
A doubly stochastic surrogate optimization scheme for nonconvex finite-sum problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

Many constrained, nonconvex and nonsmooth 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, 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 benchmark datasets.

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) 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 [17] and some variational Bayes inference techniques [8]; see for example [24] and [12] 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 [16], 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 [26]. An extended analysis of MISO has been proposed in [23].

The success of the MISO method rests upon the efficient minimization of surrogates such as convex functions, see [16, Section 2.3]. A notable application of MISO-like algorithms is described in [18] where the authors build upon the stochastic majorization-minimization framework of [16] to introduce a method for sparse matrix factorization. Yet, in many applications of interest, the

34 natural surrogate functions are intractable, yet they are defined as expectation of tractable functions.
 35 For instance, this is the case for inference in latent variable models via maximum likelihood [17].
 36 Another application is variational inference [5], in which the goal is to approximate the posterior
 37 distribution of parameters given the observations; see for example [20; 3; 22; 25; 15].

38 This paper fills the gap in the literature by proposing a method called *Minimization by Incremental*
 39 *Stochastic Surrogate Optimization (MISSO)*, designed for the nonconvex and nonsmooth finite sum
 40 optimization, with a finite-time convergence guarantee. Our work aims at formulating a *generic*
 41 *class* of incremental stochastic surrogate methods for nonconvex optimization and building the the-
 42 ory to understand its behavior. In particular, we provide convergence guarantees for stochastic EM
 43 and Variational Inference-type methods, under mild conditions. In summary, our contributions are:

- 44 • we propose a *unifying framework* of analysis for incremental stochastic surrogate optimiza-
 45 tion when the surrogates are defined as expectations of tractable functions. The proposed
 46 MISSO method is built on the Monte Carlo integration of the intractable surrogate function,
 47 i.e., a doubly stochastic surrogate optimization scheme.
- 48 • we present an incremental update of the commonly used variational inference and Monte
 49 Carlo EM methods as special cases of our newly introduced framework. The analysis of
 50 those two algorithms is thus conducted under this unifying framework of analysis.
- 51 • we establish both asymptotic and non-asymptotic convergence for the MISSO method. In
 52 particular, the MISSO method converges almost surely to a stationary point and in $\mathcal{O}(n/\epsilon)$
 53 iterations to an ϵ -stationary point, see Theorem 1.
- 54 • we relax the class of surrogate functions used in MISO [16] and allow for intractable surro-
 55 gates that can only be evaluated by Monte-Carlo approximations. We show the advantages
 56 of handling such surrogate functions on several *Latent Data* models. Working at the cross-
 57 roads of *Optimization* and *Sampling* constitutes what we believe to be the novelty and the
 58 technicality of our framework and theoretical results.

59 In Section 2, we review the techniques for incremental minimization of finite sum functions based
 60 on the MM principle; specifically, we review the MISO method [16], and present a class of surrogate
 61 functions expressed as an expectation over a latent space. The MISSO method is then introduced
 62 for the latter class of intractable surrogate functions requiring approximation. In Section 3, we pro-
 63 vide the asymptotic and non-asymptotic convergence analysis for the MISSO method (and of the
 64 MISO [16] one as a special case). Section 4 presents numerical applications including parameter in-
 65 ference for logistic regression with missing data and variational inference for two types of Bayesian
 66 neural networks. The proofs of theoretical results are reported as Supplement.

67 **Notations.** We denote $\llbracket 1, n \rrbracket = \{1, \dots, n\}$. Unless otherwise specified, $\|\cdot\|$ denotes the standard
 68 Euclidean norm and $\langle \cdot | \cdot \rangle$ is the inner product in the Euclidean space. For any function $f : \Theta \rightarrow \mathbb{R}$,
 69 $f'(\theta, d)$ is the directional derivative of f at θ along the direction d , i.e.,

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

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

71 2 Incremental Minimization of Finite Sum Nonconvex Functions

72 The objective function in (1) is composed of a finite sum of possibly nonsmooth and nonconvex
 73 functions. A popular approach here is to apply the MM method, which tackles (1) through alter-
 74 nating between two steps — (i) minimizing a *surrogate* function which upper bounds the original
 75 objective function; and (ii) updating the surrogate function to tighten the upper bound.

76 As mentioned in the introduction, the MISO method [16] is developed as an iterative scheme that
 77 only updates the surrogate functions *partially* at each iteration. Formally, for any $i \in \llbracket 1, n \rrbracket$, we
 78 consider a surrogate function $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ which satisfies the assumptions **(H1, H2)**:

79 **H1.** For all $i \in \llbracket 1, n \rrbracket$ and $\bar{\theta} \in \Theta$, $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ is convex w.r.t. θ , and it holds

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

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

81 **H2.** For any $\bar{\theta}_i \in \Theta$, $i \in \llbracket 1, n \rrbracket$ and some $\epsilon > 0$, the difference function $\widehat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n) :=$
82 $\frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}_i(\theta; \bar{\theta}_i) - \mathcal{L}(\theta)$ is defined for all $\theta \in \Theta_\epsilon$ and differentiable for all $\theta \in \Theta$, where
83 $\Theta_\epsilon = \{\theta \in \mathbb{R}^d, \inf_{\theta' \in \Theta} \|\theta - \theta'\| < \epsilon\}$ is an ϵ -neighborhood set of Θ . Moreover, for some constant
84 L , the gradient satisfies

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

Algorithm 1 The MISO method [16].

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

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

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

85 We remark that H1 is a common assumption used for surrogate functions, see [16, Section 2.3]. H2
86 can be satisfied when the difference function $\widehat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n)$ is L -smooth, i.e., \widehat{e} is differentiable on
87 Θ and its gradient $\nabla \widehat{e}$ is L -Lipschitz, $\forall \theta \in \Theta$. H2 can be implied by applying [24, Proposition 1].

88 The inequality (3) implies $\widehat{\mathcal{L}}_i(\theta; \bar{\theta}) \geq \mathcal{L}_i(\theta) > -\infty$ for any $\theta \in \Theta$. The MISO method is an
89 incremental version of the MM method, as summarized by Algorithm 1, which shows that the MISO
90 method maintains an iteratively updated set of upper-bounding surrogate functions $\{\mathcal{A}_i^k(\theta)\}_{i=1}^n$ and
91 updates the iterate via minimizing the average of the surrogate functions.

92 Particularly, only one out of the n surrogate functions is updated at each iteration [cf. Line 5] and
93 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
94 a sum of quadratic functions. As such, the MISO method is suitable for large-scale optimization as
95 the computation cost per iteration is independent of n . Under H1, H2, it was shown that the MISO
96 method converges almost surely to a stationary point of (1) [16, Prop. 3.1].

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

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

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

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

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

109 and we summarize the proposed MISSO method in Algorithm 2. Compared to the MISO method,
110 there is a crucial difference in that the MISSO method involves two types of randomness. The first

Algorithm 2 The MISSO method.

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

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

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

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

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

level of randomness comes from the selection of i_k in Line 5. The second level of randomness stems from the set of Monte Carlo approximated functions $\tilde{\mathcal{A}}_i^k(\theta)$ used in lieu of $\mathcal{A}_i^k(\theta)$ in Line 6 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 [1] 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.*, the joint likelihood of the observations and latent variables. Let

$$g_i(\theta) := \int_{\mathcal{Z}} f_i(z_i, \theta) \mu_i(dz_i), \quad i \in \llbracket 1, n \rrbracket, \quad \theta \in \Theta$$

denote the incomplete data likelihood, *i.e.*, the marginal likelihood of the observations y_i . For ease of notations, the dependence on the observations is made implicit. The maximum likelihood (ML) estimation problem sets the individual objective function $\mathcal{L}_i(\theta)$ 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 observations y_i . A surrogate function $\hat{\mathcal{L}}_i(\theta; \bar{\theta})$ satisfying H1 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:

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) = \int_{\mathcal{Z}} \underbrace{\log(p_i(z_i, \bar{\theta})/f_i(z_i, \theta))}_{=r_i(\theta; \bar{\theta}, z_i)} p_i(z_i, \bar{\theta}) \mu_i(dz_i). \quad (7)$$

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

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 data $x = (x_i, i \in \llbracket 1, n \rrbracket)$, the joint distribution of $y = (y_i, i \in \llbracket 1, n \rrbracket)$ and w is given by:

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

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 [2] boils down to minimizing the Kullback-Leibler (KL) divergence between $q(w, \theta)$ and the posterior distribution $p(w|y, x)$:

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

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

$$\mathcal{L}_i(\boldsymbol{\theta}) := -\mathbb{E}_{q(w;\boldsymbol{\theta})} [\log p(y_i|x_i, w)] + \frac{1}{n} \mathbb{E}_{q(w;\boldsymbol{\theta})} [\log q(w; \boldsymbol{\theta})/\pi(w)] := r_i(\boldsymbol{\theta}) + d(\boldsymbol{\theta}). \quad (10)$$

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:

$$\hat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}) := \mathcal{L}_i(\bar{\boldsymbol{\theta}}) + \langle \nabla_{\boldsymbol{\theta}} \mathcal{L}_i(\bar{\boldsymbol{\theta}}) | \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \rangle + \frac{L}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|^2, \quad (\boldsymbol{\theta}, \bar{\boldsymbol{\theta}}) \in \Theta^2. \quad (11)$$

It is easily checked that the quadratic function $\hat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}})$ satisfies H1, H2. To compute the gradient $\nabla \mathcal{L}_i(\bar{\boldsymbol{\theta}})$, we apply the re-parametrization technique suggested in [21; 10; 3]. 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 $w = t(z, \bar{\boldsymbol{\theta}})$ is $q(\cdot, \bar{\boldsymbol{\theta}})$, where $z \sim \mathcal{N}_d(0, \mathbf{I})$. By [3, Proposition 1], the gradient of $-r_i(\cdot)$ in (10) is:

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

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 as the gradient of the KL between the prior distribution $\pi(\cdot)$ and the variational candidate $q(\cdot, \bar{\boldsymbol{\theta}})$.

$$r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z) := \langle \nabla_{\boldsymbol{\theta}} d(\bar{\boldsymbol{\theta}}) - \mathbf{J}_{\boldsymbol{\theta}}^t(z, \bar{\boldsymbol{\theta}}) \nabla_w \log p(y_i|x_i, w)|_{w=t(z, \bar{\boldsymbol{\theta}})} | \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \rangle + \frac{L}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|^2. \quad (13)$$

Finally, using (11) and (13), 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 $\{z_m\}_{m=1}^M$ are i.i.d samples drawn from $\mathcal{N}(0, \mathbf{I})$.

3 Convergence Analysis

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

H3. For all $i \in \llbracket 1, n \rrbracket$, $\bar{\boldsymbol{\theta}} \in \Theta$, $z_i \in \mathcal{Z}$, $r_i(\cdot; \bar{\boldsymbol{\theta}}, z_i)$ is convex on Θ and is lower bounded.

We are particularly interested in the *constrained optimization* setting where Θ is a bounded set. We thus control the supremum norm of the MC approximation, in (6), as:

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

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

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

where 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 $\Pi_{i, \bar{\boldsymbol{\theta}}}$, and stationary distribution $p_i(\cdot; \bar{\boldsymbol{\theta}})$.

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 $r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z_{i,m}) - \hat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}})$ for all $z \in \mathcal{Z}$ and for fixed $(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}}) \in \Theta^2$. To characterize such control, we will have recourse to the notion of metric entropy (or bracketing number) as developed in [28; 29; 30]. A collection of results from those references gives intuition behind our assumption

H4, which is classical in empirical processes. In [29, Theorem 8.2.3], the authors recall the uniform law of large numbers:

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

for all $z_{i,m}, i \in \llbracket 1, M \rrbracket$ and where \mathcal{F} is a class of L -Lipschitz functions. Moreover, in [29, Theorem 8.1.3] and [30, Theorem 5.22], the application of the Dudley inequality yields:

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

159 where $\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)$ is the bracketing number and ε denotes the level of approximation (the
160 bracketing number goes to infinity when $\varepsilon \rightarrow 0$). Finally, in [28, p.271, Example], $\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)$
161 is bounded from above for a class of parametric functions $\mathcal{F} = f_\theta : \theta \in \Theta$:

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

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

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

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

168 where $g_+(\bar{\theta}) := \max\{0, g(\bar{\theta})\}$, $g_-(\bar{\theta}) := -\min\{0, g(\bar{\theta})\}$ denote the positive and negative part
169 of $g(\bar{\theta})$, respectively. Note that $\bar{\theta}$ is a stationary point if and only if $g_-(\bar{\theta}) = 0$. Furthermore,
170 suppose that the sequence $\{\theta^{(k)}\}_{k \geq 0}$ has a limit point $\bar{\theta}$ that is a stationary point, then one has
171 $\lim_{k \rightarrow \infty} g_-(\theta^{(k)}) = 0$. Thus, the sequence $\{\theta^{(k)}\}_{k \geq 0}$ is said to satisfy an *asymptotic stationary*
172 *point condition*. This is equivalent to [16, Definition 2.4].

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

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

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

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

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

178 Then we have following non-asymptotic bounds:

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

179 Note that $\Delta_{(K_{\max})}$ is finite for any $K_{\max} \in \mathbb{N}$.

180 **Iteration Complexity of MISSO.** As expected, the MISSO method converges to a stationary point
181 of (1) asymptotically and at a sublinear rate $\mathbb{E}[g_-(\theta^{(K)})] \leq \mathcal{O}(\sqrt{\Delta_{(K_{\max})}/K_{\max}})$. In other terms, MISSO

requires $\mathcal{O}(nL/\epsilon)$ iterations to reach an ϵ -stationary point when the suboptimality condition, that characterizes stationarity, is $\mathbb{E}[\|g_-(\theta^{(K)})\|^2]$. Note that this stationarity criterion are similar to the usual quantity used in stochastic nonconvex optimization, *i.e.*, $\mathbb{E}[\|\nabla \mathcal{L}(\theta^{(K)})\|^2]$. In fact, when the optimization problem (1) is unconstrained, *i.e.*, $\Theta = \mathbb{R}^p$, then $\mathbb{E}[g(\theta^{(K)})] = \mathbb{E}[\nabla \mathcal{L}(\theta^{(K)})]$.

Sample Complexity of MISSO. Regarding the sample complexity of our method, setting $M_{(k)} = k^2/n^2$, as a non-decreasing sequence of integers satisfying $\sum_{k=0}^{\infty} M_{(k)}^{-1/2} < \infty$, in order to keep $\Delta_{(K_{\max})} \asymp nL$, then the MISSO method requires $\sum_{k=0}^{nL/\epsilon} k^2/n^2 = nL^3/\epsilon^3$ samples to reach an ϵ -stationary point.

Furthermore, we remark that the MISO method can be analyzed in Theorem 1 as a special case of the MISSO method satisfying $C_r = C_{gr} = 0$. In this case, while the asymptotic convergence is well known from [16] [cf. H4], Eq. (16) gives a non-asymptotic rate of $\mathbb{E}[g_-^{(K)}] \leq \mathcal{O}(\sqrt{nL/K_{\max}})$ 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 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.*, $\lim_{k \rightarrow \infty} g_-(\theta^{(k)}) \stackrel{a.s.}{=} 0$.*
2. *the objective value $\mathcal{L}(\theta^{(k)})$ converges a.s. to a finite number $\underline{\mathcal{L}}$, *i.e.*, $\lim_{k \rightarrow \infty} \mathcal{L}(\theta^{(k)}) \stackrel{a.s.}{=} \underline{\mathcal{L}}$.*

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

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 $z_i = [z_{i,\text{mis}}, z_{i,\text{obs}}]$ is not fully observed where 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 definite $p \times p$ matrix. We define the conditional distribution of the observations y_i given $z_i = (z_{i,\text{mis}}, z_{i,\text{obs}})$ as:

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

where for $u \in \mathbb{R}$, $S(u) = 1/(1+e^{-u})$, $\delta = (\delta_0, \dots, \delta_p)$ are the logistic parameters and $\bar{z}_i = (1, z_i)$. Here, $\theta = (\delta, \beta, \Omega)$ is the parameter to estimate. For $i \in \llbracket n \rrbracket$, the complete log-likelihood reads:

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

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 [7], 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 $p(\cdot|z_{i,\text{obs}}, y_i; \theta^{(k)})$. We compare in Figure 1 the convergence behavior of the estimated parameters δ and β using SAEM [4] (with stepsize $\gamma_k = 1/k^\alpha$ where $\alpha = 0.6$ after tuning), MCEM [31] and the proposed MISSO method. For the MISSO method, we set the batch size to $M_{(k)} = 10 + k^2$ and we

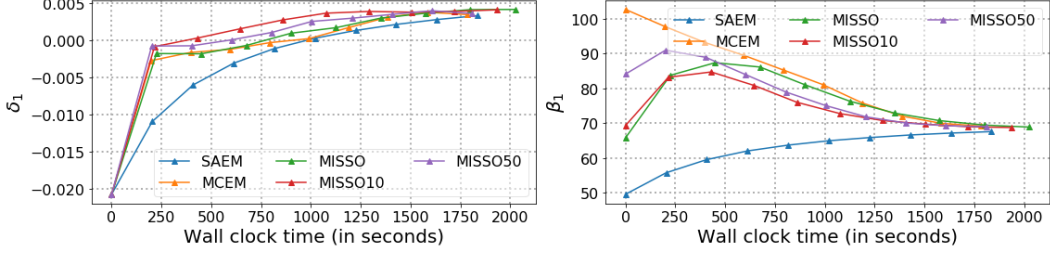


Figure 1: Convergence of parameters δ and β for the SAEM, the MCEM and the MISSO methods. The convergence is plotted against the wall-clock time.

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. Though wall clock times are similar for all methods, they are reported in the appendix for completeness.

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 $\llbracket n \rrbracket$; **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, with $\tilde{w} = t(\theta^{(k-1)}, z_m^{(k)})$, as

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

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

Bayesian LeNet-5 on MNIST [14]: We apply the MISSO method to fit a Bayesian variant of LeNet-5 [14]. We train this network on the MNIST dataset [13]. 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 [6] on CIFAR-10 [11]: We train here the Bayesian variant of the ResNet-18 neural network introduced in [6] 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 efficiency purpose and lower variance, the Flipout estimator [32] is used.

Experiment Results: We compare the convergence of the *Monte Carlo variants* of the following state of the art optimization algorithms — the ADAM [9], the Momentum [27] and the SAG [26] methods versus the *Bayes by Backprop* (BBB) [3] 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 re-parametrization described above. The mini-batch and MC batch size are respectively set to 128 and $M_{(k)} = k$. Learning rates are set to 10^{-3} for LeNet-5 and 10^{-4} for Resnet-18.

Figure 2(a) shows the convergence of the negated evidence lower bound against the wall clocks for fair comparison. As observed, the proposed MISSO method outperforms *Bayes by Backprop*

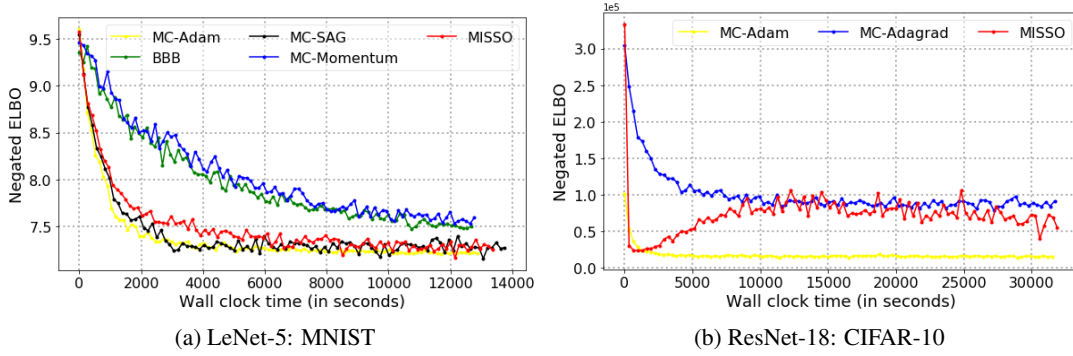


Figure 2: Negated ELBO versus time elapsed for fitting (a) Bayesian LeNet-5 on MNIST and (b) Bayesian ResNet-18 on CIFAR-10. The solid curve is obtained from averaging over 5 independent runs of the methods, and the shaded area represents the standard deviation.

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 lower 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. We report plots against the epochs lapsed and absolute values of running times for all methods in the supplementary for completeness.

Table 1: Runtime in seconds for training a Logistic regression (10 epochs) and BNNs (100 epochs)

	SAEM	MCEM	MISSO	MISSO10	MISSO50
Logistic Regression	2033.2	1972.4	2244.8	2139.4	2005.2
	MC-Adam	MC-Momentum	BBB	MC-SAG	MISSO
LeNet-5 on MNIST	12889	12816	12690	13822	13367

We provide Table 1, the running times for each method on both models, exhibiting the similar runtimes of our method compared to baselines. Those values are reported in the x-axis of each plot for fair comparison. Plots against epochs elapsed are available in the supplementary material.

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

- [1] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [2] David M. Blei, Alp Kucukelbir, and Jon 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>.
- [3] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International Conference on Machine Learning*, pages 1613–1622, 2015.
- [4] Bernard Delyon, Marc Lavielle, and Eric 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>.
- [5] 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.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] Wei Jiang, Julie Josse, and Marc Lavielle. Logistic regression with missing covariates—parameter estimation, model selection and prediction. 2018.
- [8] Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction to variational methods for graphical models. *Mach. Learn.*, 37(2):183–233, November 1999. ISSN 0885-6125. doi: 10.1023/A:1007665907178. URL <https://doi.org/10.1023/A:1007665907178>.
- [9] Diederik P. Kingma and Jimmy 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>.
- [10] Diederik P. Kingma and Max 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>.
- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [12] Kenneth Lange. *MM Optimization Algorithms*. SIAM-Society for Industrial and Applied Mathematics, USA, 2016. ISBN 1611974399, 9781611974393.
- [13] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [14] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [15] Yingzhen Li and Yarin 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.
- [16] Julien Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.

- [17] Geoffrey J. McLachlan and Thriyambakam 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>.
- [18] Arthur Mensch, Julien Mairal, Bertrand Thirion, and Gaël Varoquaux. Stochastic subsampling for factorizing huge matrices. *IEEE Transactions on Signal Processing*, 66(1):113–128, 2017.
- [19] Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.
- [20] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [21] J.W. Paisley, D.M. Blei, and M.I. Jordan. Variational bayesian inference with stochastic search. In *ICML*. icml.cc / Omnipress, 2012.
- [22] Nicholas G Polson, Vadim Sokolov, et al. Deep learning: a bayesian perspective. *Bayesian Analysis*, 12(4):1275–1304, 2017.
- [23] Xun Qian, Alibek Sailanbayev, Konstantin Mishchenko, and Peter Richtárik. Miso is making a comeback with better proofs and rates. *arXiv preprint arXiv:1906.01474*, 2019.
- [24] Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Luo. A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization*, 23(2):1126–1153, 2013.
- [25] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pages 1278–1286, 2014.
- [26] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- [27] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- [28] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [29] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [30] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- [31] Greg C. G. Wei and Martin 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>.
- [32] Yeming Wen, Paul Vicol, Jimmy Ba, Dustin Tran, and Roger Grosse. Flipout: Efficient pseudo-independent weight perturbations on mini-batches. *arXiv preprint arXiv:1803.04386*, 2018.

Checklist

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
- (b) Did you describe the limitations of your work? [Yes]
- (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- (b) Did you include complete proofs of all theoretical results? [Yes]

3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See supplementary material for details on implementation.
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Results are averaged over 5 runs in the main paper and variance is plotted in the supplementary material
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] Running time and total compute are provided. GPUs used for our experiments are Nvidia A100 GPU cards.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? [N/A]
- (b) Did you mention the license of the assets? [N/A]
- (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes]

5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]