^{(r,m)}\}_{m=1}^M \leftarrow \{z_i^{(K,m)}\}_{m=1}^M$.
 - 8: **Output:** latent data $z_{i,m}^{(k)}$ under model $\hat{\theta}_i^{(t,k)}$
-

86 **Line 7 – Compression MCMC output:** We use the notorious **Top- k** operator that we define as
87 $\mathcal{C}(x)_i = x_i$, if $i \in \mathcal{S}$; $\mathcal{C}(x)_i = 0$ otherwise and where \mathcal{S} is defined as the size- k set of $i \in [p]$.
88 Recall that after Line 6 we compute the local statistics $\tilde{S}_i^{(k+1)}$ using the output latent variables from
89 Algorithm 3. We now use those statistics and compress them using Algorithm 4 as follows:

Algorithm 4 Sparsified Statistics with **Top- k**

- 1: **Input:** Current local statistics $\tilde{S}_i^{(k+1)}$ for worker $i \in \llbracket 1, n \rrbracket$. Sparsification level k .
2: Apply **Top- k** :

$$\ddot{S}_i^{(k+1)} = \mathcal{C} \left(\tilde{S}_i^{(k+1)} \right) \quad (9)$$

- 3: **Output:** Compressed local statistics for worker i denoted $\ddot{S}_i^{(k+1)}$.
-

90 We present our final method in Algorithm 5, that performs SAEM under the federated settings.

Algorithm 5 fl-SAEM: Quantized and Compressed FL-SAEM with Periodic Statistics Averaging

- 1: **Input:** Compression operator $\mathcal{C}(\cdot)$, number of rounds R , initial parameter θ_0 .
2: **for** $r = 1$ to R **do**
3: **for** parallel for device $i \in D^r$ **do**
4: Set $\hat{\theta}_i^{(0,r)} = \hat{\theta}^{(r)}$. {Initialize each worker with current global model}
5: Draw M samples $z_{i,m}^{(r)}$ under model $\hat{\theta}_i^{(r)}$ via Quantized LD: {Local Quantized MCMC step}
6: **for** $k = 1$ to K **do**
7: Compute the quantized gradient of $\nabla \log p(z_i | y_i; \hat{\theta}_i^{(k)})$:

$$g_i(k, m) = \mathcal{C}_j^{(\ell)} \left(\nabla_j f_{\theta_t}(z_i^{(k-1,m)}), \xi_j^{(k)} \right) \quad \text{where} \quad \xi_j^{(k)} \sim \mathcal{U}_{[a,b]}$$

- 8: Sample the latent data using the following chain:

$$z_i^{(k,m)} = z_i^{(k-1,m)} + \frac{\gamma_k}{2} g_i(k, m) + \sqrt{\gamma_k} \mathbf{B}_k,$$

where \mathbf{B}_t denotes the Brownian motion and $m \in [M]$ denotes the MC sample.

- 9: **end for**
10: Assign $\{z_i^{(r,m)}\}_{m=1}^M \leftarrow \{z_i^{(K,m)}\}_{m=1}^M$.
11: Compute $\tilde{S}_i^{(r+1)}$ and its **Top- k** variant $\ddot{S}_i^{(r+1)} = \mathcal{C} \left(\tilde{S}_i^{(r+1)} \right)$. {Compressed local statistics}
12: Worker send local statistics $\ddot{S}_i^{(r+1)}$ to server. {Single round of communication}
13: **end for**
14: Server computes **global model**: {(Global) M-Step using aggregated statistics}

$$\hat{\theta}^{(r+1)} = \bar{\theta}(\ddot{S}^{(r+1)})$$

where $\ddot{S}^{(r+1)} = (\ddot{S}_i^{(r+1)}, i \in D_r)$ and send global model back to the devices.

- 15: **end for**
-

91 4 Theoretical Analysis

92 The following assumptions are required for the analysis.

93 **H1.** *The sets Z, S are compact. There exist C_S, C_Z such that:*

$$C_S := \max_{s, s' \in S} \|s - s'\| < \infty,$$

$$C_Z := \max_{i \in \llbracket 1, n \rrbracket} \int_Z |S(z, y_i)| \mu(dz) < \infty.$$

94 **H2.** *For any $i \in \llbracket 1, n \rrbracket$, $z \in Z$, $\theta, \theta' \in \text{int}(\theta)^2$ (the interior of θ), we have $|p(z|y_i; \theta) - p(z|y_i; \theta')| \leq$*
 95 $L_p \|\theta - \theta'\|$.

96 We also recall that we consider curved exponential family models such that the objective function
 97 satisfies:

98 **H3.** *For any $s \in S$, the function $\theta \mapsto L(s, \theta) := R(\theta) + \psi(\theta) - \langle s | \phi(\theta) \rangle$ admits a unique global*
 99 *minimum $\bar{\theta}(s) \in \text{int}(\theta)$. In addition, $J_\phi^\theta(\bar{\theta}(s))$, the Jacobian of the function ϕ at θ , is full rank,*
 100 *L_p -Lipschitz and $\bar{\theta}(s)$ is L_t -Lipschitz.*

101 The Monte Carlo noise of $\tilde{S}_i^{(k+1)}$ at iteration k is defined as:

$$\eta_i^{(k)} := \tilde{S}_i^{(k)} - \bar{s}_i(\vartheta^{(k)}) \quad \text{for all } i \in \llbracket 1, n \rrbracket \quad \text{and } k > 0 \quad (10)$$

102 and is controlled

103 **H4.** *For all $k > 0$, $i \in \llbracket 1, n \rrbracket$, it holds: $\mathbb{E}[\|\eta_i^{(k)}\|^2] < \infty$ and $\mathbb{E}[\|\mathbb{E}[\eta_i^{(k)} | \mathcal{F}_k]\|^2] < \infty$.*

104 Note that typically, the controls exhibited above are vanishing when the number of MC samples M_k
 105 increases with k .

106 4.1 Finite-time convergence analysis of the d-SAEM

107 4.2 Finite-time convergence analysis of the fl-SAEM

108 **5 Numerical Experiments**

109 **5.1 Nonlinear Mixed Models under Distributed Settings**

110 Compare SAEM, MCEM, dist-SAEM and maybe one distributed Gradient Descent as baseline

111 Same for Private settings with Sketched SGD or another good baseline

112 Fitting a linear mixed model on Oxford boys dataset [[Pinheiro and Bates, 2006](#)]

113 Fitting a nonlinear mixed model on Warfarin dataset [[Consortium, 2009](#)]

114 **5.2 Probabilistic Latent Dirichlet Allocation**

115 **5.3 Bi-factor models under the Federated Learning settings**

References

- D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.
- C. Baey, S. Trevezas, and P.-H. Cournède. A non linear mixed effects model of plant growth and estimation via stochastic variants of the em algorithm. *Communications in Statistics-Theory and Methods*, 45(6):1643–1669, 2016.
- 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}.
- S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- 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.
- A. Chakraborty and K. Das. Inferences for joint modelling of repeated ordinal scores and time to event data. *Computational and mathematical methods in medicine*, 11(3):281–295, 2010.
- N. Chopin and G. Ducrocq. Fast compression of mcmc output. *Entropy*, 23(8):1017, 2021.
- I. W. P. Consortium. Estimation of the warfarin dose with clinical and pharmacogenetic data. *New England Journal of Medicine*, 360(8):753–764, 2009.
- 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>.
- 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.
- P. Di Lorenzo and G. Scutari. Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136, 2016.
- J. C. Duchi, A. Agarwal, and M. J. Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2011.
- B. Efron et al. Defining the curvature of a statistical problem (with applications to second order efficiency). *The Annals of Statistics*, 3(6):1189–1242, 1975.
- J. P. Hughes. Mixed effects models with censored data with application to hiv rna levels. *Biometrics*, 55(2):625–629, 1999.
- B. Karimi and P. Li. Two-timescale stochastic em algorithms. In *IEEE International Symposium on Information Theory (ISIT)*, 2020.
- E. Kuhn, C. Matias, and T. Rebafka. Properties of the stochastic approximation em algorithm with mini-batch sampling. *arXiv preprint arXiv:1907.09164*, 2019.
- C. E. McCulloch. Maximum likelihood algorithms for generalized linear mixed models. *Journal of the American statistical Association*, 92(437):162–170, 1997.
- G. Morral, P. Bianchi, and J. Jakubowicz. On-line gossip-based distributed expectation maximization algorithm. In *2012 IEEE Statistical Signal Processing Workshop (SSP)*, pages 305–308. IEEE, 2012.

- 161 A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE*
162 *Transactions on Automatic Control*, 54(1):48, 2009.
- 163 S. Ng and G. McLachlan. On the choice of the number of blocks with the incremental EM algorithm
164 for the fitting of normal mixtures. *Statistics and Computing*, 13(1):45–55, FEB 2003. ISSN 0960-
165 3174. doi: {10.1023/A:1021987710829}.
- 166 J. Pinheiro and D. Bates. *Mixed-effects models in S and S-PLUS*. Springer Science & Business
167 Media, 2006.
- 168 H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statis-*
169 *tics*, pages 400–407, 1951.
- 170 W. Shi, Q. Ling, G. Wu, and W. Yin. Extra: An exact first-order algorithm for decentralized consen-
171 sus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- 172 S. Srivastava, G. DePalma, and C. Liu. An asynchronous distributed expectation maximization al-
173 gorithm for massive data: the dem algorithm. *Journal of Computational and Graphical Statistics*,
174 28(2):233–243, 2019.
- 175 M. Vono, V. Plassier, A. Durmus, A. Dieuleveut, and E. Moulines. Qlsd: Quantised langevin
176 stochastic dynamics for bayesian federated learning. *arXiv preprint arXiv:2106.00797*, 2021.
- 177 G. C. Wei and M. A. Tanner. A monte carlo implementation of the em algorithm and the poor man’s
178 data augmentation algorithms. *Journal of the American statistical Association*, 85(411):699–704,
179 1990.
- 180 C. J. Wu. On the convergence properties of the EM algorithm. *The Annals of statistics*, pages
181 95–103, 1983.
- 182 R. Zhu, L. Wang, C. Zhai, and Q. Gu. High-dimensional variance-reduced stochastic gradient
183 expectation-maximization algorithm. In *Proceedings of the 34th International Conference on*
184 *Machine Learning-Volume 70*, pages 4180–4188. JMLR. org, 2017.