

Supplemental File

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

In this supplementary file, we thoroughly provide the proofs for our theorems in the main paper, both the finite-time analysis and the asymptotic convergence result in Section A. In particular, we develop an intermediary lemma, see Lemma 1, extending the Robbins-Siegmund Theorem to non negative sequence of random variables.

Section B and Section C contain additional formulations regarding respectively the logistic regression with missing values and the BNNs examples. In particular, we explicitly provide the MISSO update boiling down from the surrogate functions design. Particularly for the Logistic regression application, where we decompose the deterministic, and intractable, surrogate function in two and employ two strategies of approximation for those two parts, leading to our MISSO surrogate, see subsection B.3 for more details. Additional plots against epochs elapsed are also provided.

A Proofs of the Theoretical Results

A.1 Proof of Theorem 1

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

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

Then we have following non-asymptotic bounds:

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

Proof We begin by recalling the definition

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

Notice that

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

Furthermore, we recall that

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

Due to H2, we have

$$\|\nabla \tilde{\mathcal{L}}^{(k)}(\theta^{(k)})\|^2 \leq 2L\hat{e}^{(k)}(\theta^{(k)}). \quad (18)$$

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

$$\begin{aligned} \tilde{\mathcal{L}}^{(k+1)}(\theta^{(k+1)}) &\leq \tilde{\mathcal{L}}^{(k+1)}(\theta^{(k)}) \\ &= \tilde{\mathcal{L}}^{(k)}(\theta^{(k)}) + \frac{1}{n} (\tilde{\mathcal{L}}_{i_k}(\theta^{(k)}; \theta^{(k)}, \{z_{i_k,m}^{(k)}\}_{m=1}^{M_{(k)}}) - \tilde{\mathcal{L}}_{i_k}(\theta^{(k)}; \theta^{(\tau_{i_k}^k)}, \{z_{i_k,m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})). \end{aligned} \quad (19)$$

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

428 We observe that the conditional expectation evaluates to

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

429 where the last inequality is due to H4. Moreover,

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

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

$$\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \mathcal{L}(\boldsymbol{\theta}^{(k)}) \leq n \mathbb{E}[\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) | \mathcal{F}_k] + \frac{C_r}{\sqrt{M_{(k)}}}. \quad (20)$$

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

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

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

$$\frac{\|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2}{2L} \leq n \mathbb{E}[\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \tilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}^{(k+1)}) | \mathcal{F}_k] + \frac{C_r}{\sqrt{M_{(k)}}} + \delta^{(k)}(\boldsymbol{\theta}^{(k)}). \quad (21)$$

434 Observe the following upper bound on the total expectations:

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

435 which is due to H4. It yields

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

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

$$\begin{aligned} \mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(K)}(\boldsymbol{\theta}^{(K)})\|^2] &= \frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}[\|\nabla \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\|^2] \\ &\leq \frac{2nL \mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\boldsymbol{\theta}^{(K_{\max})})]}{K_{\max}} + \frac{2LC_r}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}\left[\frac{1}{\sqrt{M_{(k)}}} + \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{M_{(\tau_i^k)}}}\right]. \end{aligned} \quad (22)$$

438 For all $i \in \llbracket 1, n \rrbracket$, the index i is selected with a probability equal to $\frac{1}{n}$ when conditioned independ-
439 dently on the past. We observe:

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

440 Taking the sum yields:

$$\begin{aligned}
\sum_{k=0}^{K_{\max}-1} \mathbb{E}[M_{(\tau_i^k)}^{-1/2}] &= \sum_{k=0}^{K_{\max}-1} \sum_{j=1}^k \frac{1}{n} \left(1 - \frac{1}{n}\right)^{j-1} M_{(k-j)}^{-1/2} = \sum_{k=0}^{K_{\max}-1} \sum_{l=0}^{k-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{k-(l+1)} M_{(l)}^{-1/2} \\
&= \sum_{l=0}^{K_{\max}-1} M_{(l)}^{-1/2} \sum_{k=l+1}^{K_{\max}-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{k-(l+1)} \leq \sum_{l=0}^{K_{\max}-1} M_{(l)}^{-1/2},
\end{aligned} \tag{24}$$

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

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

443 This concludes our proof for the first inequality in (16).

444 To prove the second inequality of (16), we define the shorthand notations $g^{(k)} := g(\boldsymbol{\theta}^{(k)})$, $g_-^{(k)} :=$
445 $-\min\{0, g^{(k)}\}$, $g_+^{(k)} := \max\{0, g^{(k)}\}$. We observe that

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

446 where the last inequality is due to the Cauchy-Schwarz inequality and we have defined
447 $\tilde{\mathcal{L}}'_i(\boldsymbol{\theta}, \boldsymbol{d}; \boldsymbol{\theta}^{(\tau_i^k)})$ as the directional derivative of $\tilde{\mathcal{L}}_i(\cdot; \boldsymbol{\theta}^{(\tau_i^k)})$ at $\boldsymbol{\theta}$ along the direction \boldsymbol{d} . Moreover,
448 for any $\boldsymbol{\theta} \in \Theta$,

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

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

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

451 We have

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

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

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

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

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

455 We note that

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

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

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

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

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

458 and concludes the proof of the theorem. □

459 A.2 Proof of Theorem 2

460 **Theorem.** Under H1-H4. In addition, assume that $\{M_{(k)}\}_{k \geq 0}$ is a non-decreasing sequence of
 461 integers which satisfies $\sum_{k=0}^{\infty} M_{(k)}^{-1/2} < \infty$. Then:

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

464 **Proof** We apply the following auxiliary lemma which proof can be found in Appendix A.3 for the
 465 readability of the current proof:

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

$$V_k \leq V_{k-1} - X_{k-1} + E_{k-1} \quad (27)$$

469 then:

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

473 We proceed from (19) by re-arranging terms and observing that

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

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

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

476 is a non-negative random variable.

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

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

478 Thirdly, we define

$$\begin{aligned} E_k &= -(\widetilde{\mathcal{L}}^{(k+1)}(\theta^{(k+1)}) - \widehat{\mathcal{L}}^{(k+1)}(\theta^{(k+1)})) + (\widetilde{\mathcal{L}}^{(k)}(\theta^{(k)}) - \widehat{\mathcal{L}}^{(k)}(\theta^{(k)})) \\ &\quad + \frac{1}{n} (\widetilde{\mathcal{L}}_{i_k}(\theta^{(k)}; \theta^{(k)}, \{z_{i_k, m}^{(k)}\}_{m=1}^{M_{(k)}}) - \widehat{\mathcal{L}}_{i_k}(\theta^{(k)}; \theta^{(k)})) \\ &\quad + \frac{1}{n} (\widehat{\mathcal{L}}_{i_k}(\theta^{(k)}; \theta^{(\tau_{i_k}^k)}) - \widetilde{\mathcal{L}}_{i_k}(\theta^{(k)}; \theta^{(\tau_{i_k}^k)}, \{z_{i_k, m}^{(\tau_{i_k}^k)}\}_{m=1}^{M_{(\tau_{i_k}^k)}})) . \end{aligned} \quad (30)$$

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

480 Under H4, we observe that

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

481

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

482

$$\mathbb{E}[|\tilde{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) - \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})|] \leq \frac{1}{n} \sum_{i=1}^n C_r \mathbb{E}[M_{(\tau_i^k)}^{-1/2}]$$

483 Therefore,

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

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

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

485 Therefore, the conclusions in Lemma 1 hold. Precisely, we have $\sum_{k=0}^{\infty} X_k < \infty$ and
 486 $\sum_{k=0}^{\infty} \mathbb{E}[X_k] < \infty$ almost surely. Note that this implies

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

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

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

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

$$\lim_{k \rightarrow \infty} g_-^{(k)} \leq \lim_{k \rightarrow \infty} \sqrt{2L} \sqrt{\hat{e}^{(k)}(\boldsymbol{\theta}^{(k)})} + \lim_{k \rightarrow \infty} \sup_{\boldsymbol{\theta} \in \Theta} |\epsilon^{(k)}(\boldsymbol{\theta})| = 0, \quad (32)$$

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

492 Finally, we prove that $\mathcal{L}(\boldsymbol{\theta}^{(k)})$ converges almost surely. As a consequence of Lemma 1, it is clear that
 493 $\{V_k\}_{k \geq 0}$ converges almost surely and so is $\{\hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)})\}_{k \geq 0}$, i.e., we have $\lim_{k \rightarrow \infty} \hat{\mathcal{L}}^{(k)}(\boldsymbol{\theta}^{(k)}) = \underline{\mathcal{L}}$.
 494 Applying (31) implies that

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

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

496 A.3 Proof of Lemma 1

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

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

500 then:

501 (i) for all $k \geq 0$, $\mathbb{E}[V_k] < \infty$ and the sequence $(V_k)_{k \geq 0}$ converges a.s. to a finite limit V_{∞} .

502 (ii) the sequence $(\mathbb{E}[V_k])_{k \geq 0}$ converges and $\lim_{k \rightarrow \infty} \mathbb{E}[V_k] = \mathbb{E}[V_{\infty}]$.

503 (iii) the series $\sum_{k=0}^{\infty} X_k$ converges almost surely and $\sum_{k=0}^{\infty} \mathbb{E}[X_k] < \infty$.

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

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

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

506 Since $0 \leq X_k \leq V_{k-1} - V_k + E_k$ we also obtain for all $k \geq 0$, $\mathbb{E}[X_k] < \infty$. Moreover, since
 507 $\mathbb{E}\left[\sum_{j=1}^{\infty} |E_j|\right] < \infty$, the series $\sum_{j=1}^{\infty} E_j$ converges a.s. We may therefore define:

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

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

$$\begin{aligned} W_k &\leq V_{k-1} - X_k + \sum_{j=k}^{\infty} E_j \leq W_{k-1} - X_k \leq W_{k-1} \\ \mathbb{E}[W_k] &\leq \mathbb{E}[W_{k-1}] - \mathbb{E}[X_k]. \end{aligned} \quad (35)$$

509 Hence the sequences $(W_k)_{k \geq 0}$ and $(\mathbb{E}[W_k])_{k \geq 0}$ are non increasing. Since for all $k \geq 0$, $W_k \geq$
 510 $-\sum_{j=1}^{\infty} |E_j| > -\infty$ and $\mathbb{E}[W_k] \geq -\sum_{j=1}^{\infty} \mathbb{E}[|E_j|] > -\infty$, the (random) sequence $(W_k)_{k \geq 0}$
 511 converges a.s. to a limit W_{∞} and the (deterministic) sequence $(\mathbb{E}[W_k])_{k \geq 0}$ converges to a limit w_{∞} .
 512 Since $|W_k| \leq V_0 + \sum_{j=1}^{\infty} |E_j|$, the Fatou lemma implies that:

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

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

514 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
 515 $\mathbb{E}[U_k] \leq \mathbb{E}[|W_0|] + \mathbb{E}[|W_k|] < \infty$ and by the monotone convergence theorem, we get:

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

516 Finally, we have:

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

517 showing that $\mathbb{E}[W_{\infty}] = w_{\infty}$ and concluding the proof of (ii). Moreover, using (35) we have that
 518 $W_k \leq W_{k-1} - X_k$ which yields:

$$\begin{aligned} \sum_{j=1}^{\infty} X_j &\leq W_0 - W_{\infty} < \infty, \\ \sum_{j=1}^{\infty} \mathbb{E}[X_j] &\leq \mathbb{E}[W_0] - w_{\infty} < \infty, \end{aligned} \quad (39)$$

519 an concludes the proof of the lemma. \square

520 B Practical Details for the Binary Logistic Regression on the Traumabase

521 B.1 Traumabase dataset quantitative variables

522 The list of the 16 quantitative variables we use in our experiments are as follows — *age, weight,*
 523 *height, BMI (Body Mass Index), the Glasgow Coma Scale, the Glasgow Coma Scale motor com-*
 524 *ponent, the minimum systolic blood pressure, the minimum diastolic blood pressure, the maximum*
 525 *number of heart rate (or pulse) per unit time (usually a minute), the systolic blood pressure at ar-*
 526 *rival of ambulance, the diastolic blood pressure at arrival of ambulance, the heart rate at arrival*
 527 *of ambulance, the capillary Hemoglobin concentration, the oxygen saturation, the fluid expansion*
 528 *colloids, the fluid expansion cristalloids, the pulse pressure for the minimum value of diastolic and*
 529 *systolic blood pressure, the pulse pressure at arrival of ambulance.*

530 B.2 Metropolis-Hastings algorithm

531 During the simulation step of the MISSO method, the sampling from the target distribution
 532 $\pi(z_{i,\text{mis}}; \theta) := p(z_{i,\text{mis}} | z_{i,\text{obs}}, y_i; \theta)$ is performed using a Metropolis-Hastings (MH) algorithm [19]
 533 with proposal distribution $q(z_{i,\text{mis}}; \delta) := p(z_{i,\text{mis}} | z_{i,\text{obs}}; \delta)$ where $\theta = (\beta, \Omega)$ and $\delta = (\xi, \Sigma)$. The
 534 parameters of the Gaussian conditional distribution of $z_{i,\text{mis}} | z_{i,\text{obs}}$ read:

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

535 where we have used the Schur Complement of $\Omega_{\text{obs},\text{obs}}$ in Ω and noted β_{mis} (resp. β_{obs}) the missing
 536 (resp. observed) elements of β . The MH algorithm is summarized in Algorithm 3.

Algorithm 3 MH algorithm

```

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

```

537 B.3 MISSO Update

538 **Choice of surrogate function for MISO:** We recall the MISO deterministic surrogate defined in
 539 (7):

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) = \int_{\mathcal{Z}} \log(p_i(z_{i,\text{mis}}, \bar{\theta}) / f_i(z_{i,\text{mis}}, \theta)) p_i(z_{i,\text{mis}}, \bar{\theta}) \mu_i(dz_i) .$$

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

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

$$\begin{aligned}
 & \hat{\mathcal{L}}_i(\theta; \bar{\theta}) \\
 &= - \int_{\mathbf{Z}} \log f_i(z_{i,\text{mis}}, z_{i,\text{obs}}, \theta) p_i(z_{i,\text{mis}}, \bar{\theta}) \mu_i(dz_{i,\text{mis}}) \\
 &= - \int_{\mathbf{Z}} \log [p_i(y_i | z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) p_i(z_{i,\text{mis}}, \beta, \Omega)] p_i(z_i, \bar{\theta}) \mu_i(dz_{i,\text{mis}}) \\
 &= \underbrace{- \int_{\mathbf{Z}} \log p_i(y_i | z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) p_i(z_i, \bar{\theta}) \mu_i(dz_{i,\text{mis}})}_{=\hat{\mathcal{L}}_i^{(1)}(\delta, \bar{\theta})} - \underbrace{\int_{\mathbf{Z}} \log p_i(z_{i,\text{mis}}, \beta, \Omega) p_i(z_i, \bar{\theta}) \mu_i(dz_{i,\text{mis}})}_{=\hat{\mathcal{L}}_i^{(2)}(\beta, \Omega, \bar{\theta})} .
 \end{aligned} \tag{40}$$

544 The mean β and the covariance Ω of the latent structure can be estimated minimizing the sum of
 545 MISSO surrogates $\tilde{\mathcal{L}}_i^{(2)}(\beta, \Omega, \bar{\theta}, \{z_m\}_{m=1}^M)$, defined as MC approximation of $\hat{\mathcal{L}}_i^{(2)}(\beta, \Omega, \bar{\theta})$, for all
 546 $i \in \llbracket n \rrbracket$, in closed-form expression.

547 We thus keep the surrogate $\hat{\mathcal{L}}_i^{(2)}(\beta, \Omega, \bar{\theta})$ as it is, and consider the following quadratic approximation
 548 of $\hat{\mathcal{L}}_i^{(1)}(\delta, \bar{\theta})$ to estimate the vector of logistic parameters δ :

$$\begin{aligned}
 & \hat{\mathcal{L}}_i^{(1)}(\bar{\delta}, \bar{\theta}) - \int_{\mathbf{Z}} \nabla \log p_i(y_i | z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) \Big|_{\delta=\bar{\delta}} p_i(z_{i,\text{mis}}, \bar{\theta}) \mu_i(dz_{i,\text{mis}}) (\delta - \bar{\delta}) \\
 & \quad - (\delta - \bar{\delta})/2 \int_{\mathbf{Z}} \nabla^2 \log p_i(y_i | z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) p_i(z_{i,\text{mis}}, \bar{\theta}) p_i(z_{i,\text{mis}}, \bar{\theta}) \mu_i(dz_{i,\text{mis}}) (\delta - \bar{\delta})^\top .
 \end{aligned}$$

549 Recall that:

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

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

552 **L1.** For all $i \in \llbracket n \rrbracket$ and $\epsilon > 0$, there exist, for all $z_i \in \mathbf{Z}$, a positive definite matrix $H_i(z_i) :=$
 553 $\frac{1}{4}(z_i z_i^\top + \epsilon I_d)$ such that for all $\delta \in \mathbb{R}^p$, $-z_i z_i^\top \dot{S}(\delta^\top z_i) \leq H_i(z_i)$.

554 Then, we use, for all $i \in \llbracket n \rrbracket$, the following surrogate function to estimate δ :

$$\bar{\mathcal{L}}_i^{(1)}(\delta, \bar{\theta}) = \hat{\mathcal{L}}_i^{(1)}(\bar{\delta}, \bar{\theta}) - D_i^\top (\delta - \bar{\delta}) + \frac{1}{2} (\delta - \bar{\delta}) H_i (\delta - \bar{\delta})^\top , \tag{41}$$

555 where:

$$\begin{aligned}
 D_i &= \int_{\mathbf{Z}} \nabla \log p_i(y_i | z_{i,\text{mis}}, z_{i,\text{obs}}, \delta) \Big|_{\delta=\bar{\delta}} p_i(z_{i,\text{mis}}, \bar{\theta}) \mu_i(dz_{i,\text{mis}}) , \\
 H_i &= \int_{\mathbf{Z}} H_i(z_{i,\text{mis}}) p_i(z_{i,\text{mis}}, \bar{\theta}) \mu_i(dz_{i,\text{mis}}) .
 \end{aligned}$$

556 Finally, at iteration k , the total surrogate is:

$$\begin{aligned}
 \tilde{\mathcal{L}}^{(k)}(\theta) &= \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{L}}_i(\theta, \theta^{(\tau_i^k)}, \{z_{i,m}\}_{m=1}^{M(\tau_i^k)}) \\
 &= \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{L}}_i^{(2)}(\beta, \Omega, \theta^{(\tau_i^k)}, \{z_{i,m}\}_{m=1}^{M(\tau_i^k)}) - \frac{1}{n} \sum_{i=1}^n \tilde{D}_i^{(\tau_i^k)} (\delta - \delta^{(\tau_i^k)}) \\
 & \quad + \frac{1}{2n} \sum_{i=1}^n (\delta - \delta^{(\tau_i^k)}) \left\{ \tilde{H}_i^{(\tau_i^k)} \right\} (\delta - \delta^{(\tau_i^k)})^\top ,
 \end{aligned} \tag{42}$$

557 where for all $i \in \llbracket n \rrbracket$:

$$\begin{aligned}\tilde{D}_i^{(\tau_i^k)} &= \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} z_{i,m}^{(\tau_i^k)} \left(y_i - S\left(\delta^{(\tau_i^k)}\right)^\top z_{i,m}(\tau_i^k) \right), \\ \tilde{H}_i^{(\tau_i^k)} &= \frac{1}{4M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} z_{i,m}^{(\tau_i^k)} (z_{i,m}^{(\tau_i^k)})^\top.\end{aligned}$$

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

561 The logistic parameters are estimated as follows:

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

562 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)
563 and which leads to the following quasi-Newton step:

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

564 with $\tilde{D}^{(k)} = \frac{1}{n} \sum_{i=1}^n \tilde{D}_i^{(\tau_i^k)}$ and $\tilde{H}^{(k)} = \frac{1}{n} \sum_{i=1}^n \tilde{H}_i^{(\tau_i^k)}$.

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

$$\begin{aligned}\beta^{(k)} &= \arg \min_{\beta \in \Theta} \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{L}}_i^{(2)}(\beta, \Omega^{(k)}, \theta^{(\tau_i^k)}, \{z_{i,m}\}_{m=1}^{M_{(\tau_i^k)}}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} z_{i,m}^{(k)}, \\ \Omega^{(k)} &= \arg \min_{\Omega \in \Theta} \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{L}}_i^{(2)}(\beta^{(k)}, \Omega, \theta^{(\tau_i^k)}, \{z_{i,m}\}_{m=1}^{M_{(\tau_i^k)}}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{M_{(\tau_i^k)}} \sum_{m=1}^{M_{(\tau_i^k)}} w_{i,m}^{(k)}, \\ \delta^{(k)} &= \frac{1}{n} \sum_{i=1}^n \delta^{(\tau_i^k)} - (\tilde{H}^{(k)})^{-1} \tilde{D}^{(k)}.\end{aligned}\tag{43}$$

570 where $z_{i,m}^{(k)} = (z_{i, \text{mis}, m}^{(k)}, z_{i, \text{obs}})$ is composed of a simulated and an observed part, $\tilde{D}^{(k)} =$
571 $\frac{1}{n} \sum_{i=1}^n \tilde{D}_i^{(\tau_i^k)}$, $\tilde{H}^{(k)} = \frac{1}{n} \sum_{i=1}^n \tilde{H}_i^{(\tau_i^k)}$ and $w_{i,m}^{(k)} = z_{i,m}^{(k)} (z_{i,m}^{(k)})^\top - \beta^{(k)} (\beta^{(k)})^\top$. Be-
572 sides, $\tilde{\mathcal{L}}_i^{(1)}(\beta, \Omega, \bar{\theta}, \{z_m\}_{m=1}^M)$ and $\tilde{\mathcal{L}}_i^{(2)}(\beta, \Omega, \bar{\theta}, \{z_m\}_{m=1}^M)$ are defined as MC approximation of
573 $\hat{\mathcal{L}}_i^{(1)}(\beta, \Omega, \bar{\theta})$ and $\hat{\mathcal{L}}_i^{(2)}(\beta, \Omega, \bar{\theta})$, for all $i \in \llbracket n \rrbracket$ as components of the surrogate function (40).

574 B.4 Wall clock time

575 We provide Table 2, the running time for each method, plotted in Figure 1, employed to train a
576 logistic regression with missing values on the TraumaBase dataset ($p = 16$ influential quantitative
577 measurements, on $n = 6384$ patients).

578 The running times are sensibly the same since for each method the computation complexity per
579 epoch is similar. We remark a slight delay using the MISSO method with a batch size of 1, as our
580 code implemented in R, is not totally optimized and parallelized. Yet, when the batch size tends to

100%, we retrieve the duration of MCEM, which is consistent with the fact that MISSO with a full batch update boils down to the MCEM algorithm.

	SAEM	MCEM	MISSO	MISSO10	MISSO50
Logistic Regression	2033.2	1972.4	2244.8	2139.4	2005.2

Table 2: Logistic Regression with missing values: running time in seconds for 10 epochs.

We plot Figure 3, the updated parameters for the Logistic regression example against the time elapsed (in seconds).

B.5 Plots against the epochs elapsed

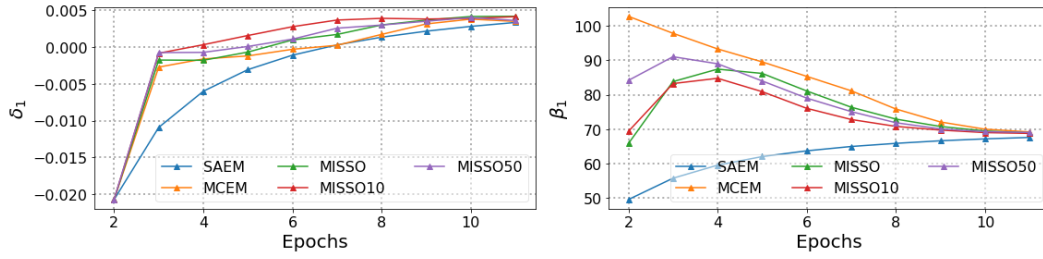


Figure 3: Convergence of parameters δ and β for the SAEM, the MCEM and the MISSO methods. The convergence is plotted against epochs elapsed.

586 C Practical Details for the Incremental Variational Inference

587 C.1 Neural Networks Architecture

588 **Bayesian LeNet-5 Architecture:** We describe in Table 3 the architecture of the Convolutional
 589 Neural Network introduced in [14] 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 3: LeNet-5 architecture

590 **Bayesian ResNet-18 Architecture:** We describe in Table 4 the architecture of the Resnet-18 we
 591 train on CIFAR-10:

layer type	Output Size	ResNet-18	nonlinearity
conv1	$112 \times 112 \times 64$	$7 \times 7, 64, \text{stride } 2$	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	$14 \times 14 \times 256$	$\begin{pmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{pmatrix} \times 2$	ReLU
conv5x	$7 \times 7 \times 512$	$\begin{pmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{pmatrix} \times 2$	ReLU
average pool	$1 \times 1 \times 512$	7×7 average pool	ReLU
fully connected	1000	512×1000 fully connections	
softmax	1000		

Table 4: ResNet-18 architecture

592 C.2 Algorithms updates

593 First, we initialize the means $\mu_\ell^{(0)}$ for $\ell \in \llbracket d \rrbracket$ and variance estimates $\sigma^{(0)}$. At iteration k , minimizing
 594 the sum of stochastic surrogates defined as in (6) and (13) yields the following MISSO update —
 595 **step (i)** pick a function index i_k uniformly on $\llbracket n \rrbracket$; **step (ii)** sample a Monte Carlo batch $\{z_m^{(k)}\}_{m=1}^{M(k)}$
 596 from $\mathcal{N}(0, \mathbf{I})$; and **step (iii)** update the parameters as

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

597 where we define the following gradient terms for all $i \in \llbracket 1, n \rrbracket$:

$$\begin{aligned} \hat{\delta}_{\mu_\ell, i}^{(k)} &= -\frac{1}{M(k)} \sum_{m=1}^{M(k)} \nabla_w \log p(y_i | x_i, w) \Big|_{w=t(\theta^{(k-1)}, z_m^{(k)})} + \nabla_{\mu_\ell} d(\theta^{(k-1)}), \\ \hat{\delta}_{\sigma, i}^{(k)} &= -\frac{1}{M(k)} \sum_{m=1}^{M(k)} z_m^{(k)} \nabla_w \log p(y_i | x_i, w) \Big|_{w=t(\theta^{(k-1)}, z_m^{(k)})} + \nabla_\sigma d(\theta^{(k-1)}). \end{aligned} \quad (45)$$

598 Note that our analysis in the main text does require the parameter to be in a compact set. For the
 599 current estimation problem considered, this can be enforced in practice by restricting the parameters

in a ball. In our simulation for the BNNs example, we did not implement the algorithms that stick closely to the compactness requirement for illustrative purposes. However, we observe empirically that the parameters are always bounded. The update rules can be easily modified to respect the requirement. For the considered VI problem, we recall the surrogate functions (11) are quadratic and indeed a simple projection step suffices to ensure boundedness of the iterates.

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

Monte Carlo SAG update: Set

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

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

Bayes By Backprop update: Set

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

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

Monte Carlo Momentum update: Set

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

where

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

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

Monte Carlo ADAM update: Set

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

where

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

and

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

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

620 C.3 Plots against the epochs elapsed

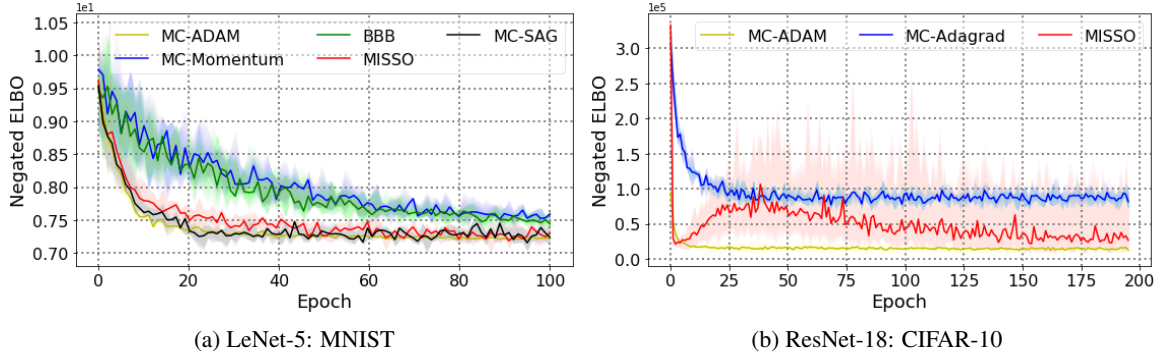


Figure 4: Negated ELBO versus epochs elapsed for fitting (a) Bayesian LeNet-5 on MNIST and (b) Bayesian ResNet-18 on CIFAR-10. The solid curve is obtained from averaging over 5 independent runs of the methods, and the shaded area represents the standard deviation.

621 C.4 Wall clock time

622 We provide Table 5, the running time for each method, plotted in Figure 2, used to train a Bayesian
623 variant of LeNet-5 on MNIST. The incremental method as MISSO and MC-SAG displays a similar
624 wall clock time, despite being a bit worse given (a) the initialization that requires to compute a vector
625 of n gradients kept in memory and updated through the iterations and (b) the average operation for
626 each parameters update to compute the aggregated drift term, see (44).

	MC-Adam	MC-Momentum	BBB	MC-SAG	MISSO
LeNet-5 on MNIST	12889	12816	12690	13822	13367

Table 5: Bayesian Deep Neural Network: running time in seconds for 100 epochs.

627 We plot Figure 5, the learning curves for the MNIST example against the time elapsed (in seconds).

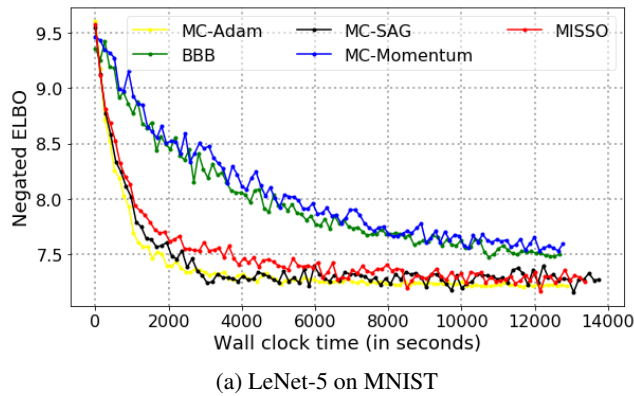


Figure 5: Negated ELBO versus wall clock time for fitting a Bayesian LeNet-5 on MNIST. Plotted on the average of the 5 repetitions.