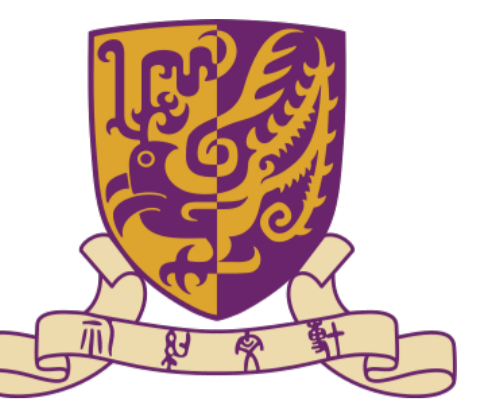
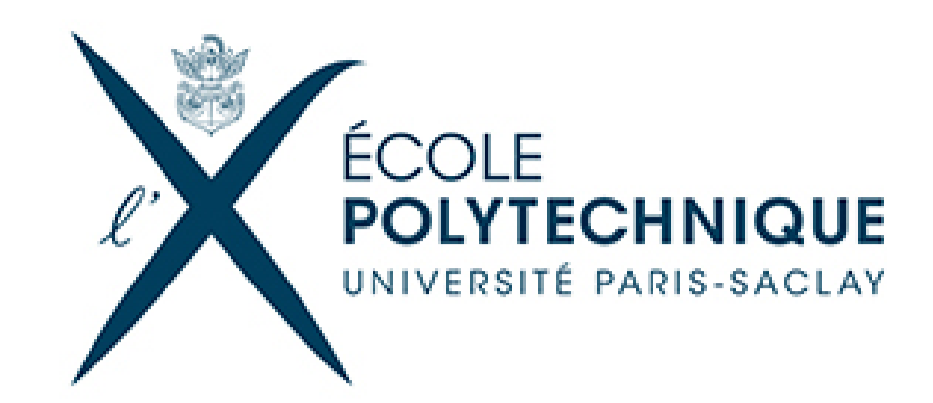


A Stochastic Path-Integrated Differential Estimator EM Algorithm

Gersende Fort¹, Eric Moulines² and Hoi-To Wai³

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



(smooth) regularized Empirical Risk Minimization

- Given a set of n observations $\{y_i, i \in [n]\}$.
- Goal: solve the possibly non-convex problem

$$\arg \min_{\theta} F, \quad F(\theta) := -\frac{1}{n} \sum_{i=1}^n \log g(y_i, \theta) + R(\theta)$$

when $\Theta \subseteq \mathbb{R}^d$ and intractable positive function g

$$g(y_i, \theta) = \int_{\mathcal{Z}} f(z, y_i; \theta) \mu(dz)$$

Setting: Maximum Likelihood Estimation in latent variable models from exponential family

- parametric model, $\theta \mapsto g(y_i, \theta)$ is the likelihood of y_i ,
- latent variable $z \in \mathcal{Z}$; $f(z, y_i; \theta)$ is complete likelihood.
- Focus on the **Exponential Family Distribution**:

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

- Tool: **Expectation Maximization (EM)** (Dempster et al., 1977) algorithm which takes advantage of the latent structure.

EM for Exponential Family

★ E-operation:

$$\bar{s}(\theta) := \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta), \quad \bar{s}_i(\theta) := \int_{\mathcal{Z}} S(z, y_i) \rho(z|y_i; \theta) \mu(dz)$$

where $z \mapsto \rho(z|y_i; \theta) \propto f(y_i, z; \theta)$ is the a posteriori distribution of the latent variable z .

★ M-operation: for any s ,

$$T(s) := \arg \min_{\theta \in \Theta} \{R(\theta) + \psi(\theta) - \langle s | \phi(\theta) \rangle\}$$

Batch EM Algorithm: given $\hat{s}^{(0)}$, repeat for $k \geq 0$,

- (M-step) $\tau^{(k+1)} = T(\hat{s}^{(k)})$
- (E-step) $\hat{s}^{(k+1)} = \bar{s}(\tau^{(k+1)}) = \bar{s} \circ T(\hat{s}^{(k)})$

- Fixed points:

$$\{s : h(s) = 0\}, \quad h(s) := \bar{s} \circ T(s) - s$$

For large n , the E-step is computationally expensive!

Variance reduced Stochastic Approximation within EM

Idea: Take advantage of the **Stochastic Approximation (SA)** context

$$h(s) = 0 \iff \mathbb{E}[\bar{s}_i \circ T(s) - s] = 0, \quad l \text{ uniform on } [n]$$

combined with **Variance Reduction (vr)**

$$h(s) = 0 \iff \mathbb{E}[\bar{s}_i \circ T(s) - s + V] = 0, \text{ with } \mathbb{E}[V] = 0.$$

vr-SA within EM (gal form):

- constant stepsize γ , mini-batch $\mathcal{B}_{t,k+1}$ in $[n]$,
- possibly nested loops (outer $\#t$, inner $\#k$),

$$\hat{s}^{(t,k+1)} = \hat{s}^{(t,k)} + \gamma (\bar{s}_{\mathcal{B}_{t,k+1}} \circ T(\hat{s}^{(t,k)}) - \hat{s}^{(t,k)} + \mathcal{S}^{(t,k+1)})$$

★ **SPIDER-EM. Init:** $\mathcal{S}^{(t,0)} = \bar{s} \circ T(\hat{s}^{(t-1,k_{in}-1)})$

$$\mathcal{S}^{(t,k+1)} = \mathcal{S}^{(t,k)} - \bar{s}_{\mathcal{B}_{t,k+1}} \circ T(\hat{s}^{(t,k-1)})$$

inspired by Fang et al., 2018; see also Nguyen et al. 2017, Wang et al. 2019

• **online-EM:** (Cappé & Moulines, 2009) unique outer, $\mathcal{S}^{(k+1)} = 0$.

• **sEM-VR:** (Chen et al., 2018) Init: $\mathcal{S}^{(t,0)} = \bar{s} \circ T(\hat{s}^{(t-1,k_{in})})$

$$\mathcal{S}^{(t,k+1)} = \mathcal{S}^{(t,0)} - \bar{s}_{\mathcal{B}_{t,k+1}} \circ T(\hat{s}^{(t-1,k_{in})})$$

• **fiEM:** (Karimi et al., 2019) unique outer, a memory $S_i^{(k)}$ for $i \in [n]$,

$$\mathcal{S}^{(k+1)} = \bar{\mathcal{S}}^{(k+1)} - S_{\mathcal{B}_{k+1}}^{(k)}, \quad \bar{\mathcal{S}}^{(k+1)} = \bar{\mathcal{S}}^{(k)} + n^{-1} \sum_{i \in \mathcal{B}_{k+1}'} (S_i^{(k+1)} - S_i^{(k)}), \quad S_i^{(k+1)} = S_i^{(k)} \mathbb{1}_{i \notin \mathcal{B}_{k+1}'} + \bar{s}_i \circ T(\hat{s}^k) \mathbb{1}_{i \in \mathcal{B}_{k+1}'}$$

SPIDER-EM: state of the art complexity bounds

• **Assumptions.** (i) $T(s)$ exists and is unique; (ii) ϕ, ψ, R and F are C^1 ; (iii) $\nabla(\phi \circ T)(s)$ is a symmetric matrix with eigenvalues uniformly bounded in $[v_{\min}, v_{\max}]$ such that $v_{\min} > 0$; (iv) for all i , $\bar{s}_i \circ T$ is globally Lipschitz with constant L_i ; (v) the function $s \mapsto \nabla(F \circ T)(s)$ is globally Lipschitz with constant L_W

• **Property.** The **Lyapunov** function $W := F \circ T$ has a globally Lipschitz gradient and $\|\nabla W(s)\|^2 \leq v_{\max}^2 \|h(s)\|^2$.

• **Upper bounds for convergence.** Minibatch size: b .

Theorem Set $L^2 := n^{-1} \sum_i L_i^2$. Choose $\gamma := \alpha/L$ where $\alpha \in (0, v_{\min}/\mu_*(k_{in}, b))$ and $\mu_*(k_{in}, b) := v_{\max} \sqrt{k_{in}/b} + L_W/(2L)$.

$$\mathbb{E}[\|h(\hat{s}^{\tau, \xi-1})\|^2] \leq \left(\frac{1}{k_{in}} + \frac{\alpha^2}{b}\right) \frac{2L}{\alpha\{v_{\min} - \alpha\mu_*(k_{in}, b)\}} \frac{1}{k_{out}} (\mathbb{E}[W(\hat{s}^{init})] - \min W)$$

• **Complexity bounds.** \mathcal{K}_{CE} : total nbr of conditional number of per-sample conditional expectations \bar{s}_i evaluations; \mathcal{K}_{opt} : total nbr of optimization steps $T(\hat{s})$.

SPIDER-EM

• Find b , k_{in} and k_{out} as a function of n, ϵ s.t.

$$\mathbb{E}[\|h(\hat{s}^{\tau, \xi-1})\|^2] \leq \epsilon$$

• Optimal strategy:

$$b = O(\sqrt{n}), \quad k_{in} = O(\sqrt{n}), \quad k_{out} = O(\epsilon^{-1}/\sqrt{n})$$

Comparison

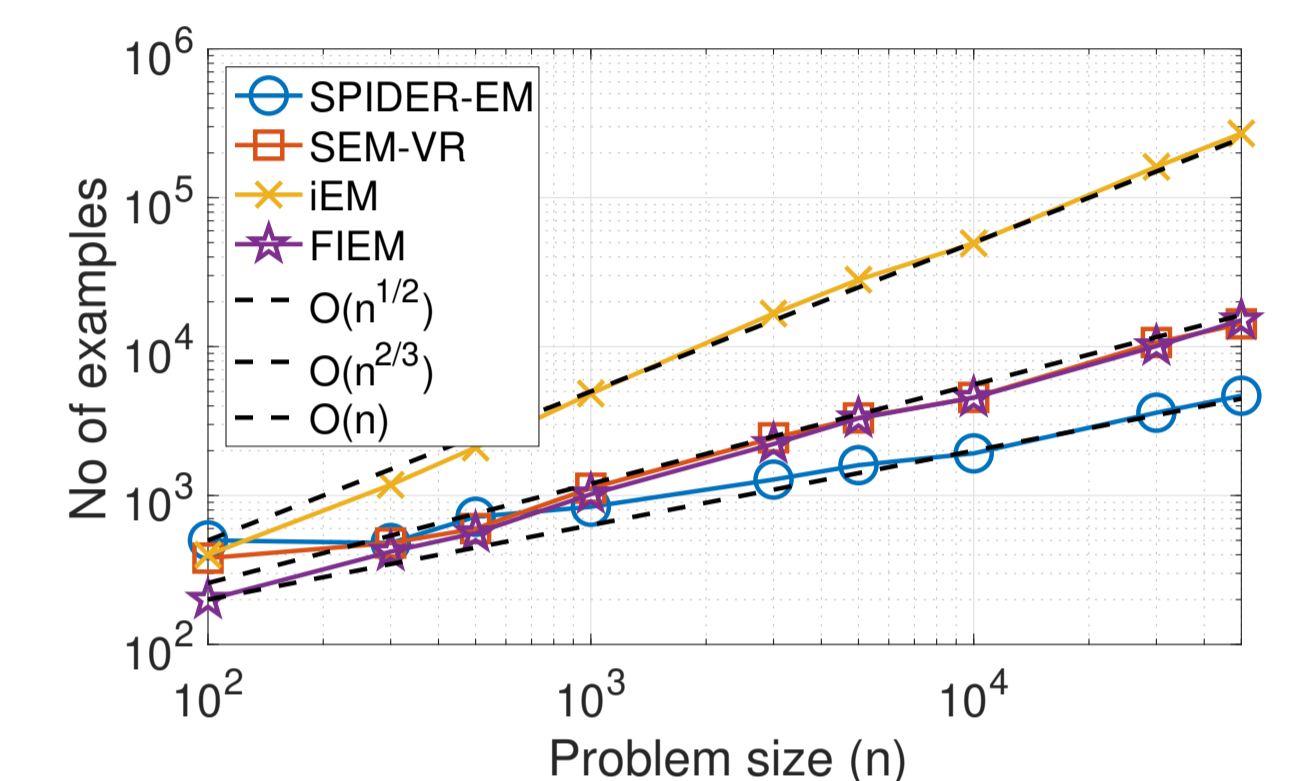
	\mathcal{K}_{CE}	\mathcal{K}_{opt}	Optimal \mathcal{K}_{CE}
EM	$n + nk_{max}$	$1 + k_{max}$	—
online-EM	$n + bk_{max}$	$1 + k_{max}$	ϵ^{-2}
fiEM	$n + 2bk_{max}$	$1 + k_{max}$	$\epsilon^{-1}n^{2/3}$
sEM-VR	$n + 2nk_{out} + bk_{in}k_{out}$	$1 + k_{in}k_{out}$	$\epsilon^{-1}n^{2/3}$
SPIDER-EM	$n + nk_{out} + 2bk_{in}k_{out}$	$1 + k_{in}k_{out}$	$\epsilon^{-1}\sqrt{n}$

$k_{max} = k_{in}k_{out}$ for algorithms with a single outer loop

Fitting a Gaussian Mixture Model (GMM)

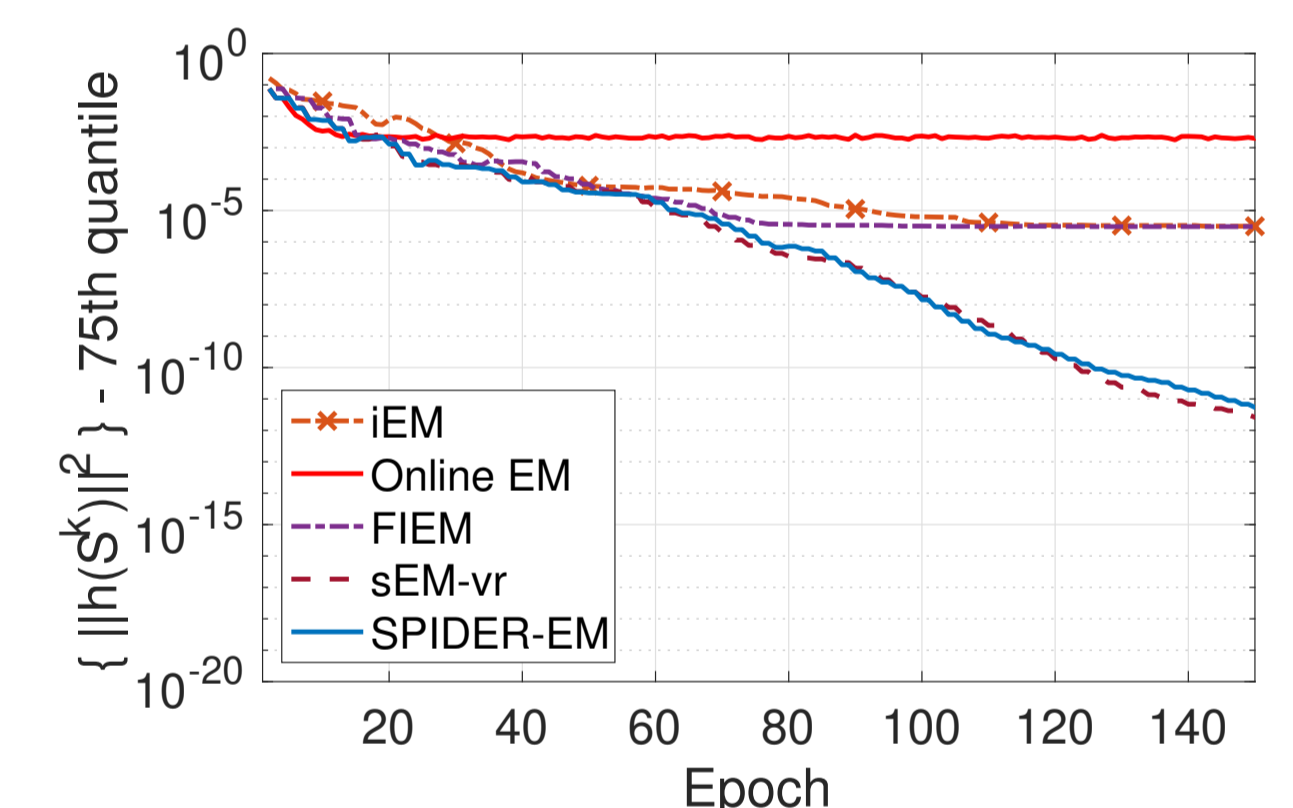
► **Synthetic data** Efficiency against the problem size n . Dataset is a GMM in \mathbb{R} with 2 components. θ collects the means.

- Median estimated nbr $\mathcal{K}_{CE} - n$ (estimated over 50 runs) needed to reach an accuracy of $2.5e - 5$.

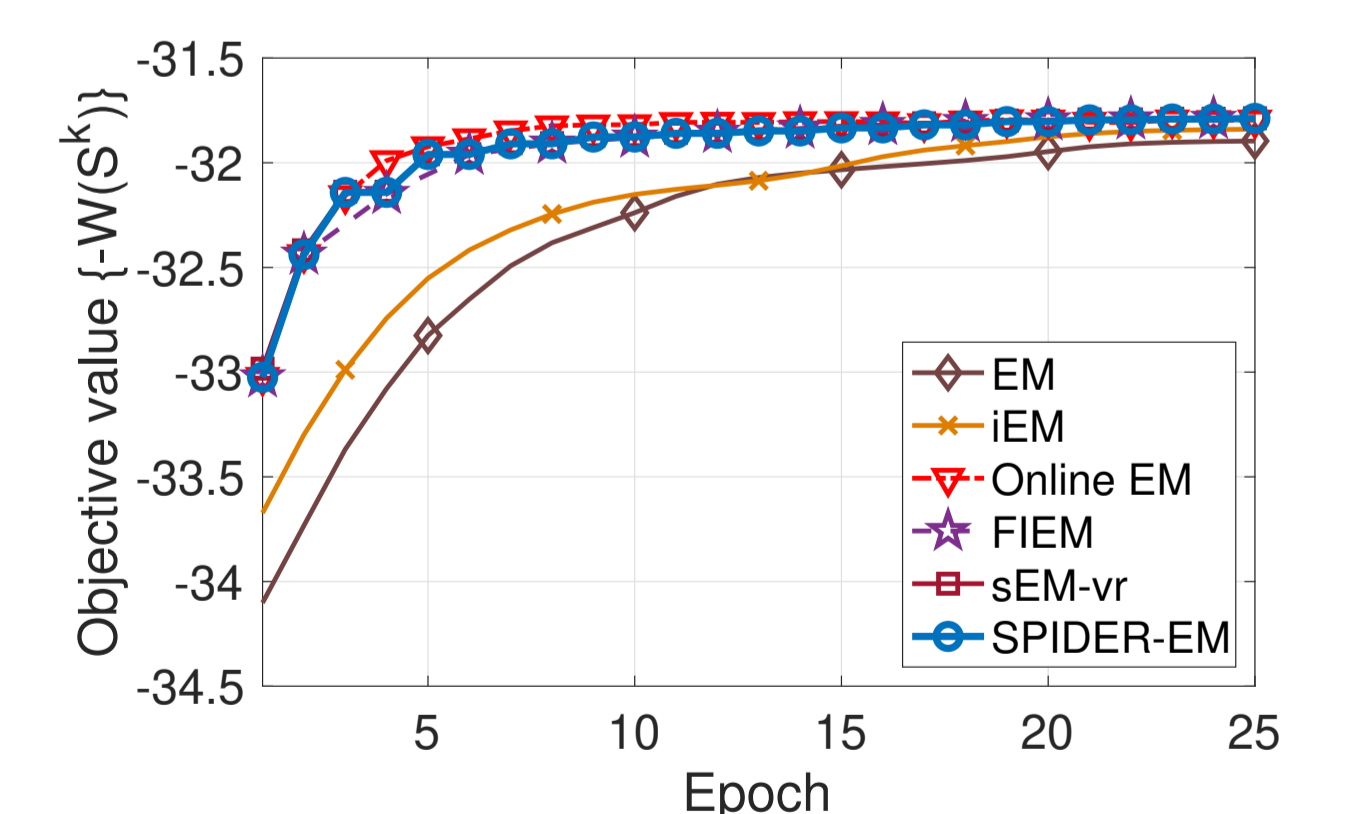


► **MNIST dataset** $n = 6e4$ images classified into 12 classes. θ : parameters of a GMM in \mathbb{R}^{20} (784 pixels reduced to 20 features through PCA)

- Quantile 0.75 of the distribution of $\|h(\hat{s}_{t,-1})\|^2$



• Evolution of the objective fct W



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.

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.

C. Fang, C.J. Li, Z. Lin, and T. Zhang. SPIDER: Near-Optimal Non-Convex Optimization via Stochastic Path-Integrated Differential Estimator. In *NeurIPS 31*, 2018.

B. Karimi, H.-T. Wai, E. Moulines, and M. Lavielle. On the Global Convergence of (Fast) Incremental Expectation Maximization Methods. In *NeurIPS 32*, 2019.

L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *ICML*, 2017.

Z. Wang, K. Ji, Y. Zhou, Y. Liang, and V. Tarokh. SpiderBoost and Momentum: Faster Stochastic Variance Reduction Algorithms. In *NeurIPS 32*, 2019.