# Minimization by Incremental Stochastic Surrogate Optimization for Large Scale Nonconvex Problems
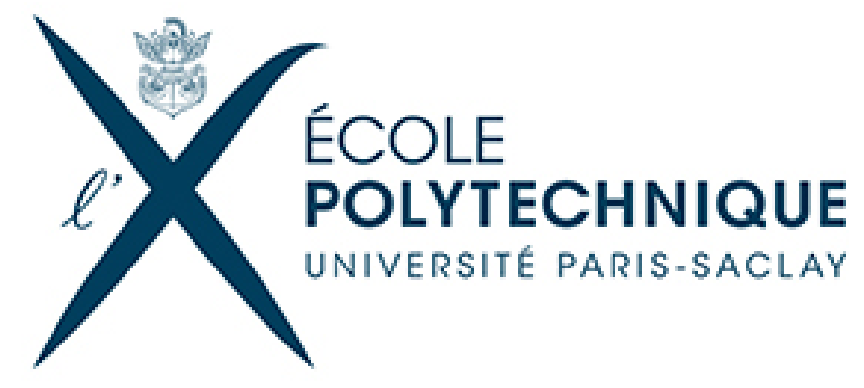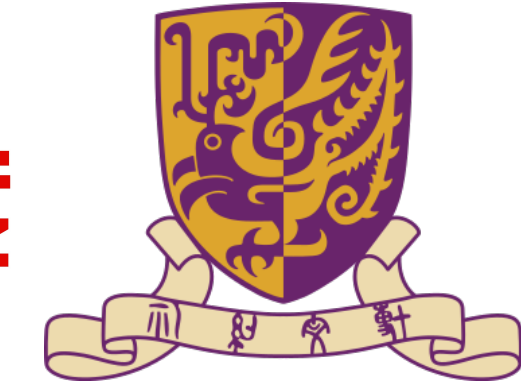
**Belhal Karimi[1], Hoi-To Wai[2], Eric Moulines[3] and Ping Li[1]**

Baidu Research[1], Chinese University of Hong Kong[2], Ecole Polytechnique[3]

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:

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

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

- The gap $\widehat{e}(\theta; \{\overline{\theta}_i\}_{i=1}^n)$ is L-smooth. Denote by $\langle \cdot \mid \cdot \rangle$ the scalar product, the stationary point condition is:

**Definition 1.** (Asymptotic Stationary Point Condition)
A sequence $(\theta^k)_{k \geq 0}$ satisfies the asymptotic stationary point condition if

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

## Majorization-Minimization Scheme

- The MISO method (Mairal, 2015)

**Algorithm 2** The MISO method (Mairal, 2015).
1: **Input:** initialization $\theta^{(0)}$.
2: Initialize the surrogate function as $\mathcal{A}_i^0(\theta) := \widehat{\mathcal{L}}_i(\theta; \theta^{(0)})$, $i \in [\![1, n]\!]$.
3: **for** $k = 0, 1, ..., K_{\max}$ **do**
4: Pick $i_k$ uniformly from $[\![1, n]\!]$.
5: Update $\mathcal{A}_i^{k+1}(\theta)$ as:

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

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

## An Inctractability 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 $z = (z_i \in \mathsf{Z}, i \in [n]) \in \mathsf{Z}[\![]\!]$.

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

- 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 $\widetilde{\mathcal{L}}_i(\theta; \overline{\theta}, \{z_m\}_{m=1}^M)$ is a Monte Carlo approximation of the surrogate function $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ defined by (3Doc-Start) such that:

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

where $\{z_i^m\}_{m=0}^{M-1}$ is a Monte Carlo batch.

## MISSO Method

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

$$\widetilde{\mathcal{A}}_i^0(\theta) := \widetilde{\mathcal{L}}_i(\theta; \theta^{(0)}, \{z_{i,m}^{(0)}\}_{m=1}^{M_{(0)}}), \ i \in [\![1, n]\!] .$$

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

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

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

## Global Convergence Analysis

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

**H 1.** For all $i \in [n]$ and $\overline{\theta} \in \Theta$, $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ is convex w.r.t. $\theta$, and it holds $\widehat{\mathcal{L}}_i(\theta; \overline{\theta}) \geq \mathcal{L}_i(\theta)$, $\forall \theta \in \Theta$ where the equality holds when $\theta = \overline{\theta}$.

**H2.** For any $\overline{\theta}_i \in \Theta$, $i \in [n]$ and some $\epsilon > 0$, the difference function $\widehat{e}(\theta; \{\overline{\theta}_i\}_{i=1}^n) := \frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}_i(\theta; \overline{\theta}_i) - \mathcal{L}(\theta)$ is defined for all $\theta \in \Theta_\epsilon$ and differentiable for all $\theta \in \Theta_\epsilon$, where $\Theta_\epsilon = \{\theta \in \mathbb{R}^d, \inf_{\theta' \in \Theta} \|\theta - \theta'\| < \epsilon\}$ is an $\epsilon$-neighborhood set of $\Theta$. Moreover, for some constant $L$, the gradient satisfies $\|\nabla \widehat{e}(\theta; \{\overline{\theta}_i\}_{i=1}^n)\|^2 \leq 2L \, \widehat{e}(\theta; \{\overline{\theta}_i\}_{i=1}^n)$, $\forall \theta \in \Theta$.

**H3.** For all $i \in [n]$, $\overline{\theta} \in \Theta$, $z_i \in \mathsf{Z}$, $r_i(\cdot; \overline{\theta}, z_i)$ is convex on $\Theta$ and is lower bounded.

**H4.** For the samples $\{z_{i,m}\}_{m=1}^M$, there exist finite constants $C_r$ and $C_{gr}$ such that for all $i \in [n]$,

$$C_r := \sup_{\overline{\theta} \in \Theta} \sup_{M > 0} \frac{1}{\sqrt{M}} \mathbb{E}_{\overline{\theta}} \left[ \sup_{\theta \in \Theta} \left| \sum_{m=1}^M \left\{ r_i(\theta; \overline{\theta}, z_{i,m}) - \widehat{\mathcal{L}}_i(\theta; \overline{\theta}) \right\} \right| \right]$$

$$C_{gr} := \sup_{\overline{\theta} \in \Theta} \sup_{M > 0} \sqrt{M} \mathbb{E}_{\overline{\theta}} \left[ \sup_{\theta \in \Theta} \frac{\left| \frac{1}{M} \sum_{m=1}^M \widehat{\mathcal{L}}'_i(\theta, \theta - \overline{\theta}; \overline{\theta}) - r'_i(\theta, \theta - \overline{\theta}; \overline{\theta}, z_{i,m}) \right|^2}{\|\overline{\theta} - \theta\|} \right]$$

where we denoted by $\mathbb{E}_{\overline{\theta}}[\cdot]$ the expectation w.r.t. a Markov chain $\{z_{i,m}\}_{m=1}^M$ with initial distribution $\xi_i(\cdot; \overline{\theta})$, transition kernel $\Pi_{i,\overline{\theta}}$, and stationary distribution $p_i(\cdot; \overline{\theta})$.

**Theorem 1** *Under H1-H4. For any $K_{\max} \in \mathbb{N}$, let $K$ be an independent discrete r.v. drawn uniformly from $\{0, ..., K_{\max} - 1\}$ and define the following quantity:*

$$\Delta_{(K_{\max})} := 2nL \mathbb{E}[\widetilde{\mathcal{L}}^{(0)}(\theta^{(0)}) - \widetilde{\mathcal{L}}^{(K_{\max})}(\theta^{(K_{\max})})] + 4LC_r \overline{M}_{(K_{\max})} .$$

*Then we have following non-asymptotic bounds:*

$$\mathbb{E}[\|\nabla \widehat{e}^{(K)}(\theta^{(K)})\|^2] \leq \frac{\Delta_{(K_{\max})}}{K_{\max}} \ and \ \mathbb{E}[g_-(\theta^{(K)})] \leq \sqrt{\frac{\Delta_{(K_{\max})}}{K_{\max}}} + \frac{C_{gr}}{K_{\max}} \overline{M}_{(K_{\max})} . \qquad (16)$$
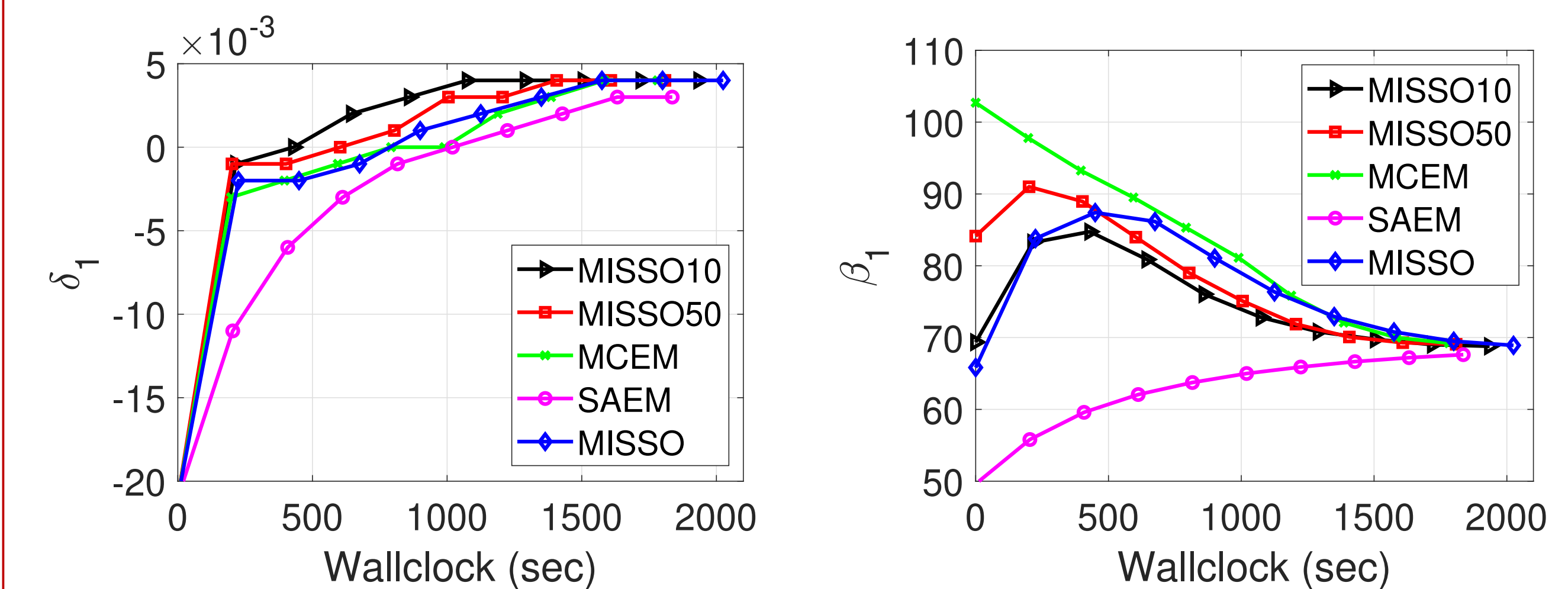
## Numerical Experiments

- Logistic Regression **with missing values** on Traumabase (severe hemorrhage):

$$p_i(y_i | z_i) = S(\delta^\top \overline{z}_i)^{y_i} \left(1 - S(\delta^\top \overline{z}_i)\right)^{1 - y_i} ,$$
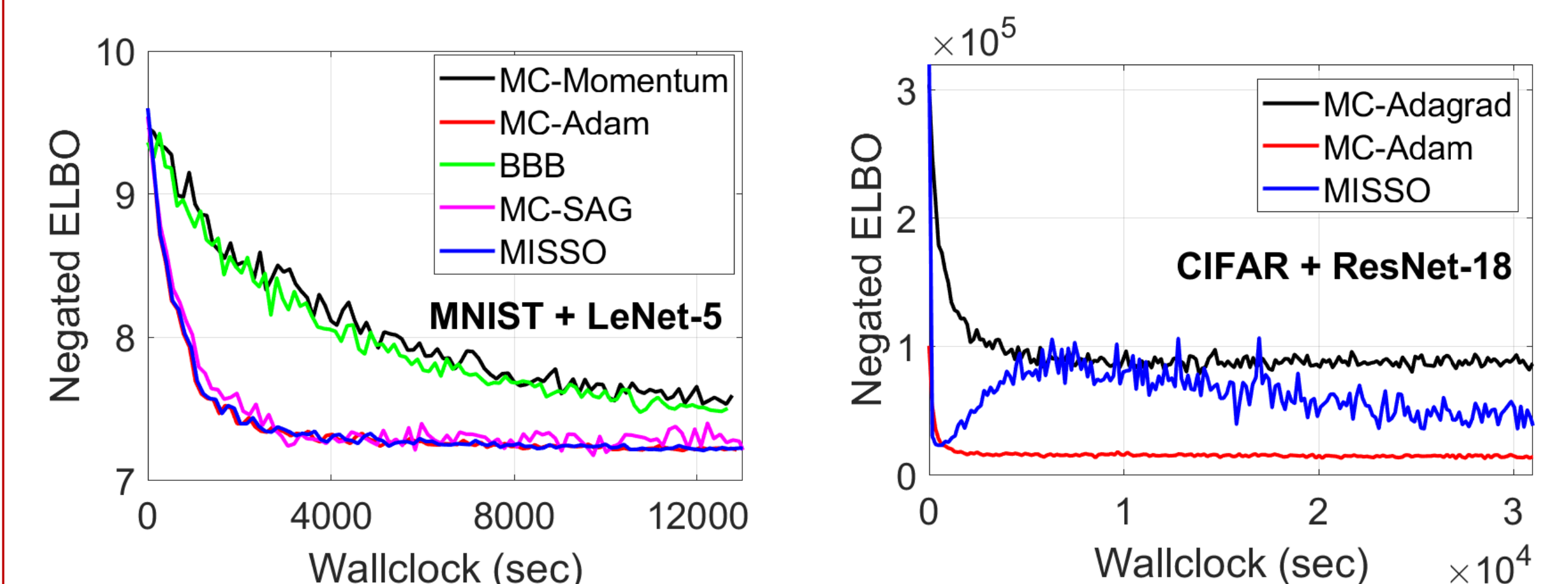
- 16 quantitative measurements, like BMI, age, blood pressure, heart rate at different stages after the accident on 6384 patients
- **MISSO is an incremental MCEM in this case**



- Bayesian variants of LeNet-5 and ResNet-18 on MNIST and CIFAR10:
- Variational inference and the ELBO loss to fit Bayesian Neural Networks on different datasets.
- **MISSO is an incremental VI in this case**



## 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

Julien Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.