MISSO: Minimization by Incremental Stochastic Surrogate Optimization for Large Scale Nonconvex and Nonsmooth Problems

Anonymous Author(s)

Affiliation Address email

Abstract

Many constrained, nonconvex and nonsmooth optimization problems can be tack-led 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, nonconvex, nonsmooth optimization setting. We apply our new framework for inference of logistic regression model with missing data and for variational inference of Bayesian variants of LeNet-5 and Resnet-18 on respectively the MNIST and CIFAR-10 datasets.

1 Introduction

2

3

8

9

10

11

12

13

14

15

20

21

22

23

24

25

28

29

30

31

32

We consider the *constrained* minimization problem of a finite sum of functions:

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

where Θ is a convex, compact, and closed subset of \mathbb{R}^p , and for any $i \in [\![1,n]\!]$, the function \mathcal{L}_i : $\mathbb{R}^p \to \mathbb{R}$ is bounded from below and is (possibly) nonconvex and nonsmooth.

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]. An extended analysis of MISO has been proposed in [Qian et al., 2019].

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. For instance, this

is the case for inference in latent variable models via maximum likelihood [McLachlan and Krishnan, 2008]. 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 method called *Minimization by Incremental Stochastic Surrogate Optimization (MISSO)*, designed for the nonconvex and nonsmooth finite sum optimization, with a finite-time convergence guarantee. Our work aims at formulating a *generic class* of incremental stochastic surrogate methods for nonconvex optimization and building the theory to understand its behavior. In particular, we provide convergence guarantees for stochastic EM and Variational Inference-type methods, under mild conditions. In summary, our contributions are:

- we propose a unifying framework of analysis for incremental stochastic surrogate optimization when the surrogates are defined as 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.
- we present an incremental update of the commonly used variational inference and Monte Carlo EM methods as special cases of our newly introduced framework. The analysis of those two algorithms is thus conducted under this unifying framework of analysis.
- 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, see Theorem 1.

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 [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 intractable surrogate functions requiring approximation. In Section 3, we provide the asymptotic and non-asymptotic convergence analysis for the MISSO method (and of the MISO [Mairal, 2015] one as a special case). Section 4 presents numerical applications including parameter inference for logistic regression with missing data and variational inference for two types of Bayesian neural networks. The proofs of theoretical results are reported as Supplement.

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

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

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

2 Incremental Minimization of Finite Sum Nonconvex Functions

The objective function in (1) is composed of a finite sum of possibly nonsmooth and nonconvex functions. A popular approach here is to apply the MM method, which 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 [Mairal, 2015] is developed as an iterative scheme that only updates the surrogate functions *partially* at each iteration. Formally, for any $i \in [1, n]$, we consider a surrogate function $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ which satisfies the assumptions (H1, H2):

H1. For all $i \in [1, n]$ and $\overline{\theta} \in \Theta$, $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ is convex w.r.t. θ , and it holds

$$\widehat{\mathcal{L}}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) \ge \mathcal{L}_i(\boldsymbol{\theta}), \ \forall \ \boldsymbol{\theta} \in \Theta \ ,$$
 (3)

where the equality holds when $oldsymbol{ heta}=\overline{oldsymbol{ heta}}.$

H2. For any $\overline{\theta}_i \in \Theta$, $i \in [\![1,n]\!]$ 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 $\theta \in \Theta$, 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 θ , the gradient satisfies

$$\|\nabla \widehat{e}(\boldsymbol{\theta}; \{\overline{\boldsymbol{\theta}}_i\}_{i=1}^n)\|^2 \le 2L\widehat{e}(\boldsymbol{\theta}; \{\overline{\boldsymbol{\theta}}_i\}_{i=1}^n), \ \forall \ \boldsymbol{\theta} \in \Theta.$$
 (4)

We remark that H1 is a common assumption used for surrogate functions, see [Mairal, 2015, 81 Section 2.3]. H2 can be satisfied when the differ-82 ence function $\widehat{e}(\boldsymbol{\theta}; \{\overline{\boldsymbol{\theta}}_i\}_{i=1}^n)$ is L-smooth, i.e., \widehat{e} is differentiable on Θ and its gradient $\nabla \hat{e}$ is L-84 Lipschitz, $\forall \theta \in \Theta$. H2 can be implied by apply-85 ing [Razaviyayn et al., 2013, Proposition 1]. 86 The inequality (3) implies $\widehat{\mathcal{L}}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) \geq \mathcal{L}_i(\boldsymbol{\theta}) > -\infty$ for any $\boldsymbol{\theta} \in \Theta$. The MISO method is 87 88 an incremental version of the MM method, as 89 summarized by Algorithm 1, which shows that 90 the MISO method maintains an iteratively updated set of upper-bounding surrogate functions 92 $\{\mathcal{A}_i^k(\boldsymbol{\theta})\}_{i=1}^n$ and updates the iterate via minimiz-

ing the average of the surrogate functions.

93

94

Algorithm 1 The MISO method [Mairal, 2015].

1: **Input:** initialization $\theta^{(0)}$.

2: Initialize the surrogate function as $\mathcal{A}_i^0(\boldsymbol{\theta}) := \widehat{\mathcal{L}}_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(0)}), i \in [1, n].$

3: **for** $k = 0, 1, ..., K_{\text{max}}$ **do**

4: Pick i_k uniformly from [1, n].

5: Update $A_i^{k+1}(\boldsymbol{\theta})$ as:

$$\mathcal{A}_i^{k+1}(oldsymbol{ heta}) = egin{cases} \widehat{\mathcal{L}}_i(oldsymbol{ heta}; oldsymbol{ heta}^{(k)}), & ext{if } i = i_k \\ \mathcal{A}_i^k(oldsymbol{ heta}), & ext{otherwise}. \end{cases}$$

6: Set
$$\boldsymbol{\theta}^{(k+1)} \in \operatorname*{arg\,min}_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \mathcal{A}_{i}^{k+1}(\boldsymbol{\theta})$$
.

7: end for

Particularly, only one out of the n surrogate functions is updated at each iteration [cf. Line 5] and 95 the sum function $\frac{1}{n}\sum_{i=1}^{n}\mathcal{A}_{i}^{k+1}(\theta)$ is designed to be 'easy to optimize', which, for example, can be a sum of quadratic functions. As such, the MISO method is suitable for large-scale optimization as 96 97 the computation cost per iteration is independent of n. Under H1, H2, it was shown that the MISO 98 method converges almost surely to a stationary point of (1) [Mairal, 2015, Prop. 3.1]. 99

We now consider the case when the surrogate functions $\widehat{\mathcal{L}}_i(m{ heta};\overline{m{ heta}})$ are intractable. Let Z be a mea-100 surable set, $p_i: \mathsf{Z} \times \Theta \to \mathbb{R}_+$ a probability density function, $r_i: \Theta \times \Theta \times \mathsf{Z} \to \mathbb{R}$ a measurable 101 function and μ_i a σ -finite measure. We consider surrogate functions which satisfy H1, H2 and that 102 can be expressed as an expectation, *i.e.*: 103

$$\widehat{\mathcal{L}}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) := \int_{\mathbf{Z}} r_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, z_i) p_i(z_i; \overline{\boldsymbol{\theta}}) \mu_i(dz_i) \quad \forall \ (\boldsymbol{\theta}, \overline{\boldsymbol{\theta}}) \in \Theta \times \Theta \ . \tag{5}$$

Plugging (5) into the MISO method is not feasible since the update step in Step 6 involves a mini-104 mization of an expectation. Several motivating examples of (1) are given in Section 2. 105

In this paper, we propose the Minimization by Incremental Stochastic Surrogate Optimization 106 (MISSO) method which replaces the expectation in (5) by Monte Carlo integration and then op-107 timizes the objective function (1) in an incremental manner. Denote by $M\in\mathbb{N}$ the Monte Carlo 108 batch size and let $\{z_m \in \mathsf{Z}\}_{m=1}^M$ be a set of samples. These samples can be drawn (Case 1) i.i.d. 109 from the distribution $p_i(\cdot; \overline{\theta})$ or (Case 2) from a Markov chain with stationary distribution $p_i(\cdot; \overline{\theta})$; 110 see Section 3 for illustrations. To this end, we define the stochastic surrogate as follows: 111

$$\widetilde{\mathcal{L}}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, \{z_m\}_{m=1}^M) := \frac{1}{M} \sum_{m=1}^M r_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, z_m) , \qquad (6)$$

and we summarize the proposed MISSO method in Algorithm 2. Compared to the MISO method, 112 there is a crucial difference in that the MISSO method involves two types of randomness. The first 113 level of randomness comes from the selection of i_k in Line 5. The second level of randomness stems 114 from the set of Monte Carlo approximated functions $\widetilde{\mathcal{A}}_i^k(\boldsymbol{\theta})$ used in lieu of $\mathcal{A}_i^k(\boldsymbol{\theta})$ in Line 6 when optimizing for the next iterate $\boldsymbol{\theta}^{(k)}$. We now discuss two applications of the MISSO method. 115

Example 1: Maximum Likelihood Estimation for Latent Variable Model. Latent variable mod-117 els [Bishop, 2006] are constructed by introducing unobserved (latent) variables which help explain 118 the observed data. We consider n independent observations $((y_i, z_i), i \in [n])$ where y_i is observed 119 and z_i is latent. In this incomplete data framework, define $\{f_i(z_i,\theta),\theta\in\Theta\}$ to be the complete 120 data likelihood models, i.e., the joint likelihood of the observations and latent variables. Let 121

$$g_i(\boldsymbol{\theta}) := \int_{\mathbf{7}} f_i(z_i, \boldsymbol{\theta}) \mu_i(\mathrm{d}z_i), \ i \in [\![1, n]\!], \ \boldsymbol{\theta} \in \Theta$$

denote the incomplete data likelihood, i.e., the marginal likelihood of the observations y_i . For ease 122 of notations, the dependence on the observations is made implicit. The maximum likelihood (ML) 123 estimation problem sets the individual objective function $\mathcal{L}_i(\theta)$ to be the *i*-th negated incomplete data log-likelihood $\mathcal{L}_i(\boldsymbol{\theta}) := -\log g_i(\boldsymbol{\theta})$.

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 $p_i(\cdot; \boldsymbol{\theta}^{(0)})$.
- 3: Initialize the surrogate function as

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

- 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 $p_i(\cdot; \boldsymbol{\theta}^{(k)})$.
- 7: Update the individual surrogate functions recursively as:

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

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

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 observations y_i . A surrogate function $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ satisfying H1 can be obtained through writing $f_i(z_i, \theta) = \frac{f_i(z_i, \theta)}{p_i(z_i, \overline{\theta})} p_i(z_i, \overline{\theta})$ and applying the Jensen inequality:

$$\widehat{\mathcal{L}}_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) = \int_{\mathsf{Z}} \underbrace{\log \left(p_{i}(z_{i}, \overline{\boldsymbol{\theta}}) / f_{i}(z_{i}, \boldsymbol{\theta}) \right)}_{=r_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, z_{i})} p_{i}(z_{i}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i}) . \tag{7}$$

We note that H2 can also be verified for common distribution models. We can apply the MISSO method following the above specification of $r_i(\theta; \overline{\theta}, z_i)$ and $p_i(z_i, \overline{\theta})$.

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

$$p(y, w|x) = \pi(w) \prod_{i=1}^{n} p(y_i|x_i, w)$$
 (8)

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 Kullback-Leibler (KL) divergence between $q(w, \theta)$ and the posterior distribution p(w|y,x):

$$\min_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}) := \mathrm{KL}\left(q(w; \boldsymbol{\theta}) || p(w|y, x)\right) := \mathbb{E}_{q(w; \boldsymbol{\theta})}\left[\log\left(q(w; \boldsymbol{\theta}) / p(w|y, x)\right)\right]. \tag{9}$$

Using (8), we decompose $\mathcal{L}(\boldsymbol{\theta}) = n^{-1} \sum_{i=1}^{n} \mathcal{L}_i(\boldsymbol{\theta}) + \text{const.}$ where:

$$\mathcal{L}_{i}(\boldsymbol{\theta}) := -\mathbb{E}_{q(w;\boldsymbol{\theta})} \left[\log p(y_{i}|x_{i},w) \right] + \frac{1}{n} \mathbb{E}_{q(w;\boldsymbol{\theta})} \left[\log q(w;\boldsymbol{\theta})/\pi(w) \right] := r_{i}(\boldsymbol{\theta}) + d(\boldsymbol{\theta}) . \quad (10)$$

Directly optimizing the finite sum objective function in (9) can be difficult. First, with $n\gg 1$, evaluating the objective function $\mathcal{L}(\boldsymbol{\theta})$ requires a full pass over the entire dataset. Second, for some complex models, the expectations in (10) can be intractable even if we assume a simple parametric model for $q(w;\boldsymbol{\theta})$. Assume that \mathcal{L}_i is L-smooth. We apply the MISSO method with a quadratic surrogate function defined as:

$$\widehat{\mathcal{L}}_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) := \mathcal{L}_{i}(\overline{\boldsymbol{\theta}}) + \left\langle \nabla_{\boldsymbol{\theta}} \mathcal{L}_{i}(\overline{\boldsymbol{\theta}}) \,|\, \boldsymbol{\theta} - \overline{\boldsymbol{\theta}} \right\rangle + \frac{L}{2} \|\overline{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^{2}, \ (\boldsymbol{\theta}, \overline{\boldsymbol{\theta}}) \in \Theta^{2}. \tag{11}$$

It is easily checked that the quadratic function $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ satisfies H1, H2. To compute the gradient $\nabla \mathcal{L}_i(\overline{\theta})$, 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. $\theta \in \Theta$

which is designed such that the law of $w=t(z,\overline{\theta})$ is $q(\cdot,\overline{\theta})$, where $z\sim\mathcal{N}_d(0,\mathbf{I})$. By [Blundell et al., 2015, Proposition 1], the gradient of $-r_i(\cdot)$ in (10) is:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(w;\overline{\boldsymbol{\theta}})} \left[\log p(y_i|x_i, w) \right] = \mathbb{E}_{z \sim \mathcal{N}_d(0, \mathbf{I})} \left[J_{\boldsymbol{\theta}}^t(z, \overline{\boldsymbol{\theta}}) \nabla_w \log p(y_i|x_i, w) \Big|_{w = t(z, \overline{\boldsymbol{\theta}})} \right], \tag{12}$$

where for each $z \in \mathbb{R}^d$, $J^t_{\theta}(z, \overline{\theta})$ is the Jacobian of the function $t(z, \cdot)$ with respect to θ evaluated at $\overline{\theta}$. In addition, for most cases, the term $\nabla d(\overline{\theta})$ can be evaluated in closed form as the gradient of the KL between the prior distribution $\pi(\cdot)$ and the variational candidate $q(\cdot, \theta)$.

$$r_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, z) := \left\langle \nabla_{\boldsymbol{\theta}} d(\overline{\boldsymbol{\theta}}) - J_{\boldsymbol{\theta}}^{t}(z, \overline{\boldsymbol{\theta}}) \nabla_{w} \log p(y_{i}|x_{i}, w) \big|_{w=t(z, \overline{\boldsymbol{\theta}})} | \boldsymbol{\theta} - \overline{\boldsymbol{\theta}} \right\rangle + \frac{L}{2} \|\boldsymbol{\theta} - \overline{\boldsymbol{\theta}}\|^{2} . \quad (13)$$

Finally, using (11) and (13), the surrogate function (6) is given by $\widetilde{\mathcal{L}}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, \{z_m\}_{m=1}^M) := M^{-1} \sum_{m=1}^M r_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, z_m)$ where $\{z_m\}_{m=1}^M$ are i.i.d samples drawn from $\mathcal{N}(0, \mathbf{I})$.

156 3 Convergence Analysis

162

157 We now provide asymptotic and non-asymptotic convergence results our method. Assume:

158 **H3.** For all
$$i \in [1, n]$$
, $\overline{\theta} \in \Theta$, $z_i \in \mathbb{Z}$, $r_i(\cdot; \overline{\theta}, z_i)$ is convex on Θ and is lower bounded.

- We are particularly interested in the *constrained optimization* setting where Θ is a bounded set. To this end, we control the supremum norm of the MC approximation, introduced in (6), as:
- 161 **H4.** For the samples $\{z_{i,m}\}_{m=1}^{M}$, there exist finite constants C_r and C_{gr} such that

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

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

for all $i \in [\![1,n]\!]$, and 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 $p_i(\cdot;\overline{\theta})$.

Some intuitions behind the controlling terms: It is common in statistical and optimization problems, to deal with the manipulation and the control of random variables indexed by sets with an infinite number of elements. Here, the controlled random variable is an image of a continuous function defined as $r_i(\theta; \overline{\theta}, z_{i,m}) - \widehat{\mathcal{L}}_i(\theta; \overline{\theta})$ for all $z \in \mathsf{Z}$ and for fixed $(\theta, \overline{\theta}) \in \Theta^2$. To characterize such control, we will have recourse to the notion of metric entropy (or bracketing number) as developed in [Van der Vaart, 2000, Vershynin, 2018, Wainwright, 2019]. A collection of results from those references gives intuition behind our assumption H4, which is classical in empirical processes. In [Vershynin, 2018, Theorem 8.2.3], the authors recall the uniform law of large numbers:

$$\mathbb{E}\left[\sup_{f\in\mathcal{F}}\left|\frac{1}{M}\sum_{i=1}^{M}f\left(z_{i,m}\right)-\mathbb{E}[f(z_{i})]\right|\right]\leq\frac{CL}{\sqrt{M}}\quad\text{for all}\quad z_{i,m},i\in\left[1,M\right],$$

where \mathcal{F} is a class of L-Lipschitz functions. Moreover, in [Vershynin, 2018, Theorem 8.1.3] and [Wainwright, 2019, Theorem 5.22], the application of the Dudley inequality yields:

$$\mathbb{E}[\sup_{f \in \mathcal{F}} |X_f - X_0|] \le \frac{1}{\sqrt{M}} \int_0^1 \sqrt{\log \mathcal{N}(\mathcal{F}, \|\cdot\|_{\infty}, \varepsilon)} d\varepsilon ,$$

where $\mathcal{N}\left(\mathcal{F}, \|\cdot\|_{\infty}, \varepsilon\right)$ is the bracketing number and ϵ denotes the level of approximation (the bracketing number goes to infinity when $\epsilon \to 0$). Finally, in [Van der Vaart, 2000, p.271, Example], $\mathcal{N}\left(\mathcal{F}, \|\cdot\|_{\infty}, \varepsilon\right)$ is bounded from above for a class of parametric functions $\mathcal{F} = f_{\theta} : \theta \in \Theta$:

$$\mathcal{N}(\mathcal{F}, \|\cdot\|_{\infty}, \varepsilon) \leq K \left(\frac{\operatorname{diam}\Theta}{\varepsilon}\right)^d$$
, every $0 < \varepsilon < \operatorname{diam}\Theta$.

The authors acknowledge that those bounds are a dramatic manifestation of the curse of dimensionality happening when sampling is needed. Nevertheless, the dependence on the dimension highly depends on the class of surrogate functions \mathcal{F} used in our scheme, as smaller bounds on these controlling terms can be derived for simpler class of functions, such as quadratic functions.

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

$$g(\overline{\boldsymbol{\theta}}) := \inf_{\boldsymbol{\theta} \in \Theta} \frac{\mathcal{L}'(\overline{\boldsymbol{\theta}}, \boldsymbol{\theta} - \overline{\boldsymbol{\theta}})}{\|\overline{\boldsymbol{\theta}} - \boldsymbol{\theta}\|} \quad \text{and} \quad g(\overline{\boldsymbol{\theta}}) = g_{+}(\overline{\boldsymbol{\theta}}) - g_{-}(\overline{\boldsymbol{\theta}}) , \tag{14}$$

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

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, $\tau_{i_k}^{k+1}=k$ for instance. We define:

$$\widehat{\mathcal{L}}^{(k)}(\theta) := \frac{1}{n} \sum_{i=1}^{n} \widehat{\mathcal{L}}_{i}(\theta; \theta^{(\tau_{i}^{k})}), \quad \widehat{e}^{(k)}(\theta) := \widehat{\mathcal{L}}^{(k)}(\theta) - \mathcal{L}(\theta), \quad \overline{M}_{(k)} := \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2} . \quad (15)$$

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

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

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

Then we have following non-asymptotic bounds:

$$\mathbb{E}\big[\|\nabla \widehat{e}^{(K)}(\boldsymbol{\theta}^{(K)})\|^2\big] \leq \frac{\Delta_{(K_{\text{max}})}}{K_{\text{max}}} \quad and \quad \mathbb{E}[g_{-}(\boldsymbol{\theta}^{(K)})] \leq \sqrt{\frac{\Delta_{(K_{\text{max}})}}{K_{\text{max}}}} + \frac{C_{\text{gr}}}{K_{\text{max}}} \overline{M}_{(k)} \ . \tag{16}$$

Note that $\Delta_{(K_{\text{max}})}$ is finite for any $K_{\text{max}} \in \mathbb{N}$. As expected, the MISSO method converges to a stationary point of (1) asymptotically and at a sublinear rate $\mathbb{E}[g_-^{(K)}] \leq \mathcal{O}(\sqrt{1/K_{\text{max}}})$. 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. H4], Eq. (16) gives a non-asymptotic rate of $\mathbb{E}[g_-^{(K)}] \leq \mathcal{O}(\sqrt{nL/K_{\text{max}}})$ 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 H1-H4. 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 a.s. to zero, i.e., $\lim_{k\to\infty} g_-(\boldsymbol{\theta}^{(k)}) \stackrel{a.s.}{=} 0$.
- 2. the objective value $\mathcal{L}(\boldsymbol{\theta}^{(k)})$ converges a.s. to a finite number $\underline{\mathcal{L}}$, i.e., $\lim_{k\to\infty} \mathcal{L}(\boldsymbol{\theta}^{(k)}) \stackrel{a.s.}{=} \underline{\mathcal{L}}$.

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

201

202

4.1 Binary logistic regression with missing values

This application follows **Example 1** described in Section 2. We consider a binary regression setup, $((y_i, z_i), i \in [\![n]\!])$ where $y_i \in \{0, 1\}$ is a binary response and $z_i = (z_{i,j} \in \mathbb{R}, j \in [\![p]\!])$ is a covariate

vector. The vector of covariates $z_i = [z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}]$ is not fully observed where we denote by $z_{i,\mathrm{mis}}$ the missing values and $z_{i,\mathrm{obs}}$ the observed covariate. It is assumed that $(z_i, i \in [n])$ are i.i.d. and marginally distributed according to $\mathcal{N}(\boldsymbol{\beta}, \boldsymbol{\Omega})$ where $\boldsymbol{\beta} \in \mathbb{R}^p$ and $\boldsymbol{\Omega}$ is a positive definite $p \times p$ matrix. We define the conditional distribution of the observations y_i given $z_i = (z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}})$ as:

$$p_i(y_i|z_i) = S(\boldsymbol{\delta}^{\top}\bar{z}_i)^{y_i} \left(1 - S(\boldsymbol{\delta}^{\top}\bar{z}_i)\right)^{1 - y_i} , \qquad (17)$$

where for $u \in \mathbb{R}$, $S(u) = 1/(1 + \mathrm{e}^{-u})$, $\boldsymbol{\delta} = (\delta_0, \dots, \delta_p)$ are the logistic parameters and $\bar{z}_i = (1, z_i)$.

Here, $\boldsymbol{\theta} = (\boldsymbol{\delta}, \boldsymbol{\beta}, \boldsymbol{\Omega})$ is the parameter to estimate. For $i \in [n]$, the complete log-likelihood reads:

$$\log f_i(z_{i,\mathrm{mis}}, \boldsymbol{\theta}) \propto y_i \boldsymbol{\delta}^\top \bar{z}_i - \log \left(1 + \exp(\boldsymbol{\delta}^\top \bar{z}_i) \right) - \frac{1}{2} \log(|\boldsymbol{\Omega}|) + \frac{1}{2} \mathrm{Tr} \left(\boldsymbol{\Omega}^{-1} (z_i - \boldsymbol{\beta}) (z_i - \boldsymbol{\beta})^\top \right).$$

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. This dataset includes information from the first stage of the trauma, namely initial observations on the patient's accident site to the last stage being intense care at the hospital and counts more than 200 variables measured for more than 7 000 patients. Since the dataset considered is heterogeneous – coming from multiple sources with frequently missed entries – we apply the latent data model described in (17) to predict the risk of a severe hemorrhage which is one of the main cause of death after a major trauma.

Similar to [Jiang et al., 2018], we select p=16 influential quantitative measurements, on n=6384 patients. For the Monte Carlo sampling of $z_{i,\text{mis}}$, required while running MISSO, we run a Metropolis-Hastings algorithm with the target distribution $p(\cdot|z_{i,\text{obs}},y_i;\boldsymbol{\theta}^{(k)})$.

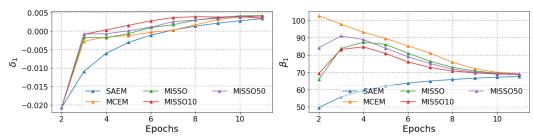


Figure 1: Convergence of first component of the vector of parameters δ and β for the SAEM, the MCEM and the MISSO methods. The convergence is plotted against No. of passes over the data.

We compare in Figure 1 the convergence behavior of the estimated parameters δ and 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 (MISSO), 10% (MISSO10) and 50% (MISSO50) minibatches 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 Training Bayesian CNN using MISSO

This application follows **Example 2** described in Section 2. We use variational inference and the ELBO loss (10) to fit Bayesian Neural Networks on different datasets. At iteration k, minimizing the sum of stochastic surrogates defined as in (6) and (13) yields the following MISSO update — step (i) pick a function index i_k uniformly on [n]; 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, with $\tilde{w}=t(\boldsymbol{\theta}^{(k-1)},z_m^{(k)})$, as

$$\mu_{\ell}^{(k)} = \hat{\mu}_{\ell}^{(\tau^k)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\mu_{\ell},i}^{(k)} \quad \text{and} \quad \hat{\delta}_{\mu_{\ell},i_k}^{(k)} = -\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} \nabla_w \log p(y_{i_k}|x_{i_k}, \tilde{w}) + \nabla_{\mu_{\ell}} d(\boldsymbol{\theta}^{(k-1)}) \;,$$

where
$$\hat{\mu}_{\ell}^{(\tau^k)} = \frac{1}{n} \sum_{i=1}^n \mu_{\ell}^{(\tau_i^k)}$$
 and $d(\boldsymbol{\theta}) = n^{-1} \sum_{\ell=1}^d \left(-\log(\sigma) + (\sigma^2 + \mu_{\ell}^2)/2 - 1/2 \right)$.

Bayesian LeNet-5 on MNIST [LeCun et al., 1998]: We apply the MISSO method to fit a Bayesian variant of LeNet-5 [LeCun et al., 1998]. 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 (8), 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 $z \sim \mathcal{N}(0, \mathbf{I})$.

Bayesian ResNet-18 [He et al., 2016] on CIFAR-10 [Krizhevsky et al., 2012]: We train here the Bayesian variant of the ResNet-18 neural network introduced in [He et al., 2016] on CIFAR-10. The latter dataset is composed of $n=60\,000$ handwritten digits, 32×32 colour images in 10 classes, with $6\,000$ images per class. As in the previous example, the weights are assumed independent and identically distributed according to $\mathcal{N}(0,\mathbf{I})$. Standard hyperparameters values found in the literature, such as the annealing constant or the number of MC samples, were used for the benchmark methods. For better efficiency and lower variance, the Flipout estimator [Wen et al., 2018] is used.

Experiment Results: 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 (10) and its gradients were computed by Monte Carlo integration based on the reparametrization described above. The mini-batch of indices and MC samples are respectively set to 128 and $M_{(k)} = k$. The learning rates are set to 10^{-3} for LeNet-5 and 10^{-4} for Resnet-18.

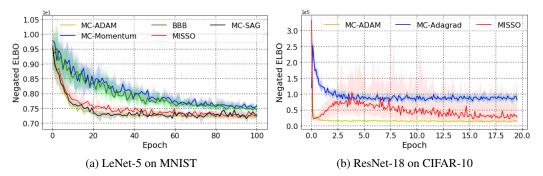


Figure 2: Negated ELBO versus epochs elapsed for fitting (a) Bayesian LeNet-5 on MNIST and (b) Bayesian ResNet-18 on CIFAR-10. The solid curve is obtained from averaging over 5 independent runs of the methods, and the shaded area represents the standard deviation.

Figure 2(a) 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 for our experiment on MNIST dataset using a Bayesian variant of LeNet-5. On the other hand, the experiment conducted on CIFAR-10 (Figure 2(b)) using a much larger network, *i.e.*, a Bayesian variant of ResNet-18 showcases the need of a well-tuned adaptive methods to reach better training loss (and also faster). Our MISSO method is similar to the Monte Carlo variant of ADAM but slower than Adagrad optimizer. Recall that the purpose of this paper is to provide a common class of optimizers, such as VI, in order to study their convergence behaviors, and not to introduce a novel method outperforming the baselines methods.

5 Conclusion

We present a unifying framework for minimizing a nonconvex and nonsmooth 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 nonconvex 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 Bayesian variants of two Convolutional Neural Networks on benchmark datasets.

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.
- 286 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.
- 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.
- R. Fletcher, N. I. Gould, S. Leyffer, P. L. Toint, and A. Wächter. Global convergence of a trustregion sqp-filter algorithm for general nonlinear programming. *SIAM Journal on Optimization*, 13(3):635–659, 2002.
- S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 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.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings* of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- 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.
- 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.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- K. Lange. *MM Optimization Algorithms*. SIAM-Society for Industrial and Applied Mathematics,
 USA, 2016. ISBN 1611974399, 9781611974393.
- 23 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.
- X. Qian, A. Sailanbayev, K. Mishchenko, and P. Richtárik. Miso is making a comeback with better proofs and rates. *arXiv preprint arXiv:1906.01474*, 2019.
- 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.
- 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.
- A. W. Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.
- R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- M. J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- 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.
- Y. Wen, P. Vicol, J. Ba, D. Tran, and R. Grosse. Flipout: Efficient pseudo-independent weight perturbations on mini-batches. *arXiv preprint arXiv:1803.04386*, 2018.

67 A Proof of Theorem 1

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

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

370 Then we have following non-asymptotic bounds:

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

Proof We begin by recalling the definition

$$\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^{n} \widetilde{\mathcal{A}}_{i}^{k}(\boldsymbol{\theta}).$$

372 Notice that

$$\begin{split} \widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^{n} \widetilde{\mathcal{L}}_{i}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(\tau_{i}^{k+1})}, \{z_{i,m}^{(\tau_{i}^{k+1})}\}_{m=1}^{M_{(\tau_{i}^{k+1})}}) \\ &= \widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) + \frac{1}{n} \big(\widetilde{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}, \{z_{i_{k},m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widetilde{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(\tau_{i_{k}}^{k})}, \{z_{i_{k},m}^{(\tau_{i_{k}}^{k})}\}_{m=1}^{M_{(\tau_{i_{k}}^{k})}}) \big). \end{split}$$

Furthermore, we recall that

$$\widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^{n} \widehat{\mathcal{L}}_{i}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(\tau_{i}^{k})}), \quad \widehat{e}^{(k)}(\boldsymbol{\theta}) := \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}).$$

Due to H2, we have

$$\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2 \le 2L\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}). \tag{18}$$

To prove the first bound in (16), using the optimality of $\theta^{(k+1)}$, one has

$$\widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) \leq \widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k)}) \\
= \widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) + \frac{1}{n} (\widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k,m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k,m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}}))$$
(19)

Let \mathcal{F}_k be the filtration of random variables up to iteration k, i.e., $\{i_{\ell-1},\{z_{i_{\ell-1},m}^{(\ell-1)}\}_{m=1}^{M_{(\ell-1)}},m{ heta}^{(\ell)}\}_{\ell=1}^k$.

We observe that the conditional expectation evaluates to

$$\begin{split} & \mathbb{E}_{i_k} \left[\mathbb{E} \big[\widetilde{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_k,m}^{(k)}\}_{m=1}^{M_{(k)}}) | \mathcal{F}_k, i_k \big] | \mathcal{F}_k \right] \\ & = \mathcal{L}(\boldsymbol{\theta}^{(k)}) + \mathbb{E}_{i_k} \big[\mathbb{E} \big[\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} r_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, z_{i_k,m}^{(k)}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}) | \mathcal{F}_k, i_k \big] | \mathcal{F}_k \big] \\ & \leq \mathcal{L}(\boldsymbol{\theta}^{(k)}) + \frac{C_{\mathsf{r}}}{\sqrt{M_{(k)}}}, \end{split}$$

where the last inequality is due to H4. Moreover,

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

Taking the conditional expectations on both sides of (19) and re-arranging terms give:

$$\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)}) \le n \mathbb{E} \left[\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) | \mathcal{F}_k \right] + \frac{C_{\mathsf{r}}}{\sqrt{M_{(k)}}}$$
(20)

Proceeding from (20), we observe the following lower bound for the left hand side

$$\begin{split} &\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)}) \stackrel{(a)}{=} \widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) + \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) \\ &\overset{(b)}{\geq} \widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) + \frac{1}{2L} \|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2 \\ &= \underbrace{\frac{1}{n} \sum_{i=1}^{n} \left\{ \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} r_i(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}, z_{i,m}^{(\tau_i^k)}) - \widehat{\mathcal{L}}_i(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_i^k)}) \right\}}_{:=-\delta^{(k)}(\boldsymbol{\theta}^{(k)})} + \frac{1}{2L} \|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2 \end{split}}$$

where (a) is due to $\hat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) = 0$ [cf. H1], (b) is due to (18) and we have defined the summation in the last equality as $-\delta^{(k)}(\boldsymbol{\theta}^{(k)})$. Substituting the above into (20) yields

$$\frac{\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2}{2L} \le n\mathbb{E}\left[\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)})|\mathcal{F}_k\right] + \frac{C_{\mathsf{r}}}{\sqrt{M_{(k)}}} + \delta^{(k)}(\boldsymbol{\theta}^{(k)}) \tag{21}$$

Observe the following upper bound on the total expectations:

$$\mathbb{E}\big[\delta^{(k)}(\boldsymbol{\theta}^{(k)})\big] \leq \mathbb{E}\Big[\frac{1}{n}\sum_{i=1}^n \frac{C_{\mathsf{r}}}{\sqrt{M_{(\tau_i^k)}}}\Big],$$

which is due to H4. It yields

$$\mathbb{E}\big[\|\nabla\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2\big] \leq 2nL\mathbb{E}\big[\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)})\big] + \frac{2LC_{\mathsf{r}}}{\sqrt{M_{(k)}}} + \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\Big[\frac{2LC_{\mathsf{r}}}{\sqrt{M_{(\tau_i^k)}}}\Big]$$

Finally, for any $K_{\text{max}} \in \mathbb{N}$, we let K be a discrete r.v. that is uniformly drawn from $\{0, 1, ..., K_{\text{max}} - 1\}$. Using H4 and taking total expectations lead to

$$\mathbb{E}[\|\nabla \widehat{e}^{(K)}(\boldsymbol{\theta}^{(K)})\|^{2}] = \frac{1}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} \mathbb{E}[\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\|^{2}] \\
\leq \frac{2nL\mathbb{E}[\widetilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \widetilde{\mathcal{L}}^{(K_{\text{max}})}(\boldsymbol{\theta}^{(K_{\text{max}})})]}{K_{\text{max}}} + \frac{2LC_{\text{r}}}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} \mathbb{E}\Big[\frac{1}{\sqrt{M_{(k)}}} + \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\sqrt{M_{(\tau_{i}^{k})}}}\Big]$$
(22)

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:

$$\mathbb{E}[M_{(\tau_i^k)}^{-1/2}] = \sum_{j=1}^k \frac{1}{n} \left(1 - \frac{1}{n}\right)^{j-1} M_{(k-j)}^{-1/2}$$
(23)

389 Taking the sum yields:

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

$$(24)$$

where the last inequality is due to upper bounding the geometric series. Plugging this back into (22) yields

$$\begin{split} & \mathbb{E} \big[\| \nabla \widehat{e}^{(K)}(\boldsymbol{\theta}^{(K)}) \|^2 \big] = \frac{1}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} \mathbb{E} [\| \nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) \|^2] \\ & \leq \frac{2nL \mathbb{E} [\widetilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \widetilde{\mathcal{L}}^{(K_{\text{max}})}(\boldsymbol{\theta}^{(K_{\text{max}})})]}{K_{\text{max}}} + \frac{1}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} \frac{4LC_{\text{r}}}{\sqrt{M_{(k)}}} = \frac{\Delta_{(K_{\text{max}})}}{K_{\text{max}}}. \end{split}$$

- This concludes our proof for the first inequality in (16).
- To prove the second inequality of (16), we define the shorthand notations $g^{(k)}:=g(\pmb{\theta}^{(k)}),\,g_-^{(k)}:=g(\pmb{\theta}^{(k)})$
- 394 $-\min\{0, g^{(k)}\}, g_+^{(k)} := \max\{0, g^{(k)}\}.$ We observe that

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

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

$$\begin{split} &\frac{1}{n}\sum_{i=1}^{n}\widehat{\mathcal{L}}_{i}^{'}(\boldsymbol{\theta}^{(k)},\boldsymbol{\theta}-\boldsymbol{\theta}^{(k)};\boldsymbol{\theta}^{(\tau_{i}^{k})})\\ &=\underbrace{\widetilde{\mathcal{L}}^{(k)'}(\boldsymbol{\theta}^{(k)},\boldsymbol{\theta}-\boldsymbol{\theta}^{(k)})}_{\geq 0}-\widehat{\mathcal{L}}^{(k)'}(\boldsymbol{\theta}^{(k)},\boldsymbol{\theta}-\boldsymbol{\theta}^{(k)})+\frac{1}{n}\sum_{i=1}^{n}\widehat{\mathcal{L}}_{i}^{'}(\boldsymbol{\theta}^{(k)},\boldsymbol{\theta}-\boldsymbol{\theta}^{(k)};\boldsymbol{\theta}^{(\tau_{i}^{k})})\\ &\geq\frac{1}{n}\sum_{i=1}^{n}\left\{\widehat{\mathcal{L}}_{i}^{'}(\boldsymbol{\theta}^{(k)},\boldsymbol{\theta}-\boldsymbol{\theta}^{(k)};\boldsymbol{\theta}^{(\tau_{i}^{k})})-\frac{1}{M_{(\tau_{i}^{k})}}\sum_{m=1}^{M_{(\tau_{i}^{k})}}r_{i}^{'}(\boldsymbol{\theta}^{(k)},\boldsymbol{\theta}-\boldsymbol{\theta}^{(k)};\boldsymbol{\theta}^{(\tau_{i}^{k})},\boldsymbol{z}_{i,m}^{(\tau_{i}^{k})})\right\} \end{split}$$

where the inequality is due to the optimality of $\theta^{(k)}$ and the convexity of $\widetilde{\mathcal{L}}^{(k)}(\theta)$ [cf. H3]. Denoting a scaled version of the above term as:

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

400 We have

$$g^{(k)} \ge -\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| + \inf_{\boldsymbol{\theta} \in \Theta} (-\epsilon^{(k)}(\boldsymbol{\theta})) \ge -\|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| - \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})|. \tag{25}$$

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

$$g_{-}^{(k)} \le \|\nabla \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| + \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})|. \tag{26}$$

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

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

404 We note that

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

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

$$\begin{split} \mathbb{E}[\sup_{\pmb{\theta} \in \Theta} \epsilon^{(K)}(\pmb{\theta})] &= \frac{1}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}} \mathbb{E}[\sup_{\pmb{\theta} \in \Theta} \epsilon^{(k)}(\pmb{\theta})] \overset{(a)}{\leq} \frac{C_{\text{gr}}}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} \mathbb{E}\Big[\frac{1}{n} \sum_{i=1}^{n} M_{(\tau_i^k)}^{-1/2}\Big] \\ &\overset{(b)}{\leq} \frac{C_{\text{gr}}}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} M_{(k)}^{-1/2} \end{split}$$

where (a) is due to H4 and (b) is due to (24). This implies

$$\mathbb{E}[\boldsymbol{g}_{-}^{(K)}] \leq \sqrt{\frac{\Delta_{(K_{\text{max}})}}{K_{\text{max}}}} + \frac{C_{\text{gr}}}{K_{\text{max}}} \sum_{k=0}^{K_{\text{max}}-1} M_{(k)}^{-1/2},$$

and concludes the proof of the theorem.

408 B Proof of Theorem 2

- Theorem. Under H1-H4. 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., $\lim_{k\to\infty} g_-(\theta^{(k)}) = 0$ a.s..
- 2. the objective value $\mathcal{L}(\boldsymbol{\theta}^{(k)})$ converges almost surely to a finite number $\underline{\mathcal{L}}$, i.e., $\lim_{k\to\infty}\mathcal{L}(\boldsymbol{\theta}^{(k)})=\underline{\mathcal{L}}$ 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 $(V_k)_{k>0}$ 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
- 419 variables such that $\sum_{k=0}^{\infty} \mathbb{E}[|E_k|] < \infty$. If for any $k \ge 1$:

$$V_k \le V_{k-1} - X_{k-1} + E_{k-1} \tag{27}$$

- 420 *then*:
- (i) for all $k \geq 0$, $\mathbb{E}[V_k] < \infty$ and the sequence $(V_k)_{k>0}$ converges a.s. to a finite limit V_{∞} .
- (ii) the sequence $(\mathbb{E}[V_k])_{k\geq 0}$ converges and $\lim_{k\to\infty}\mathbb{E}[V_k]=\mathbb{E}[V_\infty]$.
- (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 (19) by re-arranging terms and observing that

$$\begin{split} \widehat{\mathcal{L}}^{(k+1)}(\pmb{\theta}^{(k+1)}) & \leq \widehat{\mathcal{L}}^{(k)}(\pmb{\theta}^{(k)}) - \frac{1}{n} \big(\widehat{\mathcal{L}}_{i_k}(\pmb{\theta}^{(k)}; \pmb{\theta}^{(\tau_{i_k}^k)}) - \widehat{\mathcal{L}}_{i_k}(\pmb{\theta}^{(k)}; \pmb{\theta}^{(k)}) \big) \\ & - \big(\widetilde{\mathcal{L}}^{(k+1)}(\pmb{\theta}^{(k+1)}) - \widehat{\mathcal{L}}^{(k+1)}(\pmb{\theta}^{(k+1)}) \big) + \big(\widetilde{\mathcal{L}}^{(k)}(\pmb{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\pmb{\theta}^{(k)}) \big) \\ & + \frac{1}{n} \big(\widetilde{\mathcal{L}}_{i_k}(\pmb{\theta}^{(k)}; \pmb{\theta}^{(k)}, \{z_{i_k,m}^k\}_{m=1}^{M_{(k)}}) - \widehat{\mathcal{L}}_{i_k}(\pmb{\theta}^{(k)}; \pmb{\theta}^{(\kappa)}; \pmb{\theta}^{(\kappa)}) \big) \\ & + \frac{1}{n} \big(\widehat{\mathcal{L}}_{i_k}(\pmb{\theta}^{(k)}; \pmb{\theta}^{(\tau_{i_k}^k)}) - \widetilde{\mathcal{L}}_{i_k}(\pmb{\theta}^{(k)}; \pmb{\theta}^{(\tau_{i_k}^k)}, \{z_{i_k,m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}} \big) \big) \end{split}$$

Our idea is to apply Lemma 1. Under H1, the finite sum of surrogate functions $\widehat{\mathcal{L}}^{(k)}(\theta)$, defined in (15), is lower bounded by a constant $c_k > -\infty$ for any θ . To this end, we observe that

$$V_k := \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \inf_{k \ge 0} c_k \ge 0$$
(28)

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

$$X_k := \frac{1}{n} \left(\widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(\tau_{i_k}^k)}; \boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}_{i_k}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}) \right) \ge 0.$$
 (29)

429 Thirdly, we define

$$E_{k} = -\left(\widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) - \widehat{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)})\right) + \left(\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\right) + \frac{1}{n}\left(\widetilde{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}, \{z_{i_{k},m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widehat{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)})\right) + \frac{1}{n}\left(\widehat{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_{k}}^{k})}) - \widetilde{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(\tau_{i_{k}}^{k})}, \{z_{i_{k},m}^{(\tau_{i_{k}}^{k})}\}_{m=1}^{M_{(\tau_{i_{k}}^{k})}})\right).$$
(30)

- Note that from the definitions (28), (29), (30), we have $V_{k+1} \leq V_k X_k + E_k$ for any $k \geq 1$.
- 431 Under H4, we observe that

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

432

433

$$\mathbb{E}\Big[\Big|\widehat{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}^{(k)};\boldsymbol{\theta}^{(\tau_{i_{k}}^{k})}) - \widetilde{\mathcal{L}}_{i_{k}}(\boldsymbol{\theta}^{(k)};\boldsymbol{\theta}^{(\tau_{i_{k}}^{k})},\{z_{i_{k},m}^{(\tau_{i_{k}}^{k})}\}_{m=1}^{M_{(\tau_{i_{k}}^{k})}})\Big|\Big] \leq C_{\mathsf{r}}\mathbb{E}\Big[M_{(\tau_{i_{k}}^{k})}^{-1/2}\Big]$$

$$\mathbb{E}\Big[|\widetilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})|\Big] \leq \frac{1}{n}\sum_{i=1}^{n}C_{\mathsf{r}}\mathbb{E}\Big[M_{(\tau_{i_{k}}^{k})}^{-1/2}\Big]$$

434 Therefore,

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

Using (24) and the assumption on the sequence $\{M_{(k)}\}_{k>0}$, we obtain that

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

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

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

$$\lim_{k \to \infty} \widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)}) = 0 \quad \text{a.s.}$$
 (31)

and subsequently applying (18), we have $\lim_{k\to\infty} \|\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})\| = 0$ almost surely. Finally, it follows from (18) and (26) that 440

$$\lim_{k \to \infty} g_{-}^{(k)} \le \lim_{k \to \infty} \sqrt{2L} \sqrt{\widehat{e}^{(k)}(\boldsymbol{\theta}^{(k)})} + \lim_{k \to \infty} \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})| = 0, \tag{32}$$

- where the last equality holds almost surely due to the fact that $\sum_{k=0}^{\infty} \mathbb{E}[\sup_{\theta \in \Theta} |\epsilon^{(k)}(\theta)|] < \infty$. This concludes the asymptotic convergence of the MISSO method. 441
- 442
- Finally, we prove that $\mathcal{L}(\theta^{(k)})$ converges almost surely. As a consequence of Lemma 1, it is clear that 443
- $\{V_k\}_{k\geq 0}$ converges almost surely and so is $\{\widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\}_{k\geq 0}$, i.e., we have $\lim_{k\to\infty}\widehat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})=\underline{\mathcal{L}}$. 444
- Applying (31) implies that 445

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

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

Proof of Lemma 1 447

- **Lemma.** Let $(V_k)_{k>0}$ be a non negative sequence of random variables such that $\mathbb{E}[V_0] < \infty$. 448
- 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 $k \geq 1$: 449
- 450

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

- then: 451
- (i) for all $k \geq 0$, $\mathbb{E}[V_k] < \infty$ and the sequence $(V_k)_{k>0}$ converges a.s. to a finite limit V_{∞} . 452
- (ii) the sequence $(\mathbb{E}[V_k])_{k>0}$ converges and $\lim_{k\to\infty}\mathbb{E}[V_k]=\mathbb{E}[V_\infty]$. 453
- (iii) the series $\sum_{k=0}^{\infty} X_k$ converges almost surely and $\sum_{k=0}^{\infty} \mathbb{E}[X_k] < \infty$. 454

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

$$0 \le V_k \le V_0 - \sum_{j=1}^k X_j + \sum_{j=1}^k E_j \le V_0 + \sum_{j=1}^k E_j$$
(33)

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

Since $0 \le X_k \le V_{k-1} - V_k + E_k$ we also obtain for all $k \ge 0$, $\mathbb{E}[X_k] < \infty$. 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:

$$W_k = V_k + \sum_{j=k+1}^{\infty} E_j \tag{34}$$

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:

$$W_{k} \leq V_{k-1} - X_{k} + \sum_{j=k}^{\infty} E_{j} \leq W_{k-1} - X_{k} \leq W_{k-1}$$

$$\mathbb{E}[W_{k}] \leq \mathbb{E}[W_{k-1}] - \mathbb{E}[X_{k}]$$
(35)

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

$$\mathbb{E}[\liminf_{k \to \infty} |W_k|] = \mathbb{E}[|W_\infty|] \le \liminf_{k \to \infty} \mathbb{E}[|W_k|] \le \mathbb{E}[V_0] + \sum_{j=1}^{\infty} \mathbb{E}[|E_j|] < \infty$$
 (36)

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:

$$\lim_{k \to \infty} \mathbb{E}[U_k] = \mathbb{E}[\lim_{k \to \infty} U_k]$$
(37)

467 Finally, we have:

$$\lim_{k \to \infty} \mathbb{E}[U_k] = \mathbb{E}[W_0] - w_{\infty} \quad \text{and} \quad \mathbb{E}[\lim_{k \to \infty} U_k] = \mathbb{E}[W_0] - \mathbb{E}[W_{\infty}]$$
 (38)

showing that $\mathbb{E}[W_{\infty}] = w_{\infty}$ and concluding the proof of (ii). Moreover, using (35) we have that $W_k \leq W_{k-1} - X_k$ which yields:

$$\sum_{j=1}^{\infty} X_j \le W_0 - W_{\infty} < \infty$$

$$\sum_{j=1}^{\infty} \mathbb{E}[X_j] \le \mathbb{E}[W_0] - w_{\infty} < \infty$$
(39)

which concludes the proof of the lemma.

71 D Details about the Numerical Experiments

D.1 Binary Logistic Regression on the Traumabase

3 D.1.1 Traumabase quantitative variables

474 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 com-475 ponent, the minimum systolic blood pressure, the minimum diastolic blood pressure, the maximum 476 number of heart rate (or pulse) per unit time (usually a minute), the systolic blood pressure at ar-477 rival of ambulance, the diastolic blood pressure at arrival of ambulance, the heart rate at arrival 478 of ambulance, the capillary Hemoglobin concentration, the oxygen saturation, the fluid expansion 479 colloids, the fluid expansion cristalloids, the pulse pressure for the minimum value of diastolic and 480 systolic blood pressure, the pulse pressure at arrival of ambulance. 481

D.1.2 Metropolis-Hastings algorithm

482

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

$$\xi = \beta_{miss} + \Omega_{mis,obs} \Omega_{obs,obs}^{-1} (z_{i,obs} - \beta_{obs}) ,$$

$$\Sigma = \Omega_{mis,mis} + \Omega_{mis,obs} \Omega_{obs,obs}^{-1} \Omega_{obs,mis}$$

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 3.

Algorithm 3 MH aglorithm

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

D.1.3 MISSO Update

491 **Choice of surrogate function for MISO:** We recall the MISO deterministic surrogate defined in 492 (7):

$$\widehat{\mathcal{L}}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) = \int_{\mathsf{Z}} \log \left(p_i(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) / f_i(z_{i,\mathrm{mis}}, \boldsymbol{\theta}) \right) p_i(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) \mu_i(\mathrm{d}z_i) .$$

where $\theta = (\delta, \beta, \Omega)$ and $\overline{\theta} = (\overline{\delta}, \overline{\beta}, \overline{\Omega})$. We adapt it to our missing covariates problem and decompose the surrogate function defined above into an observed and a missing part.

Surrogate function decomposition We adapt it to our missing covariates problem and decompose the term depending on θ , while $\bar{\theta}$ is fixed, in two following parts leading to

$$\begin{split} \widehat{\mathcal{L}}_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}) &= -\int_{\mathsf{Z}} \log f_{i}(z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}, \boldsymbol{\theta}) \, p_{i}(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}}) \\ &= -\int_{\mathsf{Z}} \log \left[p_{i}(y_{i} | z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}, \delta) p_{i}(z_{i,\mathrm{mis}}, \beta, \Omega) \right] \, p_{i}(z_{i}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}}) \\ &= \underbrace{-\int_{\mathsf{Z}} \log p_{i}(y_{i} | z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}, \delta) p_{i}(z_{i}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}})}_{=\widehat{\mathcal{L}}_{i}^{(1)}(\delta, \overline{\boldsymbol{\theta}})} \underbrace{-\int_{\mathsf{Z}} \log p_{i}(z_{i,\mathrm{mis}}, \beta, \Omega) p_{i}(z_{i}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}})}_{=\widehat{\mathcal{L}}_{i}^{(2)}(\beta, \Omega, \overline{\boldsymbol{\theta}})} \end{split}$$

The mean β and the covariance Ω of the latent structure can be estimated minimizing the sum of

MISSO surrogates $\tilde{\mathcal{L}}_i^{(2)}(\beta, \Omega, \overline{\boldsymbol{\theta}}, \{z_m\}_{m=1}^M)$, defined as MC approximation of $\hat{\mathcal{L}}_i^{(2)}(\beta, \Omega, \overline{\boldsymbol{\theta}})$, for all $i \in [n]$, in closed-form expression.

We thus keep the surrogate $\hat{\mathcal{L}}_i^{(2)}(\beta,\Omega,\overline{\pmb{\theta}})$ as it is, and consider the following quadratic approximation

of $\hat{\mathcal{L}}_{i}^{(1)}(\delta, \overline{\boldsymbol{\theta}})$ to estimate the vector of logistic parameters δ :

$$\begin{aligned} \hat{\mathcal{L}}_{i}^{(1)}(\bar{\delta}, \overline{\boldsymbol{\theta}}) &- \int_{\mathsf{Z}} \nabla \log p_{i}(y_{i}|z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}, \delta) \big|_{\delta = \bar{\delta}} p_{i}(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}}) (\delta - \bar{\delta}) \\ &- (\delta - \bar{\delta}) / 2 \int_{\mathsf{Z}} \nabla^{2} \log p_{i}(y_{i}|z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}, \delta) p_{i}(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) p_{i}(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}}) (\delta - \bar{\delta})^{\top} \end{aligned}$$

Recall that:

$$\nabla \log p_i(y_i|z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) = z_i \left(y_i - S(\delta^\top z_i) \right)$$
$$\nabla^2 \log p_i(y_i|z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) = -z_i z_i^\top \dot{S}(\delta^\top z_i)$$

where $\dot{S}(u)$ is the derivative of S(u). Note that $\dot{S}(u) \leq 1/4$ and since, for all $i \in [n]$, the $p \times p$ matrix $z_i z_i^{\top}$ is semi-definite positive we can assume:

L1. For all $i \in [n]$ and $\epsilon > 0$, there exist, for all $z_i \in \mathbb{Z}$, a positive definite matrix $H_i(z_i) := \frac{1}{4}(z_i z_i^\top + \epsilon I_d)$ such that for all $\delta \in \mathbb{R}^p$, $-z_i z_i^\top \dot{S}(\delta^\top z_i) \leq H_i(z_i)$.

Then, we use, for all $i \in [n]$, the following surrogate function to estimate δ : 507

$$\bar{\mathcal{L}}_{i}^{(1)}(\delta, \overline{\boldsymbol{\theta}}) = \hat{\mathcal{L}}_{i}^{(1)}(\bar{\delta}, \overline{\boldsymbol{\theta}}) - D_{i}^{\mathsf{T}}(\delta - \bar{\delta}) + \frac{1}{2}(\delta - \bar{\delta})H_{i}(\delta - \bar{\delta})^{\mathsf{T}}$$
(41)

508 where:

$$D_{i} = \int_{\mathsf{Z}} \nabla \log p_{i}(y_{i}|z_{i,\mathrm{mis}}, z_{i,\mathrm{obs}}, \delta) \big|_{\delta = \bar{\delta}} p_{i}(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}})$$

$$H_{i} = \int_{\mathsf{Z}} H_{i}(z_{i,\mathrm{mis}}) p_{i}(z_{i,\mathrm{mis}}, \overline{\boldsymbol{\theta}}) \mu_{i}(\mathrm{d}z_{i,\mathrm{mis}})$$

Finally, at iteration k, the total surrogate is:

$$\tilde{\mathcal{L}}^{(k)}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \tilde{\mathcal{L}}_{i}(\theta, \theta^{(\tau_{i}^{k})}, \{z_{i,m}\}_{m=1}^{M_{(\tau_{i}^{k})}})$$

$$= \frac{1}{n} \sum_{i=1}^{n} \tilde{\mathcal{L}}_{i}^{(2)}(\beta, \Omega, \theta^{(\tau_{i}^{k})}, \{z_{i,m}\}_{m=1}^{M_{(\tau_{i}^{k})}}) - \frac{1}{n} \sum_{i=1}^{n} \tilde{D}_{i}^{(\tau_{i}^{k})}(\delta - \delta^{(\tau_{i}^{k})})$$

$$+ \frac{1}{2n} \sum_{i=1}^{n} (\delta - \delta^{(\tau_{i}^{k})}) \left\{ \tilde{H}_{i}^{(\tau_{i}^{k})} \right\} (\delta - \delta^{(\tau_{i}^{k})})^{\top}$$
(42)

where for all $i \in [n]$:

$$\tilde{D}_{i}^{(\tau_{i}^{k})} = \frac{1}{M_{(\tau_{i}^{k})}} \sum_{m=1}^{M_{(\tau_{i}^{k})}} z_{i,m}^{(\tau_{i}^{k})} \left(y_{i} - S(\left(\delta^{(\tau_{i}^{k})}\right)^{\top} z_{i,m}(\tau_{i}^{k})) \right)$$

$$\tilde{H}_{i}^{(\tau_{i}^{k})} = \frac{1}{4M_{(\tau_{i}^{k})}} \sum_{m=1}^{M_{(\tau_{i}^{k})}} z_{i,m}^{(\tau_{i}^{k})} (z_{i,m}^{(\tau_{i}^{k})})^{\top}$$

Minimizing the total surrogate (42) boils down to performing a quasi-Newton step. It is perhaps sensible to apply some diagonal loading which is perfectly compatible with the surrogate interpretation we just gave.

The logistic parameters are estimated as follows:

$$\boldsymbol{\delta}^{(k)} = \arg\min_{\delta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \tilde{\mathcal{L}}_{i}^{(1)}(\delta, \theta^{(\tau_{i}^{k})}, \{z_{i,m}\}_{m=1}^{M_{(\tau_{i}^{k})}})$$

where $\tilde{\mathcal{L}}_i^{(1)}(\delta, \theta^{(\tau_i^k)}, \{z_{i,m}\}_{m=1}^{M_{(\tau_i^k)}})$ is the MC approximation of the MISO surrogate defined in (41)and which leads to the following quasi-Newton step:

$$\boldsymbol{\delta}^{(k)} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\delta}^{(\tau_i^k)} - (\tilde{H}^{(k)})^{-1} \tilde{D}^{(k)}$$

sith
$$\tilde{D}^{(k)} = \frac{1}{n} \sum_{i=1}^{n} \tilde{D}_{i}^{(\tau_{i}^{k})}$$
 and $\tilde{H}^{(k)} = \frac{1}{n} \sum_{i=1}^{n} \tilde{H}_{i}^{(\tau_{i}^{k})}$.

MISSO updates: At the k-th iteration, and after the initialization, for all $i \in [n]$, of the latent variables $(z_i^{(0)})$, the MISSO algorithm consists in picking an index i_k uniformly on [n], completing the observations by sampling a Monte Carlo batch $\{z_{i_k, \min, m}^{(k)}\}_{m=1}^{M_{(k)}}$ of missing values from the conditional distribution $p(z_{i_k, \min}|z_{i_k, \text{obs}}, y_{i_k}; \boldsymbol{\theta}^{(k-1)})$ using an MCMC sampler and computing the estimated parameters as follows:

$$\boldsymbol{\beta}^{(k)} = \arg\min_{\beta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \tilde{\mathcal{L}}_{i}^{(2)}(\beta, \Omega^{(k)}, \theta^{(\tau_{i}^{k})}, \{z_{i,m}\}_{m=1}^{M_{(\tau_{i}^{k})}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{M_{(\tau_{i}^{k})}} \sum_{m=1}^{M_{(\tau_{i}^{k})}} z_{i,m}^{(k)}$$

$$\boldsymbol{\Omega}^{(k)} = \arg\min_{\Omega \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \tilde{\mathcal{L}}_{i}^{(2)}(\beta^{(k)}, \Omega, \theta^{(\tau_{i}^{k})}, \{z_{i,m}\}_{m=1}^{M_{(\tau_{i}^{k})}}) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{M_{(\tau_{i}^{k})}} \sum_{m=1}^{M_{(\tau_{i}^{k})}} w_{i,m}^{(k)}$$

$$\boldsymbol{\delta}^{(k)} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\delta}^{(\tau_{i}^{k})} - (\tilde{H}^{(k)})^{-1} \tilde{D}^{(k)}.$$

$$(43)$$

where $z_{i,m}^{(k)}=(z_{i,\text{mis},m}^{(k)},z_{i,\text{obs}})$ is composed of a simulated and an observed part, $\tilde{D}^{(k)}=\frac{1}{n}\sum_{i=1}^{n}\tilde{D}_{i}^{(\tau_{i}^{k})}$, $\tilde{H}^{(k)}=\frac{1}{n}\sum_{i=1}^{n}\tilde{H}_{i}^{(\tau_{i}^{k})}$ and $w_{i,m}^{(k)}=z_{i,m}^{(k)}(z_{i,m}^{(k)})^{\top}-\beta^{(k)}(\beta^{(k)})^{\top}$. Besides, $\tilde{\mathcal{L}}_{i}^{(1)}(\beta,\Omega,\overline{\theta},\{z_{m}\}_{m=1}^{M})$ and $\tilde{\mathcal{L}}_{i}^{(2)}(\beta,\Omega,\overline{\theta},\{z_{m}\}_{m=1}^{M})$ are defined as MC approximation of $\hat{\mathcal{L}}_{i}^{(1)}(\beta,\Omega,\overline{\theta})$ and $\hat{\mathcal{L}}_{i}^{(2)}(\beta,\Omega,\overline{\theta})$, for all $i\in[n]$ as components of the surrogate function (40).

27 D.2 Incremental Variational Inference

528 D.2.1 Bayesian LeNet-5 Architecture

We describe in Table 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

B1 D.2.2 Bayesian ResNet-18 Architecture

We describe in Table 2 the architecture of the Resnet-18 we train on CIFAR-10:

layer type	Output Size	ResNet-18	nonlinearity
conv1	$112 \times 112 \times 64$	7×7 , 64, stride 2	ReLU
conv2x	$56\times 56\times 64$	$\begin{pmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{pmatrix} \times 2$	ReLU
conv3x	$28 \times 28 \times 128$	$\begin{pmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{pmatrix} \times 2$	ReLU
conv4x	$14\times14\times256$	$\begin{pmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{pmatrix} \times 2$	ReLU
conv5x	$7\times7\times512$	$\begin{pmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{pmatrix} \times 2$	ReLU
average pool	$1 \times 1 \times 512$	7×7 average pool	ReLU
fully connected	1000	512×1000 fully connections	
softmax	1000		

Table 2: ResNet-18 architecture

D.2.3 Algorithms updates

533

First, we initialize the means $\mu_\ell^{(0)}$ for $\ell \in \llbracket d \rrbracket$ and variance estimates $\sigma^{(0)}$. At iteration k, minimizing the sum of stochastic surrogates defined as in (6) and (13) 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

$$\mu_{\ell}^{(k)} = \frac{1}{n} \sum_{i=1}^{n} \mu_{\ell}^{(\tau_{i}^{k})} - \frac{\gamma}{n} \sum_{i=1}^{n} \hat{\delta}_{\mu_{\ell},i}^{(k)} \quad \text{and} \quad \sigma^{(k)} = \frac{1}{n} \sum_{i=1}^{n} \sigma^{(\tau_{i}^{k})} - \frac{\gamma}{n} \sum_{i=1}^{n} \hat{\delta}_{\sigma,i}^{(k)}, \tag{44}$$

where we define the following gradient terms for all $i \in [1, n]$:

$$\hat{\delta}_{\mu_{\ell},i}^{(k)} = -\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} \nabla_{w} \log p(y_{i}|x_{i}, w) \Big|_{w=t(\boldsymbol{\theta}^{(k-1)}, z_{m}^{(k)})} + \nabla_{\mu_{\ell}} d(\boldsymbol{\theta}^{(k-1)}),$$

$$\hat{\delta}_{\sigma,i}^{(k)} = -\frac{1}{M_{(k)}} \sum_{m=1}^{M_{(k)}} z_{m}^{(k)} \nabla_{w} \log p(y_{i}|x_{i}, w) \Big|_{w=t(\boldsymbol{\theta}^{(k-1)}, z_{m}^{(k)})} + \nabla_{\sigma} d(\boldsymbol{\theta}^{(k-1)}).$$
(45)

For all benchmark algorithms, we pick, at iteration k, a function index i_k uniformly on [n] 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 $\ell \in [d]$ and σ break down as follows:

42 Monte Carlo SAG update: Set

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

where $\hat{\delta}^{(k)}_{\mu_\ell,i}=\hat{\delta}^{(k-1)}_{\mu_\ell,i}$ and $\hat{\delta}^{(k)}_{\sigma,i}=\hat{\delta}^{(k-1)}_{\sigma,i}$ for $i\neq i_k$ and are defined by (45) for $i=i_k$. The learning rate is set to $\gamma=10^{-3}$.

Bayes By Backprop update: Set

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

where the learning rate $\gamma = 10^{-3}$. 546

Monte Carlo Momentum update: Set 547

$$\mu_{\ell}^{(k)} = \mu_{\ell}^{(k-1)} + \hat{\boldsymbol{v}}_{\mu_{\ell}}^{(k)}$$
 and $\sigma^{(k)} = \sigma^{(k-1)} + \hat{\boldsymbol{v}}_{\sigma}^{(k)}$

where 548

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

where α and γ , respectively the momentum and the learning rates, are set to 10^{-3} . 549

Monte Carlo ADAM update: Set 550

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

where 551

$$\begin{split} \hat{\boldsymbol{m}}_{\mu_{\ell}}^{(k)} &= \boldsymbol{m}_{\mu_{\ell}}^{(k-1)}/(1-\rho_{1}^{k}) \quad \text{with} \quad \boldsymbol{m}_{\mu_{\ell}}^{(k)} &= \rho_{1} \boldsymbol{m}_{\mu_{\ell}}^{(k-1)} + (1-\rho_{1}) \hat{\boldsymbol{\delta}}_{\mu_{\ell}, i_{k}}^{(k)} \; , \\ \hat{\boldsymbol{v}}_{\mu_{\ell}}^{(k)} &= \boldsymbol{v}_{\mu_{\ell}}^{(k-1)}/(1-\rho_{2}^{k}) \quad \text{with} \quad \boldsymbol{v}_{\mu_{\ell}}^{(k)} &= \rho_{2} \boldsymbol{v}_{\mu_{\ell}}^{(k-1)} + (1-\rho_{1}) \left(\hat{\boldsymbol{\delta}}_{\sigma, i_{k}}^{(k)}\right)^{2} \end{split}$$

and 552

$$\begin{split} \hat{\boldsymbol{m}}_{\sigma}^{(k)} &= \boldsymbol{m}_{\sigma}^{(k-1)}/(1-\rho_{1}^{k}) \quad \text{with} \quad \boldsymbol{m}_{\sigma}^{(k)} = \rho_{1}\boldsymbol{m}_{\sigma}^{(k-1)} + (1-\rho_{1})\hat{\boldsymbol{\delta}}_{\sigma,i_{k}}^{(k)} \\ \hat{\boldsymbol{v}}_{\sigma}^{(k)} &= \boldsymbol{v}_{\sigma}^{(k-1)}/(1-\rho_{2}^{k}) \quad \text{with} \quad \boldsymbol{v}_{\sigma}^{(k)} = \rho_{2}\boldsymbol{v}_{\sigma}^{(k-1)} + (1-\rho_{1})\big(\hat{\boldsymbol{\delta}}_{\sigma,i_{k}}^{(k)}\big)^{2} \;. \end{split}$$

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