
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, $\Theta \subset \mathbb{R}^d$ is the parameters set
5 and $R : \Theta \rightarrow \mathbb{R}$ 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 $\{S(z_i, y_i) \in$
13 $\mathbb{R}^k\}_{i=1}^n$ is the vector of sufficient statistics. We refer the readers to [Efron, 1975] for details on this
14 subclass of problems which is of high interest given the broad range of problems that fall under this
15 assumption.

16 2 Algorithms

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

Algorithm 1 FL-SAEM with Periodic Statistics Averaging

```

1: Input: TO COMPLETE
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 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

```

22 2.1 Challenges with Algorithm 1

23 While Algorithm 1 is a distributed variant of the SAEM, it is neither (a) private nor (b)
 24 communication-efficient.

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

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

30 2.2 Algorithmic solutions

31 **Line 6 – Quantization:** The first step is to quantize the gradient in the Stochastic Langevin Dynam-
 32 ics step used in our sampling scheme Line 6 of Algorithm 1. Inspired by [Alistarh et al., 2017], we
 33 use an extension of the QSGD algorithm for our latent samples. Define the quantization operator as
 34 follows:

$$C_j^{(\ell)}(g, \xi_j) = \|v\| \cdot \text{sign}(g_j) \cdot (\lfloor \ell |g_j| / \|v\| \rfloor + \mathbf{1} \{ \xi_j \leq \ell |g_j| / \|v\| - \lfloor \ell |g_j| / \|v\| \rfloor \}) / \ell \quad (4)$$

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

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

Algorithm 2 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 $\nabla \log p(z_i|y_i; \hat{\theta}_i^{(k)})$:

$$g_i(k, m) = C_j^{(\ell)} \left(\nabla_j f_{\theta_t}(z_i^{(k-1,m)}), \xi_j^{(k)} \right) \quad (5)$$

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

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

$$z_i^{(k,m)} = z_i^{(k-1,m)} + \frac{\gamma_k}{2} g_i(k, m) + \sqrt{\gamma_k} B_k, \quad (6)$$

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 $\hat{\theta}_i^{(t,k)}$
-

38 **Line 7 – Compression MCMC output:** We use the notorious **Top- k** operator that we define as
 39 $\mathcal{C}(x)_i = x_i$, if $i \in \mathcal{S}$; $\mathcal{C}(x)_i = 0$ otherwise and where \mathcal{S} is defined as the size- k set of $i \in [p]$.
 40 Recall that after Line 6 we compute the local statistics $\tilde{S}_i^{(k+1)}$ using the output latent variables from
 41 Algorithm 2. We now use those statistics and compress them using Algorithm 3 as follows:

Algorithm 3 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** :

$$\ddot{S}_i^{(k+1)} = \mathcal{C} \left(\tilde{S}_i^{(k+1)} \right) \quad (7)$$

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

3 FL and Distributed Algorithmic Contributions

43 Final method:

Algorithm 4 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 $\check{S}_i^{(r+1)} = \mathcal{C}(\tilde{S}_i^{(r+1)})$. {Compressed local statistics}
12: Worker send local statistics $\check{S}_i^{(r+1)}$ to server. {Single round of communication}
13: **end for**
14: Server computes **global model**: {(Global) M-Step using aggregated statistics}
$$\hat{\theta}^{(r+1)} = \bar{\theta}(\check{S}^{(r+1)})$$

 where $\check{S}^{(r+1)} = (\check{S}_i^{(r+1)}, i \in D_r)$ and send global model back to the devices.
15: **end for**

44 We can also consider the plain distributed version of the sEM which does not tackle any privacy or
45 communication bottlenecks. It goes as follows:

Algorithm 5 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**

47 **5 Numerical Experiments**

48 **5.1 Nonlinear Mixed Models under Distributed Settings**

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

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

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

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

53 **5.2 Probabilistic Latent Dirichlet Allocation**

54 **5.3 Bi-factor models under the Federated Learning settings**

55 **References**

- 56 D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. Qsgd: Communication-efficient sgd
57 via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*,
58 pages 1709–1720, 2017.
- 59 I. W. P. Consortium. Estimation of the warfarin dose with clinical and pharmacogenetic data. *New*
60 *England Journal of Medicine*, 360(8):753–764, 2009.
- 61 B. Efron. Defining the curvature of a statistical problem (with applications to second order effi-
62 ciency). *The Annals of Statistics*, 3(6):1189–1242, 1975.
- 63 J. Pinheiro and D. Bates. *Mixed-effects models in S and S-PLUS*. Springer Science & Business
64 Media, 2006.