
A Doubly Stochastic Surrogate Optimization Scheme for Non-convex Finite-sum Problems

Belhal Karimi

Center of Applied Mathematics
Ecole Polytechnique
Palaiseau, FR
belhal.karimi@polytechnique.fr

Hoi-To Wai

Department of SEEM
The Chinese University of Hong Kong
Hong Kong
htwai@se.cuhk.edu.hk

Eric Moulines

Center of Applied Mathematics
Ecole Polytechnique
Palaiseau, FR
eric.moulines@polytechnique.fr

Abstract

Many constrained, non-convex optimization problems can be tackled using the Majorization-Minimization (MM) method which alternates between constructing a surrogate function which upper bounds the objective function, and then minimizing this surrogate. For problems which minimize a finite sum of functions, a stochastic version of the MM method selects a batch of functions at random at each iteration and optimizes the accumulated surrogate. However, in many cases of interest such as variational inference for latent variable models, the surrogate functions are expressed as an expectation. In this contribution, we propose a doubly stochastic MM method based on Monte Carlo approximation of these stochastic surrogates. We establish asymptotic and non-asymptotic convergence of our scheme in a constrained, non-convex, non-smooth optimization setting. We apply our new framework for inference of logistic regression model with missing covariates and for variational inference of autoencoder on the MNIST dataset.

1 Introduction

We consider the *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 Θ is a convex, compact, and closed subset of \mathbb{R}^p , and for any $i \in \llbracket 1, n \rrbracket$, the function $\mathcal{L}_i : \mathbb{R}^p \rightarrow \mathbb{R}$ is bounded from below and is (possibly) non-convex and non-smooth.

To tackle the optimization problem (1), a popular approach is to apply the majorization-minimization (MM) method which iteratively minimizes a majorizing surrogate function. A large number of existing procedures fall into this general framework, for instance gradient-based or proximal methods or the Expectation-Maximization (EM) algorithm [McLachlan and Krishnan, 2008] and some variational Bayes inference techniques [Jordan et al., 1999]; see for example [Razaviyayn et al., 2013] and [Lange, 2016] and the references therein. When the number of terms n in (1) is large, the vanilla MM method may be intractable because it requires to construct a surrogate function for all the n terms \mathcal{L}_i at each iteration. Here, a remedy is to apply the Minimization by Incremental Surrogate Optimization (MISO) method proposed by Mairal [2015], where the surrogate functions are

updated incrementally. The MISO method can be interpreted as a combination of MM and ideas which have emerged for variance reduction in stochastic gradient methods [Schmidt et al., 2017].

The success of the MISO method rests upon the efficient minimization of surrogates such as convex functions, see [Mairal, 2015, Section 2.3]. In many applications of interest, the natural surrogate functions are intractable, yet they are defined as expectation of tractable functions. This for example the case for inference in latent variable models. Another application is variational inference, [Ghahramani, 2015], in which the goal is to approximate the posterior distribution of parameters given the observations; see for example [Neal, 2012, Blundell et al., 2015, Polson et al., 2017, Rezende et al., 2014, Li and Gal, 2017].

This paper fills the gap in the literature by proposing a new method called *Minimization by Incremental Stochastic Surrogate Optimization (MISSO)* which is designed for the finite sum optimization with a finite-time convergence guarantee. Our contributions can be summarized as follows.

- We propose a unifying framework of analysis for incremental stochastic surrogate optimization when the surrogates are defined by expectations of tractable functions. The proposed MISSO method is built on the Monte Carlo integration of the intractable surrogate function, *i.e.*, a doubly stochastic surrogate optimization scheme. In addition, we present an incremental variational inference and Monte-Carlo EM methods as two special cases of this framework.
- We establish both asymptotic and non-asymptotic convergence for the MISSO method. In particular, the MISSO method converges almost surely to a stationary point and in $\mathcal{O}(n/\epsilon)$ iterations to an ϵ -stationary point.

In Section 2, we review the techniques for incremental minimization of finite sum functions based on the MM principle; specifically, we review the MISO method as introduced in [Mairal, 2015], and present a class of surrogate functions expressed as an expectation over a latent space. The MISSO method is then introduced for the latter class of surrogate functions. In Section 3, we provide the asymptotic and non-asymptotic convergence analysis for the MISSO method. Finally, Section 4 presents numerical applications to illustrate our findings including parameter inference for logistic regression with missing covariates and variational inference for Bayesian neural network.

Notations We denote $\llbracket 1, n \rrbracket = \{1, \dots, n\}$. Unless otherwise specified, $\|\cdot\|$ denotes the standard Euclidean norm and $\langle \cdot | \cdot \rangle$ is the inner product in Euclidean space. For any function $f : \Theta \rightarrow \mathbb{R}$, $f'(\theta, d)$ is the directional derivative of f at θ along the direction d , *i.e.*,

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

The directional derivative is assumed to exist for the functions introduced throughout this paper.

2 Incremental Minimization of Finite Sum Non-convex Functions

The objective function in (1) is composed of a finite sum of possibly non-smooth and non-convex functions. A popular approach here is to apply the MM method. The MM method tackles (1) through alternating between two steps — (i) minimizing a *surrogate* function which upper bounds the original objective function; and (ii) updating the surrogate function to tighten the upper bound.

As mentioned in the Introduction, the MISO method proposed by Mairal [2015] is developed as an iterative scheme that only updates the surrogate functions *partially* at each iteration. Formally, for any $i \in \llbracket 1, n \rrbracket$, we consider a surrogate function $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ which satisfies

S1. For all $i \in \llbracket 1, n \rrbracket$ and $\bar{\theta} \in \Theta$, the function $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ is convex w.r.t. θ , and it holds

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) \geq \mathcal{L}_i(\theta), \quad \forall \theta \in \Theta, \quad (3)$$

where the equality holds when $\theta = \bar{\theta}$.

S2. For any $\bar{\theta}_i \in \Theta$, $i \in \llbracket 1, n \rrbracket$ and some $\epsilon > 0$, the difference function $\hat{\epsilon}(\theta; \{\bar{\theta}_i\}_{i=1}^n) := \frac{1}{n} \sum_{i=1}^n \hat{\mathcal{L}}_i(\theta; \bar{\theta}_i) - \mathcal{L}(\theta)$ is defined for all $\theta \in \Theta_\epsilon$ and differentiable for all $\theta \in \Theta$, where

$\Theta_\epsilon = \{\theta \in \mathbb{R}^d, \inf_{\theta' \in \Theta} \|\theta - \theta'\| < \epsilon\}$ is an ϵ -neighborhood set of Θ . Moreover, for some constant L , the gradient satisfies

$$\|\nabla \hat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n)\|^2 \leq 2L \hat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n), \quad \forall \theta \in \Theta. \quad (4)$$

S1 is a common condition used for surrogate optimization, see [Mairal, 2015, Section 2.3]. Meanwhile, S2 can be satisfied when the difference function $\hat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n)$ is L -smooth for all $\theta \in \mathbb{R}^d$, where the condition can be implied through applying [Razaviyayn et al., 2013, Proposition 1].

The inequality (3) implies $\hat{\mathcal{L}}_i(\theta; \bar{\theta}) \geq \mathcal{L}_i(\theta) > -\infty$ for any $\theta \in \Theta$. The MISO method is an incremental version of the MM method, as summarized by Algorithm 1. As seen in the pseudo code, the MISO method maintains an iteratively updated set of surrogate upper-bound functions $\{\mathcal{A}_i^k(\theta)\}_{i=1}^n$ and updates the iterate through minimizing the average of the surrogate functions.

Particularly, only one out of the n surrogate functions is updated at each iteration [cf. Line 5] and the sum function $\frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^{k+1}(\theta)$ is designed to be ‘easy to optimize’, for example, it can be a sum of quadratic functions. As such, the MISO method is suitable for large-scale optimization as the computation cost per iteration is independent of n . Moreover, under S1, S2, it was shown that the MISO method converges almost surely to a stationary point of (1) [Mairal, 2015, Proposition 3.1].

Algorithm 1 MISO method [Mairal, 2015]

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

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

- 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**

We now consider the case when the surrogate functions $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ are intractable. Let Z be a measurable set, $p_i : Z \times \Theta \rightarrow \mathbb{R}_+$ be a pdf, $r_i : \Theta \times \Theta \times Z \rightarrow \mathbb{R}$ be a measurable function and μ_i be a σ -finite measure, we consider surrogate functions which satisfy S1, S2 that can be expressed as an expectation:

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) := \int_Z r_i(\theta; \bar{\theta}, z_i) p_i(z_i; \bar{\theta}) \mu_i(dz_i) \quad \forall (\theta, \bar{\theta}) \in \Theta \times \Theta. \quad (5)$$

Plugging (5) into the MISO method is not feasible since the update step in Step 6 involves a minimization of an expectation. Several motivating examples of (1) are given in Section 2.

We propose the *Minimization by Incremental Stochastic Surrogate Optimization* (MISSO) method which replaces the expectation in (5) by *Monte Carlo* integration and then optimizes (1) incrementally. Denote by $M \in \mathbb{N}$ the Monte Carlo batch size and let $z_m \in Z$, $m = 1, \dots, M$ be a set of samples. These samples can be drawn (Case 1) i.i.d. from the distribution $p_i(\cdot; \bar{\theta})$ or (Case 2) from a Markov chain with the stationary distribution $p_i(\cdot; \bar{\theta})$; see Section 3 for illustrations. To this end, we define

$$\tilde{\mathcal{L}}_i(\theta; \bar{\theta}, \{z_m\}_{m=1}^M) := \frac{1}{M} \sum_{m=1}^M r_i(\theta; \bar{\theta}, z_m) \quad (6)$$

and we summarize the proposed MISSO method in Algorithm 2. As seen, the method is similar to the MISO method but it involves two types of randomness. The first randomness comes from the selection of i_k in Line 5. The second randomness is that a set of Monte-Carlo approximated functions $\tilde{\mathcal{A}}_i^k(\theta)$ is used in lieu of $\mathcal{A}_i^k(\theta)$ when optimizing for the next iterate $\theta^{(k)}$. We now discuss two applications of the MISSO method.

Example 1: Maximum Likelihood Estimation for Latent Variable Model Latent variable models [Bishop, 2006] are constructed by introducing unobserved (latent) variables which help explain the observed data. We consider n independent observations $((y_i, z_i), i \in \llbracket n \rrbracket)$ where y_i is observed and z_i is latent. In this incomplete data framework, define $\{f_i(z_i, \theta), \theta \in \Theta\}$ to be the complete data likelihood models, i.e., joint likelihood of the observations and latent variables. Let

$$g_i(\theta) := \int_Z f_i(z_i, \theta) \mu_i(dz_i), \quad i \in \llbracket 1, n \rrbracket \quad (9)$$

Algorithm 2 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

$$\tilde{\mathcal{A}}_i^0(\theta) := \tilde{\mathcal{L}}_i(\theta; \theta^{(0)}, \{z_{i,m}^{(0)}\}_{m=1}^{M_{(0)}}), \quad i \in \llbracket 1, n \rrbracket. \quad (7)$$

- 4: **for** $k = 0, 1, \dots$ **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:

$$\tilde{\mathcal{A}}_i^{k+1}(\theta) = \begin{cases} \tilde{\mathcal{L}}_i(\theta; \theta^{(k)}, \{z_{i,m}^{(k)}\}_{m=1}^{M_{(k)}}), & \text{if } i = i_k \\ \tilde{\mathcal{A}}_i^k(\theta), & \text{otherwise.} \end{cases} \quad (8)$$

- 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**
-

denote the incomplete data likelihood, *i.e.*, the marginal likelihood of the observations. For ease of notations, the dependence on the observations is made implicit. The maximum likelihood (ML) estimation problem takes $\mathcal{L}_i(\theta)$ to be the i th negated incomplete data log-likelihood $\mathcal{L}_i(\theta) := -\log g_i(\theta)$.

Assume without loss of generality that $g_i(\theta) \neq 0$ for all $\theta \in \Theta$, we define by $p_i(z_i, \theta) := f_i(z_i, \theta)/g_i(\theta)$ the conditional distribution of the latent variable z_i given the observation y_i . A surrogate function $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ satisfying S1 can be obtained through writing $f_i(z_i, \theta) = \frac{f_i(z_i, \theta)}{p_i(z_i, \bar{\theta})} p_i(z_i, \bar{\theta})$ and applying the Jensen inequality:

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) = \int_{\mathcal{Z}} \underbrace{\log(p_i(z_i, \bar{\theta})/f_i(z_i, \theta))}_{=r_i(\theta; \bar{\theta}, z_i)} p_i(z_i, \bar{\theta}) \mu_i(dz_i), \quad (10)$$

We note that S2 can also be verified for common distribution models. We can apply the MISSO method following the above specification of $r_i(\theta; \bar{\theta}, z_i), p_i(z_i, \bar{\theta})$.

Example 2: Variational Inference Let $((x_i, y_i), i \in \llbracket 1, n \rrbracket)$ be i.i.d. input-output pairs and $w \in \mathcal{W} \subseteq \mathbb{R}^d$ be a latent variable. When conditioned on the input $x = (x_i, i \in \llbracket 1, n \rrbracket)$, the joint distribution of $y = (y_i, i \in \llbracket 1, n \rrbracket)$ and w is given by:

$$p(y, w|x) = \pi(w) \prod_{i=1}^n p(y_i|x_i, w). \quad (11)$$

Our goal is to compute the posterior distribution $p(w|y, x)$. In most cases, the posterior distribution $p(w|y, x)$ is intractable and is approximated using a family of parametric distributions, $\{q(w, \theta), \theta \in \Theta\}$. The variational inference (VI) problem [Blei et al., 2017] boils down to minimizing the KL divergence between $q(w, \theta)$ and the posterior distribution $p(w|y, x)$, as follows:

$$\min_{\theta \in \Theta} \mathcal{L}(\theta) := \text{KL}(q(w; \theta) || p(w|y, x)) := \mathbb{E}_{q(w; \theta)} [\log(q(w; \theta)/p(w|y, x))] . \quad (12)$$

Using (11), we decompose $\mathcal{L}(\theta) = n^{-1} \sum_{i=1}^n \mathcal{L}_i(\theta) + \text{const.}$ where:

$$\mathcal{L}_i(\theta) := -\mathbb{E}_{q(w; \theta)} [\log p(y_i|x_i, w)] + \frac{1}{n} \mathbb{E}_{q(w; \theta)} [\log q(w; \theta)/\pi(w)] = r_i(\theta) + d(\theta) . \quad (13)$$

Directly optimizing the finite sum objective function in (12) can be difficult. First, with $n \gg 1$, evaluating the objective function $\mathcal{L}(\theta)$ requires a full pass over the entire dataset. Second, for some complex models, the expectations in (13) can be intractable even if we assume a simple parametric model for $q(w; \theta)$. Assume that \mathcal{L}_i is L-smooth, *i.e.*, \mathcal{L}_i is differentiable on Θ and its gradient $\nabla \mathcal{L}_i$ is L-Lipschitz. We apply the MISSO method with a quadratic surrogate function defined as:

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) := \mathcal{L}_i(\bar{\theta}) + \langle \nabla_{\theta} \mathcal{L}_i(\bar{\theta}) | \theta - \bar{\theta} \rangle + \frac{L}{2} \|\bar{\theta} - \theta\|^2 . \quad (14)$$

It is easily checked that $\hat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}})$ satisfies S1, S2. To compute the gradient $\nabla \mathcal{L}_i(\bar{\boldsymbol{\theta}})$, we apply the re-parametrization technique suggested in [Paisley et al., 2012, Kingma and Welling, 2014, Blundell et al., 2015]. Let $t : \mathbb{R}^d \times \Theta \mapsto \mathbb{R}^d$ be a differentiable function w.r.t. $\boldsymbol{\theta} \in \Theta$ which is designed such that the law of $w = t(z, \boldsymbol{\theta})$, where $z \sim \mathcal{N}_d(0, \mathbf{I})$, is $q(\cdot, \bar{\boldsymbol{\theta}})$. By [Blundell et al., 2015, Proposition 1], the gradient of $-r_i(\cdot)$ in (13) is:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(w; \bar{\boldsymbol{\theta}})} [\log p(y_i | x_i, w)] = \mathbb{E}_{z \sim \mathcal{N}_d(0, \mathbf{I})} [\mathbf{J}_{\boldsymbol{\theta}}^t(z, \bar{\boldsymbol{\theta}}) \nabla_w \log p(y_i | x_i, w) \big|_{w=t(z, \bar{\boldsymbol{\theta}})}], \quad (15)$$

where for each $z \in \mathbb{R}^d$, $\mathbf{J}_{\boldsymbol{\theta}}^t(z, \bar{\boldsymbol{\theta}})$ is the Jacobian of the function $t(z, \cdot)$ with respect to $\boldsymbol{\theta}$ evaluated at $\bar{\boldsymbol{\theta}}$. In addition, for most cases, the term $\nabla d(\bar{\boldsymbol{\theta}})$ can be evaluated in closed form.

$$r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z) := \left\langle \nabla_{\boldsymbol{\theta}} d(\bar{\boldsymbol{\theta}}) - \mathbf{J}_{\boldsymbol{\theta}}^t(z, \bar{\boldsymbol{\theta}}) \nabla_w \log p(y_i | x_i, w) \big|_{w=t(z, \bar{\boldsymbol{\theta}})} \mid \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \right\rangle + \frac{L}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|^2. \quad (16)$$

Finally, using (14) and (16), the surrogate function (6) is given by $\tilde{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, \{z_m\}_{m=1}^M) := M^{-1} \sum_{m=1}^M r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z_m)$ where $\{z_m\}_{m=1}^M$ is an i.i.d sample from $\mathcal{N}(0, \mathbf{I})$.

3 Convergence Analysis

We provide non-asymptotic convergence bound for the MISSO method and show that it converges asymptotically to a stationary point. Consider the following assumptions.

H1. For all $i \in \llbracket 1, n \rrbracket$, $\bar{\boldsymbol{\theta}} \in \Theta$, $z_i \in \mathbb{Z}$, the measurable function $r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z_i)$ is convex in $\boldsymbol{\theta}$ and is lower bounded.

H2. For the samples $\{z_{i,m}\}_{m=1}^M$, there exists finite constants C_r and C_{gr} such that

$$C_r := \sup_{\bar{\boldsymbol{\theta}} \in \Theta} \sup_{M > 0} \frac{1}{\sqrt{M}} \mathbb{E}_{\bar{\boldsymbol{\theta}}} \left[\sup_{\boldsymbol{\theta} \in \Theta} \left| \sum_{m=1}^M \left\{ r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z_{i,m}) - \hat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}) \right\} \right| \right] \quad (17)$$

$$C_{gr} := \sup_{\bar{\boldsymbol{\theta}} \in \Theta} \sup_{M > 0} \sqrt{M} \mathbb{E}_{\bar{\boldsymbol{\theta}}} \left[\sup_{\boldsymbol{\theta} \in \Theta} \left| \frac{1}{M} \sum_{m=1}^M \frac{\hat{\mathcal{L}}'_i(\boldsymbol{\theta}, \boldsymbol{\theta} - \bar{\boldsymbol{\theta}}; \bar{\boldsymbol{\theta}}) - r'_i(\boldsymbol{\theta}, \boldsymbol{\theta} - \bar{\boldsymbol{\theta}}; \bar{\boldsymbol{\theta}}, z_{i,m})}{\|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|} \right|^2 \right] \quad (18)$$

for all $i \in \llbracket 1, n \rrbracket$, and we denoted by $\mathbb{E}_{\bar{\boldsymbol{\theta}}}[\cdot]$ the expectation w.r.t. a Markov chain $\{z_{i,m}\}_{m=1}^M$ with initial distribution $\xi_i(\cdot; \bar{\boldsymbol{\theta}})$, transition kernel $P_{i, \bar{\boldsymbol{\theta}}}$, and stationary distribution $p_i(\cdot; \bar{\boldsymbol{\theta}})$.

H2 essentially requires to control the expectation of the supremum of an empirical process [Shapiro et al., 2009, Boucheron et al., 2013]. In particular, if $M \rightarrow \infty$, the surrogate function's value and its directional derivative approximate that of $\hat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}})$ uniformly for all $\boldsymbol{\theta} \in \Theta$. As discussed before, there are two relevant cases here:

Case 1: When the samples $\{z_m\}_{m=1}^M$ used to construct the approximation $\tilde{\mathcal{L}}_i(\cdot; \cdot)$ are drawn i.i.d. directly from $p_i(\cdot; \bar{\boldsymbol{\theta}})$ and Θ is bounded, then H2 can be implied by the concentration of measure under certain additional regularity conditions.

Case 2: When the samples are generated by an MCMC procedure, H2 can be achieved through an maximal inequality for beta-mixing sequences obtained in [Doukhan et al., 1995]. The condition may also be implied by a number of drift and minorization conditions [Meyn and Tweedie, 2012].

Stationarity measure As problem (1) is a constrained optimization, we consider the following stationarity measure:

$$g(\bar{\boldsymbol{\theta}}) := \inf_{\boldsymbol{\theta} \in \Theta} \frac{\mathcal{L}'(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta} - \bar{\boldsymbol{\theta}})}{\|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|} \quad \text{and} \quad g(\bar{\boldsymbol{\theta}}) = g_+(\bar{\boldsymbol{\theta}}) - g_-(\bar{\boldsymbol{\theta}}), \quad (19)$$

where $g_+(\bar{\boldsymbol{\theta}}) := \max\{0, g(\bar{\boldsymbol{\theta}})\}$, $g_-(\bar{\boldsymbol{\theta}}) := -\min\{0, g(\bar{\boldsymbol{\theta}})\}$ denote the positive and negative part of $g(\bar{\boldsymbol{\theta}})$, respectively. Note that $\bar{\boldsymbol{\theta}}$ is a stationary point if and only if $g_-(\bar{\boldsymbol{\theta}}) = 0$ [Fletcher et al., 2002]. Furthermore, suppose that the sequence $\{\boldsymbol{\theta}^{(k)}\}_{k \geq 0}$ has a limit point $\bar{\boldsymbol{\theta}}$ that is a stationary point, then one has $\lim_{k \rightarrow \infty} g_-(\boldsymbol{\theta}^{(k)}) = 0$. In this sense, the sequence $\{\boldsymbol{\theta}^{(k)}\}_{k \geq 0}$ is said to satisfy an *asymptotic stationary point condition*. This is equivalent to [Mairal, 2015, Definition 2.4].

To explain the condition (19), observe that if $\bar{\theta} \in \text{int}(\Theta)$, the directional derivative can be replaced by the inner product between the gradient $\nabla \mathcal{L}(\bar{\theta})$ and $\theta - \bar{\theta}$, i.e., $\mathcal{L}'(\bar{\theta}, \theta - \bar{\theta}) = \langle \nabla \mathcal{L}(\bar{\theta}) | \theta - \bar{\theta} \rangle$. Therefore, from the definition we have $g(\bar{\theta}) = -\|\nabla \mathcal{L}(\bar{\theta})\| = -g_-(\bar{\theta})$. If in addition $g_-(\bar{\theta}) = 0$, then $\bar{\theta}$ is a stationary point to (1) in the same sense as in unconstrained optimization.

To facilitate our analysis, we define τ_i^k as the iteration index where the i th function is last accessed in the MISSO method prior to iteration k . For example, we have $\tau_{i_k}^{k+1} = k$. We define:

$$\hat{\mathcal{L}}^{(k)}(\theta) := \frac{1}{n} \sum_{i=1}^n \hat{\mathcal{L}}_i(\theta; \theta^{\tau_i^k}), \quad \hat{e}^{(k)}(\theta) := \hat{\mathcal{L}}^{(k)}(\theta) - \mathcal{L}(\theta). \quad (20)$$

We first establish a non-asymptotic convergence rate for the MISSO method:

Theorem 1. *Under S1, S2, H1, H2. For any $K_{\max} \in \mathbb{N}$, let K be an independent discrete r.v. drawn uniformly from $\{0, \dots, K_{\max} - 1\}$ and define the following quantity:*

$$\Delta_{(K_{\max})} := 2nL\mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\theta^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\theta^{(K_{\max})})] + \sum_{k=0}^{K_{\max}-1} \frac{4LC_r}{\sqrt{M_{(k)}}}, \quad (21)$$

Then we have following non-asymptotic bounds:

$$\mathbb{E}[\|\nabla \hat{e}^{(K)}(\theta^{(K)})\|^2] \leq \frac{\Delta_{(K_{\max})}}{K_{\max}}, \quad \mathbb{E}[g_-(\theta^{(K)})] \leq \sqrt{\frac{\Delta_{(K_{\max})}}{K_{\max}}} + \frac{C_{\text{gr}}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2}. \quad (22)$$

Note that $\Delta_{(K_{\max})}$ is finite for any $K_{\max} \in \mathbb{N}$. As expected, the MISSO method converges to a stationary point of (1) asymptotically and at a sublinear rate $\mathbb{E}[g_-(\theta^{(K)})] \leq -\mathcal{O}(\sqrt{1/K_{\max}})$. Furthermore, we remark that the MISO method can be analyzed in Theorem 1 as a special case of the MISSO method satisfying $C_r = C_{\text{gr}} = 0$. In this case, while the asymptotic convergence is well known from [Mairal, 2015] [cf. H2], Eq. (22) gives a non-asymptotic rate of $\mathbb{E}[g_-(\theta^{(K)})] \leq -\mathcal{O}(\sqrt{nL/K_{\max}})$ which is new to our best knowledge.

Next, we show that under an additional assumption on the sequence of batch size $M_{(k)}$, the MISSO method converges almost surely to a stationary point:

Theorem 2. *Under S1, S2, H1, H2. In addition, assume that $\{M_{(k)}\}_{k \geq 0}$ is a non-decreasing sequence of integers which satisfies $\sum_{k=0}^{\infty} M_{(k)}^{-1/2} < \infty$. Then:*

1. *the negative part of the stationarity measure converges almost surely to zero, i.e., $\lim_{k \rightarrow \infty} g_-(\theta^{(k)}) = 0$ a.s..*
2. *the objective value $\mathcal{L}(\theta^{(k)})$ converges almost surely to a finite number $\underline{\mathcal{L}}$, i.e., $\lim_{k \rightarrow \infty} \mathcal{L}(\theta^{(k)}) = \underline{\mathcal{L}}$ a.s..*

In particular, the first result above shows that the sequence $\{\theta^{(k)}\}_{k \geq 0}$ produced by the MISSO method satisfies an *asymptotic stationary point condition*.

4 Numerical Experiments

4.1 Binary logistic regression with missing values

This application follows **Example 1** described in Section 2. We consider a binary regression setup, $((y_i, z_i), i \in \llbracket n \rrbracket)$ where $y_i \in \{0, 1\}$ is a binary response and $z_i = (z_{i,j} \in \mathbb{R}, j \in \llbracket p \rrbracket)$ is a covariate vector. The vector of covariates z_i is not fully observed. We denote by $z_{i,\text{mis}}$ the missing values and $z_{i,\text{obs}}$ the observed covariate. It is assumed that $(z_i, i \in \llbracket n \rrbracket)$ are i.i.d. and marginally distributed according to $\mathcal{N}(\beta, \Omega)$ where $\beta \in \mathbb{R}^p$ and Ω is a positive semidefinite $p \times p$ matrix. Here, $\theta = (\beta, \Omega)$ is the parameter to be estimated.

We are interested in finding the latent structure of the covariates z_i . We define the conditional distribution of the observations y_i given $z_i = (z_{i,\text{mis}}, z_{i,\text{obs}})$ as:

$$p_i(y_i | z_i) = \frac{\exp(-y_i \delta^\top \bar{z}_i)}{1 + \exp(-\delta^\top \bar{z}_i)}, \quad (23)$$

where $\delta = (\delta_0, \dots, \delta_p)$ are, for simplicity, assumed to be known and $\bar{z}_i = (1, z_i)$. For $i \in \llbracket n \rrbracket$, the complete data log-likelihood is expressed as:

$$\log f_i(z_i, \theta) \propto -y_i \delta^\top \bar{z}_i - \log(1 + \exp(-\delta^\top \bar{z}_i)) - \frac{1}{2} \log(|\Omega|) + \frac{1}{2} \text{Tr}(\Omega^{-1}(z_i - \beta)(z_i - \beta)^\top).$$

This model belongs to the curved exponential family [Keener, 2010] where for all $i \in \llbracket n \rrbracket$, the complete data sufficient statistics are given by $\tilde{S}_i(z_i) \triangleq (z_i, z_i z_i^\top)$, see Appendix D.1.2.

At the k -th iteration, and after the initialization, for all $i \in \llbracket n \rrbracket$, of the statistics $(s_i^{1,(0)}, s_i^{2,(0)})$, the MISSO algorithm consists in picking an index i_k uniformly on $\llbracket n \rrbracket$, sampling a Monte Carlo batch $\{z_{i_k, mis, m}^{(k)}\}_{m=1}^{M_{(k)}}$ from the conditional distribution $p(z_{i_k, mis} | z_{i_k, obs}, y_{i_k}; \theta^{(k-1)})$ using an MCMC sampler and computing the quantities $(s_i^{1,(k)}, s_i^{2,(k)})$ as follows:

$$(s_i^{1,(k)}, s_i^{2,(k)}) = \begin{cases} \left(\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} z_{i, m}^{(k)}, \frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} z_{i, m}^{(k)} (z_{i, m}^{(k)})^\top \right) & \text{if } i = i_k \\ (s_i^{1,(k-1)}, s_i^{2,(k-1)}) & \text{otherwise} \end{cases} \quad (24)$$

where $z_{i, m}^{(k)} = (z_{i, mis, m}^{(k)}, z_{i, obs})$ is composed of a simulated and an observed part. Finally, set

$$\beta^{(k)} = \frac{1}{n} \sum_{i=1}^n s_i^{1,(k)} \quad \text{and} \quad \Omega^{(k)} = \frac{1}{n} \sum_{i=1}^n s_i^{2,(k)} - \beta^{(k)} (\beta^{(k)})^\top. \quad (25)$$

Fitting a logistic regression model on the TraumaBase dataset We apply the MISSO method to fit a logistic regression model on the TraumaBase (<http://traumabase.eu>) dataset, which consists of data collected from 15 trauma centers in France, covering measurements on patients from the initial to last stage of trauma.

Similar to [Jiang et al., 2018], we select $p = 16$ influential quantitative measurements, described in Appendix D.1.1, on $n = 6384$ patients, and we adopt the logistic regression model with missing covariates in (23) to predict the risk of a severe hemorrhage which is one of the main cause of death after a major trauma. Note as the dataset considered is heterogeneous – coming from multiple sources with frequently missed entries – we apply the latent data model described in the above. For the Monte-Carlo sampling of $z_{i, mis}$, we run a Metropolis Hastings algorithm with the target distribution $p(\cdot | z_{i, obs}, y_i; \theta^{(k)})$ whose procedure is detailed in Appendix D.1.3.

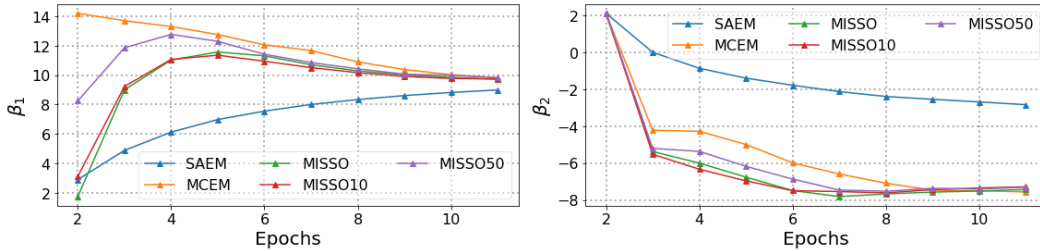


Figure 1: Convergence of two components of the vector of parameters β for the SAEM, the MCEM and the MISSO methods. The convergence is plotted against the number of passes over the data.

We compare in Figure 1 the convergence behavior of the estimated parameters β using SAEM [Delyon et al., 1999] (with stepsize $\gamma_k = 1/k$), MCEM [Wei and Tanner, 1990] and the proposed MISSO method. For the MISSO method, we set the batch size to $M_{(k)} = 10 + k^2$ and we examine with selecting different number of functions in Line 5 in the method – the default settings with 1 function (MISSO), 10% (MISSO10) and 50% (MISSO50) of the functions per iteration. From Figure 1, the MISSO method converges to a static value with less number of epochs than the MCEM, SAEM methods. It is worth noting that the difference among the MISSO runs for different number of selected functions demonstrates a variance-cost tradeoff.

4.2 Fitting Bayesian LeNet-5 on MNIST

This application follows Example 2 described in Section 2. We apply the MISSO method to fit a Bayesian variant of LeNet-5 [LeCun et al., 1998] (see Appendix D.2.1). We train this network on

the MNIST dataset [LeCun, 1998]. The training set is composed of $N = 55\,000$ handwritten digits, 28×28 images. Each image is labelled with its corresponding number (from zero to nine). Under the prior distribution π , see (11), the weights are assumed independent and identically distributed according to $\mathcal{N}(0, 1)$. We also assume that $q(\cdot; \theta) \equiv \mathcal{N}(\mu, \sigma^2 \mathbf{I})$. The variational posterior parameters are thus $\theta = (\mu, \sigma)$ where $\mu = (\mu_\ell, \ell \in \llbracket d \rrbracket)$ where d is the number of weights in the neural network. We use the re-parametrization as $w = t(\theta, z) = \mu + \sigma z$ with $z \sim \mathcal{N}(0, \mathbf{I})$.

At iteration k , minimizing the sum of stochastic surrogates defined as in (6) and (16) yields the following MISSO update — **step (i)** pick a function index i_k uniformly on $\llbracket n \rrbracket$; **step (ii)** sample a Monte Carlo batch $\{z_m^{(k)}\}_{m=1}^{M(k)}$ from $\mathcal{N}(0, \mathbf{I})$; and **step (iii)** update the parameters as

$$\mu_\ell^{(k)} = \frac{1}{n} \sum_{i=1}^n \mu_\ell^{(\tau_i^k)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\mu_\ell, i}^{(k)} \quad \text{and} \quad \sigma^{(k)} = \frac{1}{n} \sum_{i=1}^n \sigma^{(\tau_i^k)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\sigma, i}^{(k)}, \quad (26)$$

where $\hat{\delta}_{\mu_\ell, i}^{(k)} = \hat{\delta}_{\mu_\ell, i}^{(k-1)}$ and $\hat{\delta}_{\sigma, i}^{(k)} = \hat{\delta}_{\sigma, i}^{(k-1)}$ for $i \neq i_k$ and:

$$\begin{aligned} \hat{\delta}_{\mu_\ell, i_k}^{(k)} &= -\frac{1}{M(k)} \sum_{m=1}^{M(k)} \nabla_w \log p(y_{i_k} | x_{i_k}, w) \Big|_{w=t(\theta^{(k-1)}, z_m^{(k)})} + \nabla_{\mu_\ell} d(\theta^{(k-1)}), \\ \hat{\delta}_{\sigma, i_k}^{(k)} &= -\frac{1}{M(k)} \sum_{m=1}^{M(k)} z_m^{(k)} \nabla_w \log p(y_{i_k} | x_{i_k}, w) \Big|_{w=t(\theta^{(k-1)}, z_m^{(k)})} + \nabla_\sigma d(\theta^{(k-1)}) \end{aligned}$$

with $d(\theta) = n^{-1} \sum_{\ell=1}^d (-\log(\sigma) + (\sigma^2 + \mu_\ell^2)/2 - 1/2)$.

We compare the convergence of the *Monte Carlo variants* of the following state of the art optimization algorithms — the ADAM [Kingma and Ba, 2015], the Momentum [Sutskever et al., 2013] and the SAG [Schmidt et al., 2017] methods versus the *Bayes by Backprop* (BBB) [Blundell et al., 2015] and our proposed MISSO method. For all these methods, the loss function (13) and its gradients were computed by Monte Carlo integration using Tensorflow Probability library [Dillon et al., 2017], based on the re-parametrization described above. Update rules for each algorithm are performed using their vanilla implementations on TensorFlow [Abadi et al., 2015] as detailed in Appendix D.2.2. We use the following hyperparameters for all runs — the learning rate is 10^{-3} , we run 100 epochs with a mini-batch size of 128 and use the batchsize of $M(k) = k$.

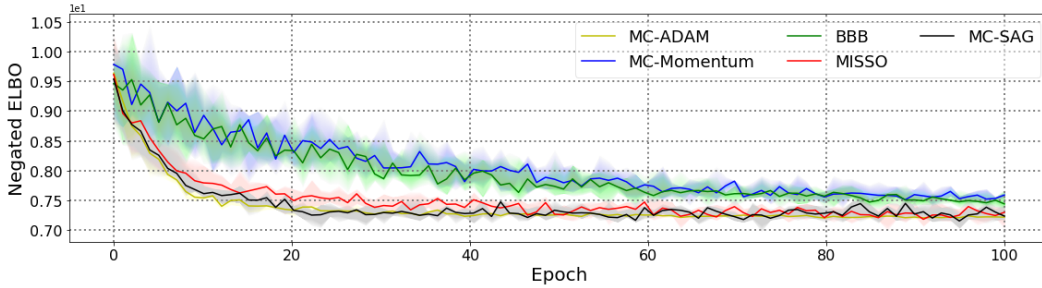


Figure 2: (Incremental Variational Inference) Negated ELBO versus epochs elapsed for fitting the Bayesian LeNet-5 on MNIST using different algorithms. The solid curve is obtained from averaging over 5 independent runs of the methods, and the shaded area represents the standard deviation.

Figure 2 shows the convergence of the negated evidence lower bound against the number of passes over data (one pass represents an epoch). As observed, the proposed MISSO method outperforms *Bayes by Backprop* and Momentum, while similar convergence rates are observed with the MISSO, ADAM and SAG methods.

5 Conclusions

We present a unifying framework for minimizing a non-convex finite-sum objective function using incremental surrogates when the latter functions are expressed as an expectation and are intractable.

Our approach covers a large class of non-convex applications in machine learning such as logistic regression with missing values and variational inference. We provide both finite-time and asymptotic guarantees of our incremental stochastic surrogate optimization technique and illustrate our findings training a binary logistic regression with missing covariates to predict hemorrhagic shock and a Bayesian variant of LeNet-5 on MNIST.

References

- M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017. doi: 10.1080/01621459.2017.1285773. URL <https://doi.org/10.1080/01621459.2017.1285773>.
- C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural network. In *International Conference on Machine Learning*, pages 1613–1622, 2015.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- B. Delyon, M. Lavielle, and E. Moulines. Convergence of a stochastic approximation version of the em algorithm. *Ann. Statist.*, 27(1):94–128, 03 1999. doi: 10.1214/aos/1018031103. URL <https://doi.org/10.1214/aos/1018031103>.
- J. V. Dillon, I. Langmore, D. Tran, E. Brevdo, S. Vasudevan, D. Moore, B. Patton, A. Alemi, M. D. Hoffman, and R. A. Saurous. Tensorflow distributions. *CoRR*, abs/1711.10604, 2017. URL <http://arxiv.org/abs/1711.10604>.
- P. Doukhan, P. Massart, and E. Rio. Invariance principles for absolutely regular empirical processes. In *Annales de l’IHP Probabilités et statistiques*, volume 31, pages 393–427, 1995.
- R. Fletcher, N. I. Gould, S. Leyffer, P. L. Toint, and A. Wächter. Global convergence of a trust-region sqp-filter algorithm for general nonlinear programming. *SIAM Journal on Optimization*, 13(3):635–659, 2002.
- Z. Ghahramani. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452–459, May 2015. doi: 10.1038/nature14541. URL <https://www.ncbi.nlm.nih.gov/pubmed/26017444/>. On Probabilistic models.
- W. Jiang, J. Josse, and M. Lavielle. Logistic regression with missing covariates—parameter estimation, model selection and prediction. 2018.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Mach. Learn.*, 37(2):183–233, Nov. 1999. ISSN 0885-6125. doi: 10.1023/A:1007665907178. URL <https://doi.org/10.1023/A:1007665907178>.
- R. Keener. *Curved Exponential Families*, pages 85–99. Springer New York, New York, NY, 2010. ISBN 978-0-387-93839-4.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL <http://arxiv.org/abs/1312.6114>.
- K. Lange. *MM Optimization Algorithms*. SIAM-Society for Industrial and Applied Mathematics, USA, 2016. ISBN 1611974399, 9781611974393.
- Y. LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.

- Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Y. Li and Y. Gal. Dropout inference in bayesian neural networks with alpha-divergences. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2052–2061. JMLR. org, 2017.
- J. Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM J. Optim.*, 25(2):829–855, 2015. ISSN 1052-6234. doi: 10.1137/140957639. URL <https://doi.org/10.1137/140957639>.
- G. J. McLachlan and T. Krishnan. *The EM algorithm and extensions*. Wiley Series in Probability and Statistics. Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, second edition, 2008. ISBN 978-0-471-20170-0. doi: 10.1002/9780470191613. URL <https://doi.org/10.1002/9780470191613>.
- S. P. Meyn and R. L. Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.
- R. M. Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- J. Paisley, D. Blei, and M. Jordan. Variational bayesian inference with stochastic search. In *ICML*. icml.cc / Omnipress, 2012.
- N. G. Polson, V. Sokolov, et al. Deep learning: a bayesian perspective. *Bayesian Analysis*, 12(4): 1275–1304, 2017.
- M. Razaviyayn, M. Hong, and Z.-Q. Luo. A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization*, 23(2):1126–1153, 2013.
- D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pages 1278–1286, 2014.
- M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2009.
- I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- G. C. G. Wei and M. A. Tanner. A monte carlo implementation of the em algorithm and the poor man’s data augmentation algorithms. *Journal of the American Statistical Association*, 85(411): 699–704, 1990. doi: 10.1080/01621459.1990.10474930. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1990.10474930>.

A Proof of Theorem 1

Theorem. Under S1, S2, H1, H2. For any $K_{\max} \in \mathbb{N}$, let K be an independent discrete r.v. drawn uniformly from $\{0, \dots, K_{\max} - 1\}$ and define the following quantity:

$$\Delta_{(K_{\max})} := 2nL\mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\boldsymbol{\theta}^{(K_{\max})})] + \sum_{k=0}^{K_{\max}-1} \frac{4LC_r}{\sqrt{M_{(k)}}},$$

Then we have following non-asymptotic bounds:

$$\mathbb{E}[\|\nabla \hat{e}^{(K)}(\boldsymbol{\theta}^{(K)})\|^2] \leq \frac{\Delta_{(K_{\max})}}{K_{\max}}, \quad \mathbb{E}[g_{-}(\boldsymbol{\theta}^{(K)})] \leq \sqrt{\frac{\Delta_{(K_{\max})}}{K_{\max}}} + \frac{C_{\text{gr}}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2}.$$

Proof We begin by recalling the definition

$$\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{A}}_i^k(\boldsymbol{\theta}). \quad (27)$$

Notice that

$$\begin{aligned} \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{L}}_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(\tau_i^{k+1})}, \{z_{i,m}^{(\tau_i^{k+1})}\}_{m=1}^{M_{(\tau_i^{k+1})}}) \\ &= \tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) + \frac{1}{n} (\tilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}, \{z_{i_k,m}^{(k)}\}_{m=1}^{M_{(k)}}) - \tilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k,m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})). \end{aligned} \quad (28)$$

Furthermore, we recall that

$$\hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \hat{\mathcal{L}}_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(\tau_i^k)}), \quad \hat{e}^{(k)}(\boldsymbol{\theta}) := \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}). \quad (29)$$

Due to S2, we have

$$\|\nabla \hat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2 \leq 2L\hat{e}^{(k)}(\boldsymbol{\theta}^{(k)}). \quad (30)$$

To prove the first bound in (22), using the optimality of $\boldsymbol{\theta}^{(k+1)}$, one has

$$\begin{aligned} \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) &\leq \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k)}) \\ &= \tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) + \frac{1}{n} (\tilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k,m}^{(k)}\}_{m=1}^{M_{(k)}}) - \tilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k,m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})) \end{aligned} \quad (31)$$

Let \mathcal{F}_k be the filtration of random variables $\{\dots, i_{k-1}, \{z_{i_{k-1},m}^{(k-1)}\}_{m=1}^{M_{(k-1)}}, \boldsymbol{\theta}^{(k)}\}$. We observe that the conditional expectation evaluates to

$$\begin{aligned} \mathbb{E}_{i_k} [\mathbb{E}[\tilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k,m}^{(k)}\}_{m=1}^{M_{(k)}}) | \mathcal{F}_k, i_k] | \mathcal{F}_k] \\ = \mathcal{L}(\boldsymbol{\theta}^{(k)}) + \mathbb{E}_{i_k} [\mathbb{E}[\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} r_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, z_{i_k,m}^{(k)}) - \hat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}) | \mathcal{F}_k, i_k] | \mathcal{F}_k] \\ \leq \mathcal{L}(\boldsymbol{\theta}^{(k)}) + \frac{C_r}{\sqrt{M_{(k)}}}, \end{aligned} \quad (32)$$

where the last inequality is due to H2. Moreover,

$$\mathbb{E}[\tilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k,m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}}) | \mathcal{F}_k] = \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{L}}_i(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}, \{z_{i,m}^{(\tau_i^k)}\}_{m=1}^{M_{(\tau_i^k)}}) = \tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}). \quad (33)$$

Taking the conditional expectations on both sides of (31) and re-arranging terms give:

$$\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)}) \leq n\mathbb{E}[\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) | \mathcal{F}_k] + \frac{C_r}{\sqrt{M_{(k)}}} \quad (34)$$

Proceeding from (34), we observe the following lower bound for the left hand side

$$\begin{aligned}
& \tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)}) \stackrel{(a)}{=} \tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) + \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) \\
& \stackrel{(b)}{\geq} \tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) + \frac{1}{2L} \|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2 \\
& = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} r_i(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}, z_{i,m}^{(\tau_i^k)}) - \hat{\mathcal{L}}_i(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}) \right\} + \frac{1}{2L} \|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2 \\
& \quad \underbrace{\hspace{10em}}_{:= -\delta^{(k)}(\boldsymbol{\theta}^{(k)})}
\end{aligned} \tag{35}$$

where (a) is due to $\hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) = 0$ [cf. S1], (b) is due to (30) and we have defined the summation in the last equality as $-\delta^{(k)}(\boldsymbol{\theta}^{(k)})$. Substituting the above into (34) yields

$$\frac{\|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2}{2L} \leq n \mathbb{E}[\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) | \mathcal{F}_k] + \frac{C_r}{\sqrt{M_{(k)}}} + \delta^{(k)}(\boldsymbol{\theta}^{(k)}) \tag{36}$$

Observe the following upper bound on the total expectations:

$$\mathbb{E}[\delta^{(k)}(\boldsymbol{\theta}^{(k)})] \leq \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \frac{C_r}{\sqrt{M_{(\tau_i^k)}}}\right], \tag{37}$$

which is due to H2. It yields

$$\mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2] \leq 2nL \mathbb{E}[\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)})] + \frac{2LC_r}{\sqrt{M_{(k)}}} + \frac{1}{n} \sum_{i=1}^n \mathbb{E}\left[\frac{2LC_r}{\sqrt{M_{(\tau_i^k)}}}\right]$$

Finally, for any $K_{\max} \in \mathbb{N}$, we let K be a discrete r.v. that is uniformly drawn from $\{0, 1, \dots, K_{\max} - 1\}$. Using H2 and taking total expectations lead to

$$\begin{aligned}
\mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(K)}(\boldsymbol{\theta}^{(K)})\|^2] &= \frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2] \\
&\leq \frac{2nL \mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\boldsymbol{\theta}^{(K_{\max})})]}{K_{\max}} + \frac{2LC_r}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}\left[\frac{1}{\sqrt{M_{(k)}}} + \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{M_{(\tau_i^k)}}}\right]
\end{aligned} \tag{38}$$

For all $i \in [1, n]$, the index i is selected with a probability equal to $\frac{1}{n}$ when conditioned independently on the past. We observe:

$$\mathbb{E}[M_{(\tau_i^k)}^{-1/2}] = \sum_{j=1}^k \frac{1}{n} \left(1 - \frac{1}{n}\right)^{j-1} M_{(k-j)}^{-1/2} \tag{39}$$

Taking the sum yields:

$$\begin{aligned}
\sum_{k=0}^{K_{\max}-1} \mathbb{E}[M_{(\tau_i^k)}^{-1/2}] &= \sum_{k=0}^{K_{\max}-1} \sum_{j=1}^k \frac{1}{n} \left(1 - \frac{1}{n}\right)^{j-1} M_{(k-j)}^{-1/2} = \sum_{k=0}^{K_{\max}-1} \sum_{l=0}^{k-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{k-(l+1)} M_{(l)}^{-1/2} \\
&= \sum_{l=0}^{K_{\max}-1} M_{(l)}^{-1/2} \sum_{k=l+1}^{K_{\max}-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{k-(l+1)} \leq \sum_{l=0}^{K_{\max}-1} M_{(l)}^{-1/2}
\end{aligned} \tag{40}$$

where the last inequality is due to upper bounding the geometric series. Plugging this back into (38) yields

$$\begin{aligned}
\mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(K)}(\boldsymbol{\theta}^{(K)})\|^2] &= \frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2] \\
&\leq \frac{2nL \mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\boldsymbol{\theta}^{(K_{\max})})]}{K_{\max}} + \frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \frac{4LC_r}{\sqrt{M_{(k)}}} = \frac{\Delta_{(K_{\max})}}{K_{\max}}.
\end{aligned} \tag{41}$$

This concludes our proof for the first inequality in (22).

To prove the second inequality of (22), we define the shorthand notations $g^{(k)} := g(\boldsymbol{\theta}^{(k)})$, $g_-^{(k)} := -\min\{0, g^{(k)}\}$, $g_+^{(k)} := \max\{0, g^{(k)}\}$. We observe that

$$\begin{aligned} g^{(k)} &= \inf_{\boldsymbol{\theta} \in \Theta} \frac{\mathcal{L}'(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)})}{\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}\|} \\ &= \inf_{\boldsymbol{\theta} \in \Theta} \left\{ \frac{\frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)})}{\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}\|} - \frac{\langle \nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) | \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)} \rangle}{\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}\|} \right\} \\ &\geq -\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| + \inf_{\boldsymbol{\theta} \in \Theta} \frac{\frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)})}{\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}\|} \end{aligned} \quad (42)$$

where the last inequality is due to the Cauchy-Schwarz inequality and we have defined $\widehat{\mathcal{L}}'_i(\boldsymbol{\theta}, \boldsymbol{d}; \boldsymbol{\theta}^{(\tau_i^k)})$ as the directional derivative of $\widehat{\mathcal{L}}_i(\cdot; \boldsymbol{\theta}^{(\tau_i^k)})$ at $\boldsymbol{\theta}$ along the direction \boldsymbol{d} . Moreover, for any $\boldsymbol{\theta} \in \Theta$,

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}) \\ &= \underbrace{\widetilde{\mathcal{L}}^{(k)'}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}) - \widetilde{\mathcal{L}}^{(k)'}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)})}_{\geq 0} + \frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}) \\ &\geq \frac{1}{n} \sum_{i=1}^n \left\{ \widehat{\mathcal{L}}'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}) - \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} r'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}, z_{i,m}^{(\tau_i^k)}) \right\} \end{aligned} \quad (43)$$

where the inequality is due to the optimality of $\boldsymbol{\theta}^{(k)}$ and the convexity of $\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta})$ [cf. H1]. Denoting a scaled version of the above term as:

$$\epsilon^{(k)}(\boldsymbol{\theta}) := \frac{\frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} r'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}, z_{i,m}^{(\tau_i^k)}) - \widehat{\mathcal{L}}'_i(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}) \right\}}{\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}\|}.$$

We have

$$g^{(k)} \geq -\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| + \inf_{\boldsymbol{\theta} \in \Theta} (-\epsilon^{(k)}(\boldsymbol{\theta})) \geq -\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| - \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})|. \quad (44)$$

Since $g^{(k)} = g_+^{(k)} - g_-^{(k)}$ and $g_+^{(k)} g_-^{(k)} = 0$, this implies

$$g_-^{(k)} \leq \|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| + \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})|. \quad (45)$$

Consider the above inequality when $k = K$, i.e., the random index, and taking total expectations on both sides gives

$$\mathbb{E}[g_-^{(K)}] \leq \mathbb{E}[\|\nabla \widehat{e}^{(K)}(\boldsymbol{\theta}^{(K)})\|] + \mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} \epsilon^{(K)}(\boldsymbol{\theta})] \quad (46)$$

We note that

$$\left(\mathbb{E}[\|\nabla \widehat{e}^{(K)}(\boldsymbol{\theta}^{(K)})\|] \right)^2 \leq \mathbb{E}[\|\nabla \widehat{e}^{(K)}(\boldsymbol{\theta}^{(K)})\|^2] \leq \frac{\Delta(K_{\max})}{K_{\max}}, \quad (47)$$

where the first inequality is due to the convexity of $(\cdot)^2$ and the Jensen's inequality, and

$$\begin{aligned} \mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} \epsilon^{(K)}(\boldsymbol{\theta})] &= \frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}} \mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} \epsilon^{(k)}(\boldsymbol{\theta})] \stackrel{(a)}{\leq} \frac{C_{\text{gr}}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n M_{(\tau_i^k)}^{-1/2}\right] \\ &\stackrel{(b)}{\leq} \frac{C_{\text{gr}}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2} \end{aligned} \quad (48)$$

where (a) is due to H2 and (b) is due to (40). This implies

$$\mathbb{E}[g_-^{(K)}] \leq \sqrt{\frac{\Delta(K_{\max})}{K_{\max}}} + \frac{C_{\text{gr}}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2}, \quad (49)$$

and concludes the proof of the theorem. \square

B Proof of Theorem 2

Theorem. Under S1, S2, H1, H2. In addition, assume that $\{M_{(k)}\}_{k \geq 0}$ is a non-decreasing sequence of integers which satisfies $\sum_{k=0}^{\infty} M_{(k)}^{-1/2} < \infty$. Then:

1. the negative part of the stationarity measure converges almost surely to zero, i.e., $\lim_{k \rightarrow \infty} g_{-}(\boldsymbol{\theta}^{(k)}) = 0$ a.s..
2. the objective value $\mathcal{L}(\boldsymbol{\theta}^{(k)})$ converges almost surely to a finite number $\underline{\mathcal{L}}$, i.e., $\lim_{k \rightarrow \infty} \mathcal{L}(\boldsymbol{\theta}^{(k)}) = \underline{\mathcal{L}}$ a.s..

Proof We apply the following auxiliary lemma which proof can be found in Appendix C for the readability of the current proof:

Lemma 1. Let $(V_k)_{k \geq 0}$ be a non negative sequence of random variables such that $\mathbb{E}[V_0] < \infty$. Let $(X_k)_{k \geq 0}$ a non negative sequence of random variables and $(E_k)_{k \geq 0}$ be a sequence of random variables such that $\sum_{k=0}^{\infty} \mathbb{E}[|E_k|] < \infty$. If for any $k \geq 1$:

$$V_k \leq V_{k-1} - X_{k-1} + E_{k-1} \quad (50)$$

then:

- (i) for all $k \geq 0$, $\mathbb{E}[V_k] < \infty$ and the sequence $(V_k)_{k \geq 0}$ converges a.s. to a finite limit V_{∞} .
- (ii) the sequence $(\mathbb{E}[V_k])_{k \geq 0}$ converges and $\lim_{k \rightarrow \infty} \mathbb{E}[V_k] = \mathbb{E}[V_{\infty}]$.
- (iii) the series $\sum_{k=0}^{\infty} X_k$ converges almost surely and $\sum_{k=0}^{\infty} \mathbb{E}[X_k] < \infty$.

We proceed from (31) by re-arranging terms and observing that

$$\begin{aligned} \widehat{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) &\leq \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \frac{1}{n} (\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})) \\ &\quad - (\widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) - \widehat{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)})) + (\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})) \\ &\quad + \frac{1}{n} (\widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k, m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})) \\ &\quad + \frac{1}{n} (\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}) - \widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k, m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})) \end{aligned} \quad (51)$$

Our idea is to apply Lemma 1. Under S1, the finite sum of surrogate functions $\widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta})$, defined in (20), is lower bounded by a constant $c_k > -\infty$ for any $\boldsymbol{\theta}$. To this end, we observe that

$$V_k := \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \inf_{k \geq 0} c_k \geq 0 \quad (52)$$

is a non-negative random variable.

Secondly, under H1, the following random variable is non-negative

$$X_k := \frac{1}{n} (\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(\tau_{i_k}^k)}; \boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})) \geq 0. \quad (53)$$

Thirdly, we define

$$\begin{aligned} E_k &= -(\widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) - \widehat{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)})) + (\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})) \\ &\quad + \frac{1}{n} (\widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k, m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})) \\ &\quad + \frac{1}{n} (\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}) - \widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k, m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})). \end{aligned} \quad (54)$$

Note that from the definitions (52), (53), (54), we have $V_{k+1} \leq V_k - X_k + E_k$ for any $k \geq 1$.

Under H2, we observe that

$$\mathbb{E}[|\widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k, m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})|] \leq C_r M_{(k)}^{-1/2} \quad (55)$$

$$\mathbb{E}\left[\left|\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}) - \widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k, m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})\right|\right] \leq C_r \mathbb{E}\left[M_{(\tau_{i_k}^k)}^{-1/2}\right] \quad (56)$$

$$\mathbb{E}\left[|\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})|\right] \leq \frac{1}{n} \sum_{i=1}^n C_r \mathbb{E}\left[M_{(\tau_i^k)}^{-1/2}\right] \quad (57)$$

Therefore,

$$\mathbb{E}[|E_k|] \leq \frac{C_r}{n} \left(M_{(k)}^{-1/2} + \mathbb{E}\left[M_{(\tau_i^k)}^{-1/2} + \sum_{i=1}^n \{M_{(\tau_i^k)}^{-1/2} + M_{(\tau_i^{k+1})}^{-1/2}\}\right] \right) \quad (58)$$

Using (40) and the assumption on the sequence $\{M_{(k)}\}_{k \geq 0}$, we obtain that

$$\sum_{k=0}^{\infty} \mathbb{E}[|E_k|] < \frac{C_r}{n} (2 + 2n) \sum_{k=0}^{\infty} M_{(k)}^{-1/2} < \infty. \quad (59)$$

Therefore, the conclusions in Lemma 1 hold. Precisely, we have $\sum_{k=0}^{\infty} X_k < \infty$ and $\sum_{k=0}^{\infty} \mathbb{E}[X_k] < \infty$ almost surely. Note that this implies

$$\begin{aligned} \infty &> \sum_{k=0}^{\infty} \mathbb{E}[X_k] = \frac{1}{n} \sum_{k=0}^{\infty} \mathbb{E}[\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})] \\ &= \frac{1}{n} \sum_{k=0}^{\infty} \mathbb{E}[\widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)})] = \frac{1}{n} \sum_{k=0}^{\infty} \mathbb{E}[\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})] \end{aligned} \quad (60)$$

Since $\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) \geq 0$, the above implies

$$\lim_{k \rightarrow \infty} \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) = 0 \quad \text{a.s.} \quad (61)$$

and subsequently applying (30), we have $\lim_{k \rightarrow \infty} \|\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| = 0$ almost surely. Finally, it follows from (30) and (45) that

$$\lim_{k \rightarrow \infty} g_-^{(k)} \leq \lim_{k \rightarrow \infty} \sqrt{2L} \sqrt{\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})} + \lim_{k \rightarrow \infty} \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})| = 0, \quad (62)$$

where the last equality holds almost surely due to the fact that $\sum_{k=0}^{\infty} \mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})|] < \infty$. This concludes the asymptotic convergence of the MISSO method.

Finally, we prove that $\mathcal{L}(\boldsymbol{\theta}^{(k)})$ converges almost surely. As a consequence of Lemma 1, it is clear that $\{V_k\}_{k \geq 0}$ converges almost surely and so is $\{\widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\}_{k \geq 0}$, i.e., we have $\lim_{k \rightarrow \infty} \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) = \underline{\mathcal{L}}$. Applying (61) implies that

$$\underline{\mathcal{L}} = \lim_{k \rightarrow \infty} \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) = \lim_{k \rightarrow \infty} \mathcal{L}(\boldsymbol{\theta}^{(k)}) \quad \text{a.s.} \quad (63)$$

This shows that $\mathcal{L}(\boldsymbol{\theta}^{(k)})$ converges almost surely to $\underline{\mathcal{L}}$. \square

C Proof of Lemma 1

Lemma. Let $(V_k)_{k \geq 0}$ be a non negative sequence of random variables such that $\mathbb{E}[V_0] < \infty$. Let $(X_k)_{k \geq 0}$ a non negative sequence of random variables and $(E_k)_{k \geq 0}$ be a sequence of random variables such that $\sum_{k=0}^{\infty} \mathbb{E}[|E_k|] < \infty$. If for any $k \geq 1$:

$$V_k \leq V_{k-1} - X_{k-1} + E_{k-1}$$

then:

- (i) for all $k \geq 0$, $\mathbb{E}[V_k] < \infty$ and the sequence $(V_k)_{k \geq 0}$ converges a.s. to a finite limit V_{∞} .
- (ii) the sequence $(\mathbb{E}[V_k])_{k \geq 0}$ converges and $\lim_{k \rightarrow \infty} \mathbb{E}[V_k] = \mathbb{E}[V_{\infty}]$.
- (iii) the series $\sum_{k=0}^{\infty} X_k$ converges almost surely and $\sum_{k=0}^{\infty} \mathbb{E}[X_k] < \infty$.

Proof We first show that for all $k \geq 0$, $\mathbb{E}[V_k] < \infty$. Note indeed that:

$$0 \leq V_k \leq V_0 - \sum_{j=1}^k X_j + \sum_{j=1}^k E_j \leq V_0 + \sum_{j=1}^k E_j \quad (64)$$

showing that $\mathbb{E}[V_k] \leq \mathbb{E}[V_0] + \mathbb{E}\left[\sum_{j=1}^k E_j\right] < \infty$.

Since $0 \leq X_k \leq V_{k-1} - V_k + E_k$ we also obtain for all $k \geq 0$, $\mathbb{E}[X_k] < \infty$. Moreover, since $\mathbb{E}\left[\sum_{j=1}^{\infty} |E_j|\right] < \infty$, the series $\sum_{j=1}^{\infty} E_j$ converges a.s. We may therefore define:

$$W_k = V_k + \sum_{j=k+1}^{\infty} E_j \quad (65)$$

Note that $\mathbb{E}[|W_k|] \leq \mathbb{E}[V_k] + \mathbb{E}\left[\sum_{j=k+1}^{\infty} |E_j|\right] < \infty$. For all $k \geq 1$, we get:

$$\begin{aligned} W_k &\leq V_{k-1} - X_k + \sum_{j=k}^{\infty} E_j \leq W_{k-1} - X_k \leq W_{k-1} \\ \mathbb{E}[W_k] &\leq \mathbb{E}[W_{k-1}] - \mathbb{E}[X_k] \end{aligned} \quad (66)$$

Hence the sequences $(W_k)_{k \geq 0}$ and $(\mathbb{E}[W_k])_{k \geq 0}$ are non increasing. Since for all $k \geq 0$, $W_k \geq -\sum_{j=1}^{\infty} |E_j| > -\infty$ and $\mathbb{E}[W_k] \geq -\sum_{j=1}^{\infty} \mathbb{E}[|E_j|] > -\infty$, the (random) sequence $(W_k)_{k \geq 0}$ converges a.s. to a limit W_{∞} and the (deterministic) sequence $(\mathbb{E}[W_k])_{k \geq 0}$ converges to a limit w_{∞} . Since $|W_k| \leq V_0 + \sum_{j=1}^{\infty} |E_j|$, the Fatou lemma implies that:

$$\mathbb{E}[\liminf_{k \rightarrow \infty} |W_k|] = \mathbb{E}[|W_{\infty}|] \leq \liminf_{k \rightarrow \infty} \mathbb{E}[|W_k|] \leq \mathbb{E}[V_0] + \sum_{j=1}^{\infty} \mathbb{E}[|E_j|] < \infty \quad (67)$$

showing that the random variable W_{∞} is integrable.

In the sequel, set $U_k \triangleq W_0 - W_k$. By construction we have for all $k \geq 0$, $U_k \geq 0$, $U_k \leq U_{k+1}$ and $\mathbb{E}[U_k] \leq \mathbb{E}[|W_0|] + \mathbb{E}[|W_k|] < \infty$ and by the monotone convergence theorem, we get:

$$\lim_{k \rightarrow \infty} \mathbb{E}[U_k] = \mathbb{E}[\lim_{k \rightarrow \infty} U_k] \quad (68)$$

Finally, we have:

$$\lim_{k \rightarrow \infty} \mathbb{E}[U_k] = \mathbb{E}[W_0] - w_{\infty} \quad \text{and} \quad \mathbb{E}[\lim_{k \rightarrow \infty} U_k] = \mathbb{E}[W_0] - \mathbb{E}[W_{\infty}] \quad (69)$$

showing that $\mathbb{E}[W_{\infty}] = w_{\infty}$ and concluding the proof of (ii). Moreover, using (66) we have that $W_k \leq W_{k-1} - X_k$ which yields:

$$\begin{aligned} \sum_{j=1}^{\infty} X_j &\leq W_0 - W_{\infty} < \infty \\ \sum_{j=1}^{\infty} \mathbb{E}[X_j] &\leq \mathbb{E}[W_0] - w_{\infty} < \infty \end{aligned} \quad (70)$$

which concludes the proof of the lemma. \square

D Details about the Numerical Experiments

D.1 Binary Logistic Regression on the Traumabase

D.1.1 Traumabase quantitative variables

The list of the 16 quantitative variables we use in our experiments are as follows — *age*, *weight*, *height*, *BMI (Body Mass Index)*, *the Glasgow Coma Scale*, *the Glasgow Coma Scale motor component*, *the minimum systolic blood pressure*, *the minimum diastolic blood pressure*, *the maximum*

number of heart rate (or pulse) per unit time (usually a minute), the systolic blood pressure at arrival of ambulance, the diastolic blood pressure at arrival of ambulance, the heart rate at arrival of ambulance, the capillary Hemoglobin concentration, the oxygen saturation, the fluid expansion colloids, the fluid expansion cristaloids, the pulse pressure for the minimum value of diastolic and systolic blood pressure, the pulse pressure at arrival of ambulance.

D.1.2 MISSO, MCEM and SAEM for Curved Exponential Family

Consider the case where the complete model $\theta \rightarrow f_i(z_i, \theta)$ belongs to the curved exponential family. For all $i \in \llbracket 1, n \rrbracket$ and $\theta \in \Theta$, one has:

$$\log f_i(z_i, \theta) = H_i(z_i) - \psi_i(\theta) + \langle \tilde{S}_i(z_i) | \phi_i(\theta) \rangle. \quad (71)$$

where $\psi_i : \Theta \rightarrow \mathbb{R}$ and $\phi_i : \Theta \rightarrow \mathbb{R}$ are twice continuously differentiable functions of θ , $H_i : \mathcal{Z} \rightarrow \mathbb{R}$ is a twice continuously differentiable function of z_i and $\tilde{S}_i : \mathcal{Z} \rightarrow \mathcal{S}_i$ is a statistic taking its values in a convex subset \mathcal{S}_i of \mathbb{R} and such that $\int_{\mathcal{Z}} |\tilde{S}_i(z_i)| p_i(z_i, \theta) \mu_i(dz_i) < \infty$. Finally, the i th loss function is given by

$$\mathcal{L}_i(\theta) = -\log \int_{\mathcal{Z}} f_i(z_i, \theta) \mu_i(dz_i). \quad (72)$$

To derive the surrogate function for this example, we follow from **Example 1** in Section 2 by taking $p(z_i, \bar{\theta}) = f(z_i, \bar{\theta})/g(\bar{\theta})$ to yield

$$\mathcal{L}_i(\theta) \leq \hat{\mathcal{L}}_i(\theta; \bar{\theta}) := \text{constant} + \int_{\mathcal{Z}} \{ -\log f_i(z_i, \theta) \} p(z_i, \bar{\theta}) \mu_i(dz_i) \quad (73)$$

Using (71), the latter integral can be further expressed as

$$\text{constant} + \psi_i(\theta) - \int_{\mathcal{Z}} \langle \tilde{S}_i(z_i) | \phi_i(\theta) \rangle p(z_i, \bar{\theta}) \mu_i(dz_i). \quad (74)$$

To summarize, we observe that the surrogate function can be written as $\hat{\mathcal{L}}_i(\theta; \bar{\theta}) = \text{constant} + \psi_i(\theta) - \langle \int_{\mathcal{Z}} \tilde{S}_i(z_i) p(z_i, \bar{\theta}) \mu_i(dz_i) | \phi_i(\theta) \rangle$. As suggested in the MISSO method, we note that the integral can be approximated using sample averages.

Now consider the following function $L(s; \theta)$:

$$L(s; \theta) := \frac{1}{n} \sum_{i=1}^n \left\{ \psi_i(\theta) - \langle s_i | \phi_i(\theta) \rangle \right\}, \quad (75)$$

for all $\theta \in \Theta$ and $s = (s_i, 1 \leq i \leq N) \in \mathcal{S} := \times_{i=1}^N \mathcal{S}_i$. As noted from (74), when s_i is obtained from an appropriately sampling procedure, (75) constitutes the summed stochastic surrogate function $\tilde{\mathcal{L}}(\theta)$ as suggested in the MISSO method [cf. Algorithm 2]. For simplicity, we consider a case when the minimization w.r.t. θ can be uniquely attained and we denote $\theta^{(s)} := \arg \min_{\theta \in \Theta} L(s; \theta)$. We note that for this special case, the MISSO method can be compactly described in Algorithm 3.

Now, in the context of the logistic regression described in section 4.1, for all $i \in \llbracket 1, n \rrbracket$, the complete log likelihood is expressed as:

$$\log f_i(z_i, \theta) \propto -y_i \delta^\top \bar{z}_i - \log(1 + \exp(-\delta^\top \bar{z}_i)) - \frac{1}{2} \log(|\Omega|) + \frac{1}{2} \text{Tr}(\Omega^{-1}(z_i - \beta)(z_i - \beta)^\top).$$

For all $i \in \llbracket 1, n \rrbracket$, the sufficient statistics is $\tilde{S}_i(z_i) \triangleq (z_i, z_i z_i^\top)$ and the minimizing map $s \rightarrow \hat{\theta}(s)$ can be derived as:

$$\hat{\theta}((s_{i,1}, s_{i,2})_{i=1}^n) = \left(\frac{1}{n} \sum_{i=1}^n s_{i,1}, \frac{1}{n} \sum_{i=1}^n s_{i,2} - \frac{1}{n^2} \sum_{i=1}^n s_{i,1} \left(\sum_{i=1}^n s_{i,1} \right)^\top \right) \quad (77)$$

Let us compare the MISSO method to the MCEM, SAEM methods as follows. First we note that the MISSO method differs from the MCEM method in that the latter picks *all the data sample* at each iteration, i.e., Line 3 to 5 is performed for every $i \in \llbracket 1, n \rrbracket$ instead of only one randomly chosen

Algorithm 3 MISSO for curved exponential family

- 1: **Initialization:** given an initial parameter estimate $\theta^{(0)}$, for all $i \in \llbracket 1, n \rrbracket$, sample a Monte Carlo batch $\{z_{i,m}^{(0)}\}_{m=1}^{M_{(0)}}$ from $p(z_i, \theta^{(0)})$ and compute $s_i^{(0)} = \frac{1}{M_{(0)}} \sum_{m=1}^{M_{(0)}} \tilde{S}_i(z_{i,m}^{(0)})$.
- 2: **for** $k = 0, 1, \dots$, **do**
- 3: Pick a function index i_k uniformly on $\llbracket 1, n \rrbracket$.
- 4: Draw $M_{(k)}$ Monte-Carlo samples $\{z_{i,m}^{(k)}\}_{m=1}^{M_{(k)}}$ from $p(z_i, \theta^{(k-1)})$.
- 5: Update the individual sufficient statistics recursively — for each $i \in \llbracket 1, n \rrbracket$,

$$s_i^{(k)} = \begin{cases} \frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} \tilde{S}_i(z_{i,m}^{(k)}) & \text{if } i = i_k \\ s_i^{(k-1)} & \text{otherwise} \end{cases} \quad (76)$$

- 6: Set $\theta^{(k)} = \theta^{(s^{(k)})} = \arg \min_{\theta \in \Theta} L(s^{(k)}; \theta)$ where $s^{(k)} = (s_i^{(k)}, 1 \leq i \leq n)$.
 - 7: **end for**
-

Algorithm 4 MH aglorithm

- 1: **Input:** initialization $z_{i,mis,0} \sim q(z_{i,mis}; \delta)$
 - 2: **for** $m = 1, \dots, M$ **do**
 - 3: Sample $z_{i,mis,m} \sim q(z_{i,mis}; \delta)$
 - 4: Sample $u \sim \mathcal{U}(\llbracket 0, 1 \rrbracket)$
 - 5: Calculate the ratio $r = \frac{\pi(z_{i,mis,m}; \theta)/q(z_{i,mis,m}; \delta)}{\pi(z_{i,mis,m-1}; \theta)/q(z_{i,mis,m-1}; \delta)}$
 - 6: **if** $u < r$ **then**
 - 7: Accept $z_{i,mis,m}$
 - 8: **else**
 - 9: $z_{i,mis,m} \leftarrow z_{i,mis,m-1}$
 - 10: **end if**
 - 11: **end for**
 - 12: **Output:** $z_{i,mis,M}$
-

$i_k \in \llbracket 1, n \rrbracket$. Second, the SAEM method consists of a different update of the individual sufficient statistics — at iteration k and for all $i \in \llbracket 1, n \rrbracket$, the SAEM updates the quantity $s_i^{(k)}$ as follows:

$$s_i^{(k)} = s_i^{(k-1)} + \gamma_k \left(\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} \tilde{S}_i(z_{i,m}^{(k)}) - s_i^{(k-1)} \right) \quad (78)$$

where $\{\gamma_k\}_{k>0}$ is a decreasing sequence of stepsizes. Note that when $\gamma_k = 1$, the above update recovers the MCEM method.

D.1.3 Metropolis Hastings algorithm

During the simulation step of the MISSO method, the sampling from the target distribution $\pi(z_{i,mis}; \theta) := p(z_{i,mis} | z_{i,obs}, y_i; \theta)$ is performed using a Metropolis Hastings (MH) algorithm [Meyn and Tweedie, 2012] with proposal distribution $q(z_{i,mis}; \delta) := p(z_{i,mis} | z_{i,obs}; \delta)$ where $\theta = (\beta, \Omega)$ and $\delta = (\xi, \Sigma)$. The parameters of the Gaussian conditional distribution of $z_{i,mis} | z_{i,obs}$ reads:

$$\begin{aligned} \xi &= \beta_{mis} + \Omega_{mis,obs} \Omega_{obs,obs}^{-1} (z_{i,obs} - \beta_{obs}), \\ \Sigma &= \Omega_{mis,mis} + \Omega_{mis,obs} \Omega_{obs,obs}^{-1} \Omega_{obs,mis} \end{aligned} \quad (79)$$

where we have used the Schur Complement of $\Omega_{obs,obs}$ in Ω and noted β_{mis} (resp. β_{obs}) the missing (resp. observed) elements of β . The MH algorithm is summarized in Algorithm 4.

D.2 Incremental Variational Inference for MNIST

D.2.1 Bayesian LeNet-5 Architecture

We describe in Table D.2.1 the architecture of the Convolutional Neural Network introduced in [LeCun et al., 1998] and trained on MNIST:

layer type	width	stride	padding	input shape	nonlinearity
convolution (5×5)	6	1	0	$1 \times 32 \times 32$	ReLU
max-pooling (2×2)		2	0	$6 \times 28 \times 28$	
convolution (5×5)	6	1	0	$1 \times 14 \times 14$	ReLU
max-pooling (2×2)		2	0	$16 \times 10 \times 10$	
fully-connected	120			400	ReLU
fully-connected	84			120	ReLU
fully-connected	10			84	

Table 1: LeNet-5 architecture

D.2.2 Algorithms updates

First, we initialize the means $\mu_\ell^{(0)}$ for $\ell \in \llbracket d \rrbracket$ and variance estimates $\sigma^{(0)}$. In the sequel, at iteration k and for all $i \in \llbracket n \rrbracket$ we define the following drift terms:

$$\begin{aligned}\hat{\delta}_{\mu_\ell, i_k}^{(k)} &= -\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} \nabla_w \log p(y_{i_k} | x_{i_k}, w) \Big|_{w=t(\theta^{(k-1)}, z_m^{(k)})} + \nabla_{\mu_\ell} d(\theta^{(k-1)}), \\ \hat{\delta}_{\sigma, i_k}^{(k)} &= -\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} z_m^{(k)} \nabla_w \log p(y_{i_k} | x_{i_k}, w) \Big|_{w=t(\theta^{(k-1)}, z_m^{(k)})} + \nabla_{\sigma} d(\theta^{(k-1)}).\end{aligned}\tag{80}$$

For all benchmark algorithms, we pick, at iteration k , a function index i_k uniformly on $\llbracket n \rrbracket$ and sample a Monte Carlo batch $\{z_m^{(k)}\}_{m=1}^{M_{(k)}}$ from the standard Gaussian distribution. The updates of the parameters μ_ℓ for all $\ell \in \llbracket d \rrbracket$ and σ break down as follows:

Monte Carlo SAG update: Set

$$\mu_\ell^{(k)} = \mu_\ell^{(k-1)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\mu_\ell, i}^{(k)} \quad \text{and} \quad \sigma^{(k)} = \sigma^{(k-1)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\sigma, i}^{(k)}, \tag{81}$$

where $\hat{\delta}_{\mu_\ell, i}^{(k)} = \hat{\delta}_{\mu_\ell, i}^{(k-1)}$ and $\hat{\delta}_{\sigma, i}^{(k)} = \hat{\delta}_{\sigma, i}^{(k-1)}$ for $i \neq i_k$ and are defined by (80) for $i = i_k$. where the learning rate $\gamma = 10^{-3}$.

Bayes By Backprop update: Set

$$\mu_\ell^{(k)} = \mu_\ell^{(k-1)} - \frac{\gamma}{n} \hat{\delta}_{\mu_\ell, i_k}^{(k)} \quad \text{and} \quad \sigma^{(k)} = \sigma^{(k-1)} - \frac{\gamma}{n} \hat{\delta}_{\sigma, i_k}^{(k)}, \tag{82}$$

where the learning rate $\gamma = 10^{-3}$.

Monte Carlo Momentum update: Set

$$\mu_\ell^{(k)} = \mu_\ell^{(k-1)} + \hat{v}_{\mu_\ell}^{(k)} \quad \text{and} \quad \sigma^{(k)} = \sigma^{(k-1)} + \hat{v}_{\sigma}^{(k)}, \tag{83}$$

where

$$\hat{v}_{\mu_\ell, i}^{(k)} = \alpha \hat{v}_{\mu_\ell, i}^{(k-1)} - \frac{\gamma}{n} \hat{\delta}_{\mu_\ell, i_k}^{(k)} \quad \text{and} \quad \hat{v}_{\sigma}^{(k)} = \alpha \hat{v}_{\sigma}^{(k-1)} - \frac{\gamma}{n} \hat{\delta}_{\sigma, i_k}^{(k)}, \tag{84}$$

where α and γ , respectively the momentum and the learning rates, are set to 10^{-3} .

Monte Carlo ADAM update: Set

$$\mu_\ell^{(k)} = \mu_\ell^{(k-1)} - \frac{\gamma}{n} \hat{m}_{\mu_\ell}^{(k)} / (\sqrt{\hat{m}_{\mu_\ell}^{(k)}} + \epsilon) \quad \text{and} \quad \sigma^{(k)} = \sigma^{(k-1)} - \frac{\gamma}{n} \hat{m}_{\sigma}^{(k)} / (\sqrt{\hat{m}_{\sigma}^{(k)}} + \epsilon), \tag{85}$$

where

$$\begin{aligned}\hat{\mathbf{m}}_{\mu_\ell}^{(k)} &= \mathbf{m}_{\mu_\ell}^{(k-1)} / (1 - \rho_1^k) \quad \text{with} \quad \mathbf{m}_{\mu_\ell}^{(k)} = \rho_1 \mathbf{m}_{\mu_\ell}^{(k-1)} + (1 - \rho_1) \hat{\boldsymbol{\delta}}_{\mu_\ell, i_k}^{(k)}, \\ \hat{\mathbf{v}}_{\mu_\ell}^{(k)} &= \mathbf{v}_{\mu_\ell}^{(k-1)} / (1 - \rho_2^k) \quad \text{with} \quad \mathbf{v}_{\mu_\ell}^{(k)} = \rho_2 \mathbf{v}_{\mu_\ell}^{(k-1)} + (1 - \rho_2) (\hat{\boldsymbol{\delta}}_{\mu_\ell, i_k}^{(k)})^2\end{aligned}\tag{86}$$

and

$$\begin{aligned}\hat{\mathbf{m}}_\sigma^{(k)} &= \mathbf{m}_\sigma^{(k-1)} / (1 - \rho_1^k) \quad \text{with} \quad \mathbf{m}_\sigma^{(k)} = \rho_1 \mathbf{m}_\sigma^{(k-1)} + (1 - \rho_1) \hat{\boldsymbol{\delta}}_{\sigma, i_k}^{(k)}, \\ \hat{\mathbf{v}}_\sigma^{(k)} &= \mathbf{v}_\sigma^{(k-1)} / (1 - \rho_2^k) \quad \text{with} \quad \mathbf{v}_\sigma^{(k)} = \rho_2 \mathbf{v}_\sigma^{(k-1)} + (1 - \rho_2) (\hat{\boldsymbol{\delta}}_{\sigma, i_k}^{(k)})^2.\end{aligned}\tag{87}$$

The hyperparameters are set as follows: $\gamma = 10^{-3}$, $\rho_1 = 0.9$, $\rho_2 = 0.999$, $\epsilon = 10^{-8}$.