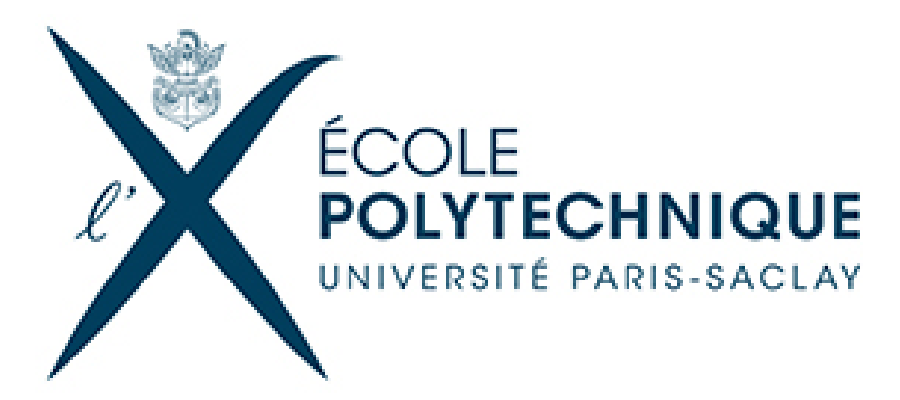
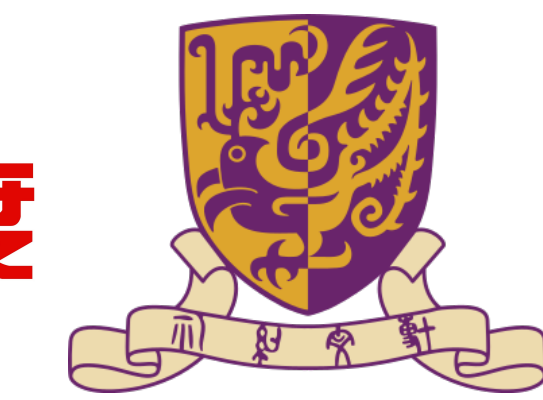


# Minimization by Incremental Stochastic Surrogate Optimization for Large Scale Nonconvex Problems

Belhal Karimi<sup>1</sup>, Hoi-To Wai<sup>2</sup>, Eric Moulines<sup>3</sup> and Ping Li<sup>1</sup>

Baidu Research<sup>1</sup>, Chinese University of Hong Kong<sup>2</sup>, Ecole Polytechnique<sup>3</sup>

belhalkarimi@baidu.com, htwai@se.cuhk.edu.hk, eric.moulines@polytechnique.edu, liping11@baidu.com@gmail.com



## Large Scale Optimization

- **Objective:** Constrained minimization problem of a finite sum of functions:

$$\min_{\theta \in \Theta} \mathcal{L}(\theta) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta), \quad (1)$$

where  $\mathcal{L}_i : \mathbb{R}^p \rightarrow \mathbb{R}$  is bounded from below and is (possibly) nonconvex and include a nonsmooth penalty.

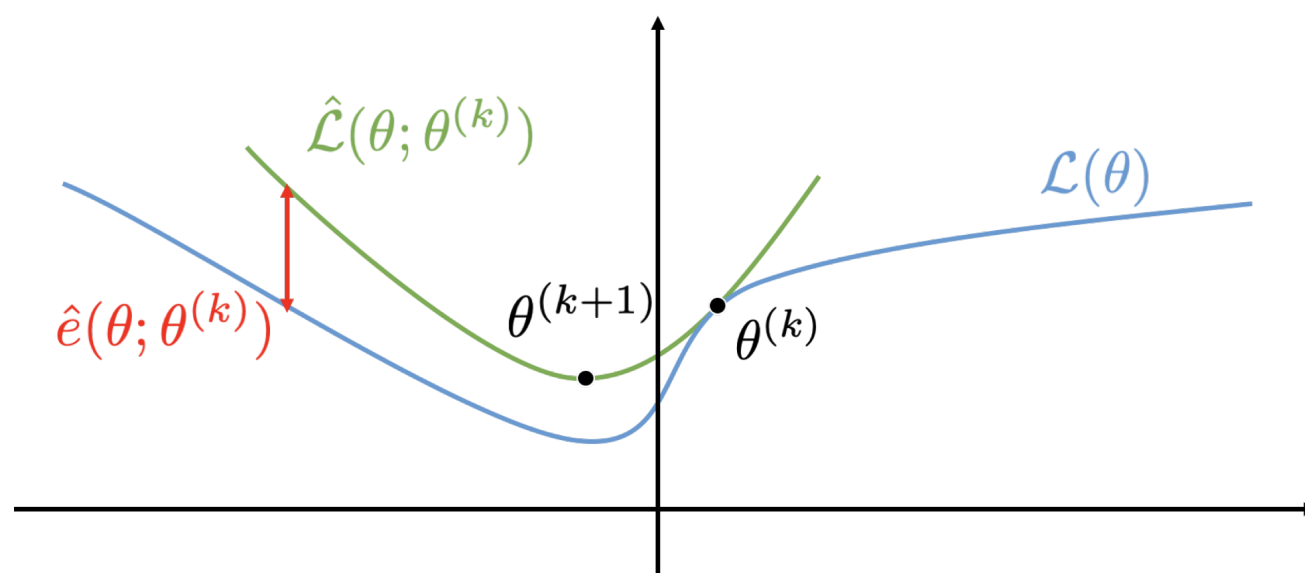
- The gap  $\hat{e}(\theta; \{\theta_i\}_{i=1}^n)$  plays a key role in the convergence analysis and we require this error to be L-smooth for some constant  $L > 0$ . Denote by  $\langle \cdot | \cdot \rangle$  the scalar product, we also introduce the following stationary point condition:

**Definition 1.** (Asymptotic Stationary Point Condition)

A sequence  $(\theta^k)_{k \geq 0}$  satisfies the asymptotic stationary point condition if

$$f'(\theta, d) := \lim_{t \rightarrow 0^+} \frac{f(\theta + td) - f(\theta)}{t} \geq 0. \quad (2)$$

## Majorization-Minimization Scheme



**Algorithm 2** The MISO method (Mairal, 2015).

- 1: **Input:** initialization  $\theta^{(0)}$ .
- 2: Initialize the surrogate function as  $\mathcal{A}_i^0(\theta) := \tilde{\mathcal{L}}_i(\theta; \theta^{(0)})$ ,  $i \in \llbracket 1, n \rrbracket$ .
- 3: **for**  $k = 0, 1, \dots, K_{\max}$  **do**
- 4: Pick  $i_k$  uniformly from  $\llbracket 1, n \rrbracket$ .
- 5: Update  $\mathcal{A}_i^{k+1}(\theta)$  as:

$$\mathcal{A}_i^{k+1}(\theta) = \begin{cases} \tilde{\mathcal{L}}_i(\theta; \theta^{(k)}), & \text{if } i = i_k \\ \mathcal{A}_i^k(\theta), & \text{otherwise.} \end{cases}$$

- 6: Set  $\theta^{(k+1)} \in \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^{k+1}(\theta)$ .
- 7: **end for**

- **MISO Method:** fix any  $n \geq 1$ , we stop the SA at a random iteration  $N$  with

## An Intractability for Latent Data Models

- Case when the surrogate functions computed in Algorithm ?? are not tractable.
- Assume that the surrogate can be expressed as an integral over a set of latent variables  $z = (z_i \in \mathcal{Z}, i \in [n]) \in \mathcal{Z}^n$ .

$$\tilde{\mathcal{L}}_i(\theta; \bar{\theta}) := \int_{\mathcal{Z}} r_i(\theta; \bar{\theta}, z_i) p_i(z_i; \bar{\theta}) \mu_i(dz_i) \quad \forall (\theta, \bar{\theta}) \in \Theta \times \Theta. \quad (3)$$

- Our scheme is based on the computation, at each iteration, of stochastic auxiliary functions for a mini-batch of components. For  $i \in [n]$ , the auxiliary function, noted  $\tilde{\mathcal{L}}_i(\theta; \bar{\theta}, \{z_m\}_{m=1}^M)$  is a Monte Carlo approximation of the surrogate function  $\tilde{\mathcal{L}}_i(\theta; \bar{\theta})$  defined by (3) such that:

$$\tilde{\mathcal{L}}_i(\theta; \bar{\theta}, \{z_m\}_{m=1}^M) := \frac{1}{M} \sum_{m=1}^M r_i(\theta; \bar{\theta}, z_m), \quad (4)$$

where  $\{z_m\}_{m=1}^M$  is a Monte Carlo batch.

## MISSO Method

- dd

**Algorithm 2** The MISSO method.

- 1: **Input:** initialization  $\theta^{(0)}$ ; a sequence of non-negative numbers  $\{M_{(k)}\}_{k=0}^\infty$ .
- 2: For all  $i \in \llbracket 1, n \rrbracket$ , draw  $M_{(0)}$  Monte Carlo samples with the stationary distribution  $p_i(\cdot; \theta^{(0)})$ .
- 3: Initialize the surrogate function as

$$\tilde{\mathcal{A}}_i^0(\theta) := \tilde{\mathcal{L}}_i(\theta; \theta^{(0)}, \{z_{i,m}^{(0)}\}_{m=1}^{M_{(0)}}), \quad i \in \llbracket 1, n \rrbracket.$$

- 4: **for**  $k = 0, 1, \dots, K_{\max}$  **do**
- 5: Pick a function index  $i_k$  uniformly on  $\llbracket 1, n \rrbracket$ .
- 6: Draw  $M_{(k)}$  Monte Carlo samples with the stationary distribution  $p_{i_k}(\cdot; \theta^{(k)})$ .
- 7: Update the individual surrogate functions recursively as:

$$\tilde{\mathcal{A}}_i^{k+1}(\theta) = \begin{cases} \tilde{\mathcal{L}}_i(\theta; \theta^{(k)}, \{z_{i,m}^{(k)}\}_{m=1}^{M_{(k)}}), & \text{if } i = i_k \\ \tilde{\mathcal{A}}_i^k(\theta), & \text{otherwise.} \end{cases}$$

- 8: Set  $\theta^{(k+1)} \in \arg \min_{\theta \in \Theta} \tilde{\mathcal{L}}^{(k+1)}(\theta) := \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{A}}_i^{k+1}(\theta)$ .
- 9: **end for**

## Global Convergence Analysis

**Assumptions:** we need a few regularity conditions in this case,

1. There exists a Borel measurable function  $\hat{H} : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{H}$ ,

$$\hat{H}_\eta(x) - P_\eta \hat{H}_\eta(x) = H_\eta(x) - h(\eta), \quad \forall \eta \in \mathcal{H}, x \in \mathcal{X}.$$

$\implies$  existence of solution to the **Poisson equation**.

2. For all  $\eta \in \mathcal{H}$  and  $x \in \mathcal{X}$ ,  $\|\hat{H}_\eta(x)\| \leq L_{PH}^{(0)}$ ,  $\|P_\eta \hat{H}_\eta(x)\| \leq L_{PH}^{(0)}$ , and

$$\sup_{x \in \mathcal{X}} \|P_\eta \hat{H}_\eta(x) - P_{\eta'} \hat{H}_{\eta'}(x)\| \leq L_{PH}^{(1)} \|\eta - \eta'\|, \quad \forall (\eta, \eta') \in \mathcal{H}^2.$$

$\implies$  **smoothness** of  $\hat{H}_\eta(x)$ , satisfied if  $P_\eta, H_\eta(x)$  are smooth w.r.t.  $\eta$ .

3. It holds that  $\sup_{\eta \in \mathcal{H}, x \in \mathcal{X}} \|H_\eta(x) - h(\eta)\| \leq \sigma$ .

$\implies$  requires the noise is **uniformly bounded** for all  $x \in \mathcal{X}$ .

**Example:** assumptions 1 & 2 are satisfied if the Markov kernel  $P_{\eta_n}$  is geometrically ergodic + smooth, and the drift term is smooth w.r.t.  $\eta$ .

**Theorem 1.** Suppose that the step sizes are decreasing and  $\gamma_1 \leq 0.5(c_1(L + C_h))^{-1}$  (+other conditions). Let  $V_{0,n} := \mathbb{E}[V(\eta_0) - V(\eta_{n+1})]$ ,

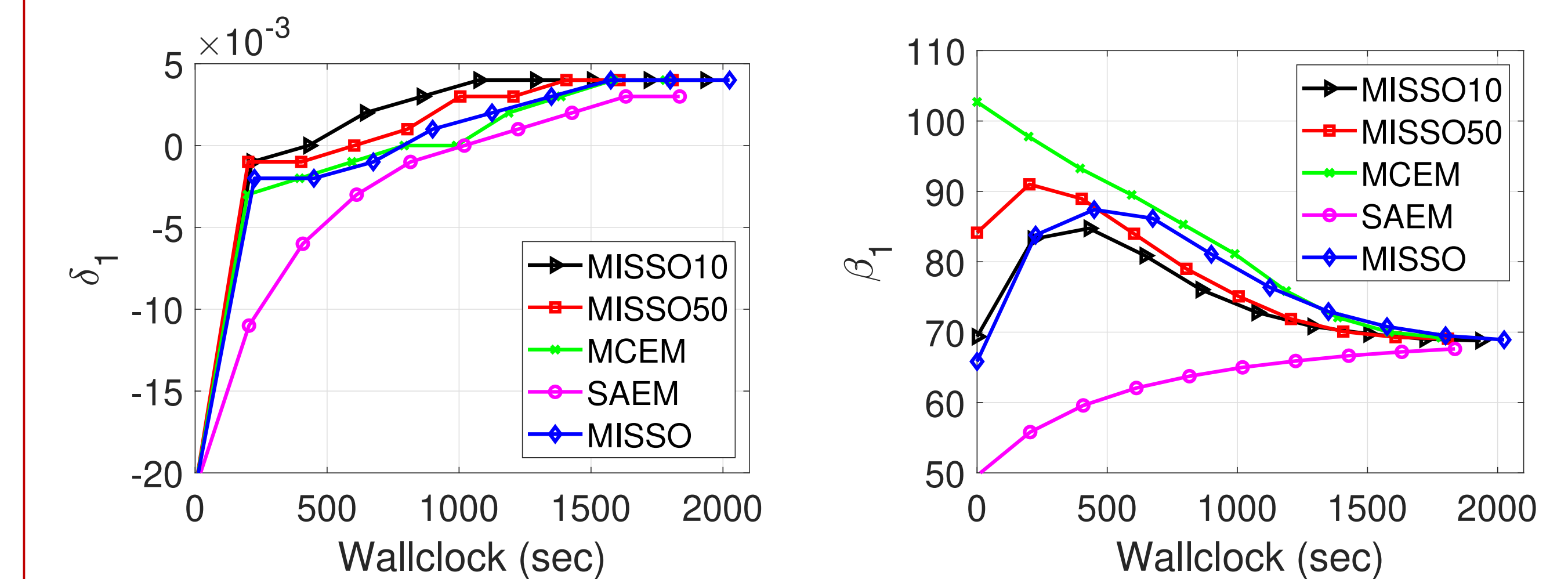
$$\mathbb{E}[\|h(\eta_N)\|^2] \leq \frac{2c_1(V_{0,n} + C_{0,n} + (\sigma^2 L + C_\gamma) \sum_{k=0}^n \gamma_{k+1}^2)}{\sum_{k=0}^n \gamma_{k+1}} + 2C_0.$$

- Set  $\gamma_k = (2c_1 L(1 + C_h) \sqrt{k})^{-1} \implies \mathbb{E}[\|h(\eta_N)\|^2] = \mathcal{O}(c_0 + \log n / \sqrt{n})$  (same as Case 1).

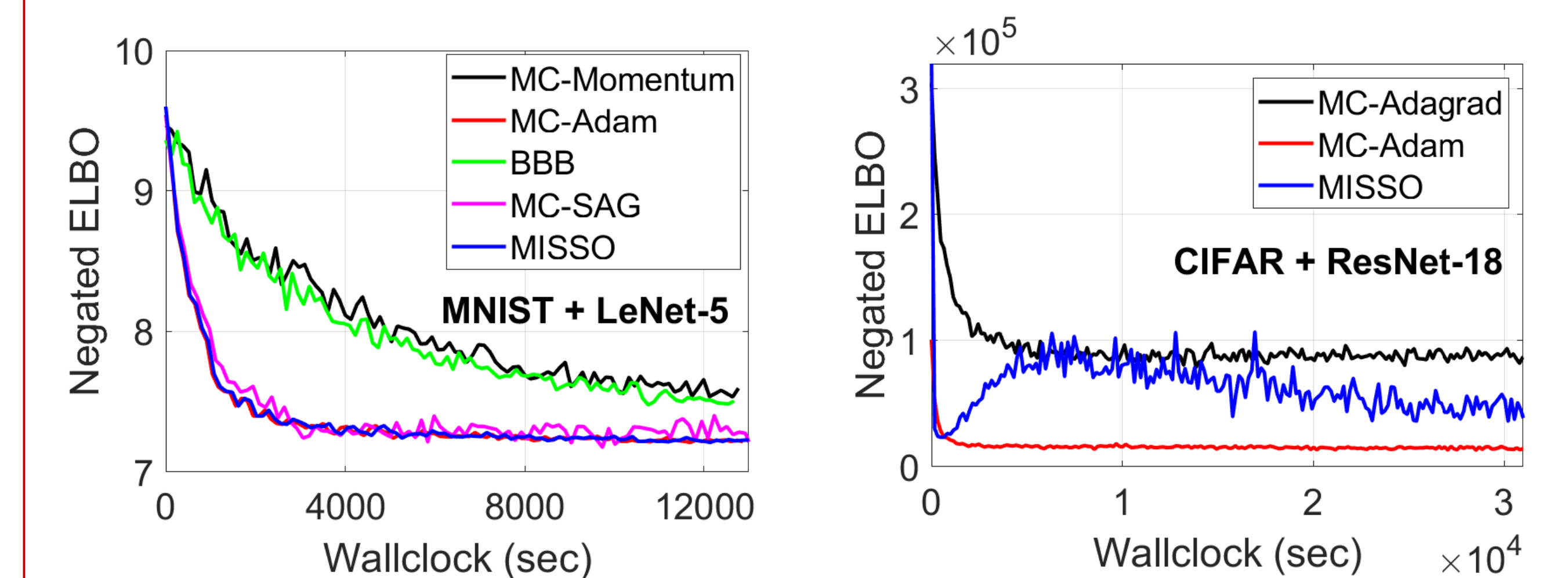
- **Proof idea:** challenge is that  $e_{n+1}$  is not zero-mean  $\implies$  bound the sum of  $\mathbb{E}[\langle \nabla V(\eta_n) | e_{n+1} \rangle]$  w/ Poisson equation + a novel decomposition (cf. [Lemma 2](#)).

## Numerical Experiments

- Logistic Regression on Traumabase dataset (severe hemorrhage):



- Bayesian variants of LeNet-5 and ResNet-18 on MNIST and CIFAR10:



## Conclusion

- [Theorem 1 & 2](#) show the non-asymptotic convergence rate of biased SA scheme with smooth (possibly non-convex) Lyapunov function.
- With appropriate step size, in  $n$  iterations the SA scheme finds  $\mathbb{E}[\|h(\eta_N)\|^2] = \mathcal{O}(c_0 + \log n / \sqrt{n})$ , where  $c_0$  is the bias and  $h(\cdot)$  is the mean field.
- Applications to online EM and online policy gradient.

## References



