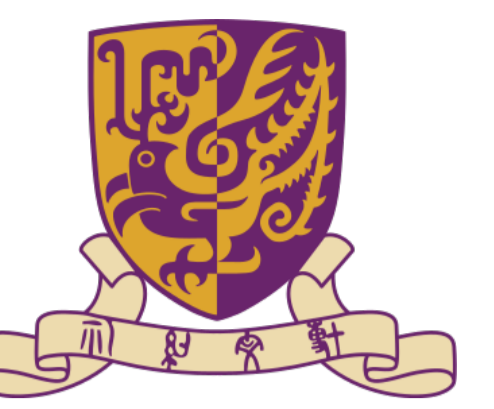
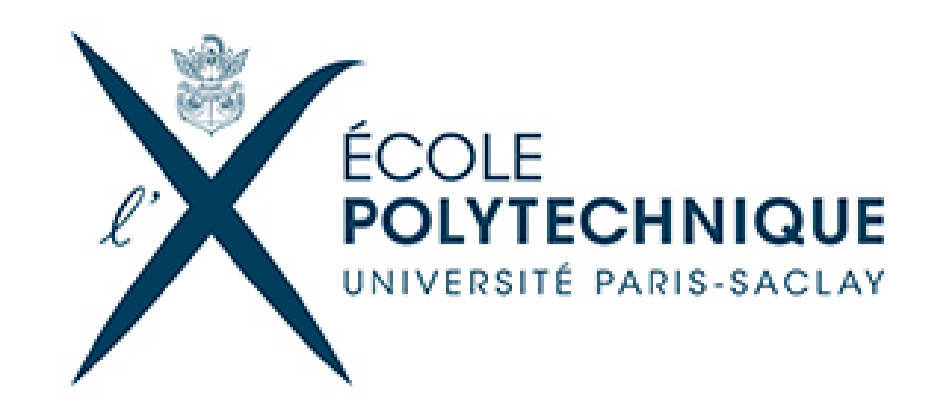


Geom-SPIDER-EM: Faster Variance Reduced Stochastic Expectation Maximization for Nonconvex Finite-Sum Optimization

Gersende Fort¹, Eric Moulines² and Hoi-To Wai³

CNRS, France¹, École Polytechnique, France², Chinese University of Hong Kong, Hong-Kong³



Contribution

- A novel EM algorithm
- Adapted to the finite sum setting
- Stochastic: it combines Stochastic Approximation and variance reduction techniques
- Same complexity as SPIDER-EM (Fort et al, 2020b) – state of the art, among the incremental EM's.

Optimization problem

- Solve on $\Theta \subseteq \mathbb{R}^d$
- $$\operatorname{argmin}_{\theta \in \Theta} - \sum_{i=1}^n \log \int_{\mathcal{Z}} p_i(z; \theta) d\mu(z) + R(\theta), \quad p_i(z; \theta) > 0$$
- Curved exponential family:
- $$- \sum_{i=1}^n \log \int_{\mathcal{Z}} h_i(z) \exp(\langle s_i(z) | \phi(\theta) \rangle) d\mu(z) + R(\theta)$$

- In computational Statistics: minimization of the (penalized) negative likelihood in *latent variable* models.

From EM to incremental EM

- **EM algorithm:** Repeat for $t = 0, \dots$

$$\text{E-step} \quad \bar{s}(\theta_t) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta_t)$$

$$\text{M-step} \quad \theta_{t+1} = T(\bar{s}(\theta_t))$$

where

$$\bar{s}_i(\theta) = \int_{\mathcal{Z}} s_i(z) \frac{p_i(z; \theta)}{\int p_i(u; \theta) d\mu(u)} d\mu(z)$$

$$T(s) = \operatorname{argmin}_{\theta \in \Theta} R(\theta) - \langle s | \phi(\theta) \rangle$$

- EM: an algorithm in the *expectation space* (Delyon et al, 1999)

$$S_{t+1} = \bar{s} \circ T(S_t) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ T(S_t)$$

- EM designed to find the roots of

$$h(s) := \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ T(s) - s = \mathbb{E}[\bar{s}_l(s) - s + V]$$

where $l \sim \mathcal{U}(\{1, \dots, n\})$ and V is a *control variate* i.e. r.v. correlated with \bar{s}_l and centered.

- **Stochastic Approximation** The algorithm

$$\hat{S}_{t+1} = \hat{S}_t + \gamma_{t+1} H_{t+1} \quad \mathbb{E}[H_{t+1} | \text{past}_t] = h(\hat{S}_t)$$

has the **same** limiting set: $\{s : h(s) = 0\}$.

Variance reduced incremental EM

$$\hat{S}_{t+1} = \hat{S}_t + \gamma_{t+1} \left(\frac{1}{b} \sum_{i \in \mathcal{B}_{t+1}} \bar{s}_i \circ T(\hat{S}_t) - \hat{S}_t + V_{t+1} \right)$$

where \mathcal{B}_{t+1} is a mini-batch of examples of size $b \ll n$.

- Online-EM (Neal and Hinton, 1998; Cappé and Moulines, 2009). NO variance reduction ($V_{t+1} = 0$).
- sEM-vr: Stochastic Expectation Maximization with Variance Reduction (Chen et al, 2018)
- FIEM: Fast Increment Expectation Maximization (Karimi et al, 2019; Fort et al, 2020a)
- SPIDER-EM (Fort et al, 2020b) and **Geom-SPIDER-EM**: Stochastic Path Integrated Differential Estimator Expectation Maximization

$$V_{t+1} = V_t + \frac{1}{b} \sum_{i \in \mathcal{B}_t} \bar{s}_i \circ T(\hat{S}_{t-1}) - \frac{1}{b} \sum_{i \in \mathcal{B}_{t+1}} \bar{s}_i \circ T(\hat{S}_{t-1})$$

Geom-SPIDER-EM (Stochastic Path Integrated Differential Estimator)

- 1: $\hat{S}_{1,0} = \hat{S}_{1,-1} = \hat{S}_{\text{init}} \quad S_{1,0} = \bar{s} \circ T(\hat{S}_{1,-1}) + \mathcal{E}_1$
- 2: **for** $t = 1, \dots, k_{\text{out}}$ **do**
- 3: **for** $k = 0, \dots, \xi_t - 1$ **do**
- 4: Sample a mini batch $\mathcal{B}_{t,k+1}$ of size b from $\{1, \dots, n\}$
- 5: $S_{t,k+1} = S_{t,k} + b^{-1} \sum_{i \in \mathcal{B}_{t,k+1}} (\bar{s}_i \circ T(\hat{S}_{t,k}) - \bar{s}_i \circ T(\hat{S}_{t,k-1}))$
- 6: $\hat{S}_{t,k+1} = \hat{S}_{t,k} + \gamma_{t,k+1} (S_{t,k+1} - \hat{S}_{t,k})$
- 7: **end for**
- 8: $\hat{S}_{t+1,-1} = \hat{S}_{t,\xi_t}$
- 9: $S_{t+1,0} = \bar{s} \circ T(\hat{S}_{t+1,-1}) + \mathcal{E}_{t+1} \quad \mathcal{E}_{t+1}$: an error (e.g. part of the data set)
- 10: $\hat{S}_{t+1,0} = \hat{S}_{t+1,-1} + \gamma_{t+1,0} (S_{t+1,0} - \hat{S}_{t+1,-1})$
- 11: **end for**

The **control variate** is refreshed at each *outer loop* # t (see Line 9)

The **length of the outer loop** is a **Geometric** random variable ξ_t (Li et al, 2020; Horvath et al., 2020)

Complexity for ϵ -approximate stationarity

We provide an **explicit** expression of an upper bound for $\mathbb{E}[\|h(\hat{S}_{\tau, \xi_\tau})\|^2]$

- in the non convex setting
- at the end of an outer loop # τ where τ is sampled uniformly in $\{1, \dots, k_{\text{out}}\}$
- as a function of k_{out} , b , n , the learning rate γ ($= \gamma_{t,k}$) and the expectation k_{in} of ξ_t .

With: $k_{\text{in}} = b = O(\sqrt{n})$, $k_{\text{out}} = O(1/(\epsilon k_{\text{in}}))$

Nbr of optimization steps: $O(1/\epsilon)$ Nbr of \bar{s}_i 's evaluations: $\mathcal{K} = O(\sqrt{n} \epsilon^{-1})$

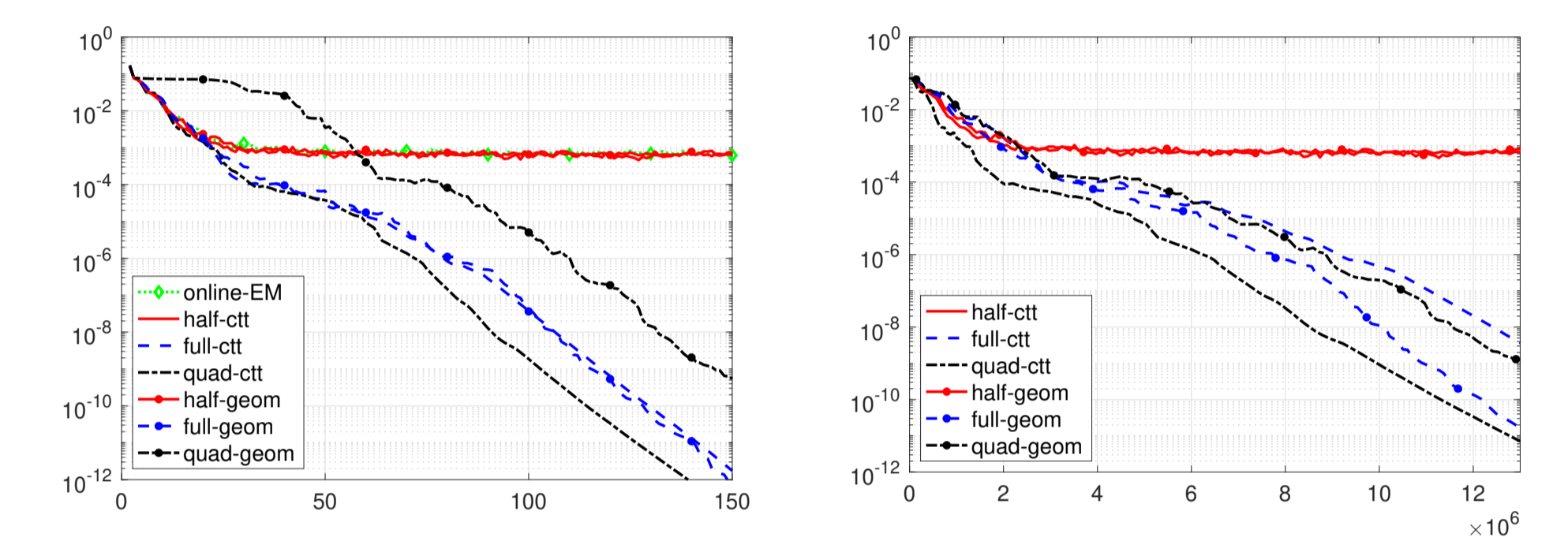
For Online EM: $\mathcal{K} = O(\epsilon^{-2})$ For sEM-vr: $\mathcal{K} = O(n^{2/3} \epsilon^{-1})$

For FIEM: $\mathcal{K} = O(n^{2/3} \epsilon^{-1} \wedge \sqrt{n} \epsilon^{-3/2})$ For SPIDER-EM: $\mathcal{K} = O(\sqrt{n} \epsilon^{-1})$

Inference in Gaussian Mixture Models (from the MNIST data set)

► Gaussian mixture models in \mathbb{R}^{20} ; $G = 12$ components; $n = 6 \cdot 10^4$ examples

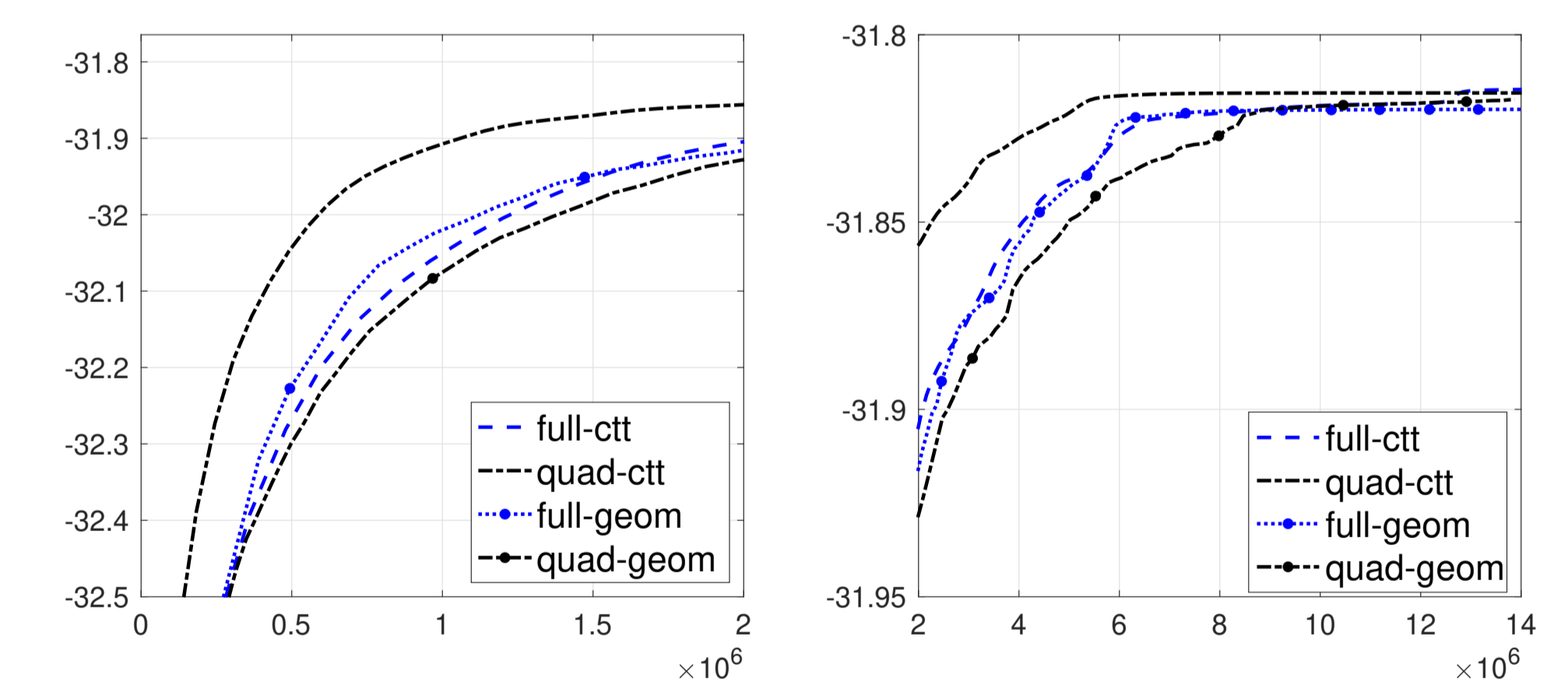
► Displayed: quantile of order 0.5 of $\|h(\hat{S}_{t, \xi_t})\|^2$ vs the number of epochs (left) and vs the number of \bar{s}_i 's evaluations (right)



Length of each outer: either constant (**ctt**) $\xi_t = k_{\text{in}}$, or a geometric r.v. (**geom**) with expectation k_{in}

When refreshing the control variate: use the full data set (**full**), or the half data set (**half**) or a quadratically increasing nbr of examples (**quad**).

► Displayed: evolution of the normalized log-likelihood vs the number of \bar{s}_i 's evaluations until $2e6$ (left) and after (right).



References

- O. Cappé and E. Moulines. On-line expectation-maximization algorithm for latent data models. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 71(3):593–613, 2009.
- J. Chen, J. Zhu, Y.-W. Teh, and T. Zhang. Stochastic expectation maximization with variance reduction. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 7967–7977. 2018.
- B. Delyon, M. Lavielle, and E. Moulines. Convergence of a Stochastic Approximation version of the EM algorithm. *Ann. Statist.*, 27(1):94–128, 1999.
- A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum Likelihood from Incomplete Data via the EM Algorithm. *J. Roy. Stat. Soc. B Met.*, 39(1):1–38, 1977.
- G. Fort, E. Moulines, and P. Gach. Fast Incremental Expectation Maximization for finite-sum optimization: nonasymptotic convergence. Technical report, HAL 02617725v2, 2020a.
- G. Fort, E. Moulines, and H.T. Wai. A Stochastic Path-Integrated Differential Estimator Expectation Maximization Algorithm. In *Advances in Neural Information Processing Systems 34*, Curran Associates, Inc., 2020b.
- S. Horvath, L. Lei, P. Richtarik, and M.I. Jordan. Adaptivity of Stochastic Gradient Methods for Nonconvex Optimization. Technical report, arXiv 2002.05359, 2020.
- B. Karimi, H.-T. Wai, E. Moulines, and M. Lavielle. On the Global Convergence of (Fast) Incremental Expectation Maximization Methods. In *NeurIPS 32*. 2019.
- Z. Li, H. Bao, X. Zhang, and P. Richtarik. PAGE: A Simple and Optimal Probabilistic Gradient Estimator for Nonconvex Optimization. Technical report, arXiv 2008.10898, 2020.
- R. M. Neal and G. E. Hinton. *A View of the EM Algorithm that Justifies Incremental, Sparse, and other Variants*, pages 355–368. Springer Netherlands, Dordrecht, 1998.