1 We would like to thank four reviewers for their feedback. Upon acceptance, we will include in the final version (a) *a*
2 *clearer presentation of the numerical results* and (b) *missing references*. We first discuss a common concern shared by
3 **reviewer 1**, **reviewer 2**, **reviewer 4**.

4 ●●● **Novelty of The Contribution**: We want to stress on the generality of our incremental framework, which tackles a
5 *constrained*, *non-convex* and *non-smooth* optimization problem. The main contribution of this paper is to propose a
6 *unifying* framework for the analysis of a large class of optimization algorithms which indeed includes well-known but
7 not so well-studied algorithms. The major goal here is to relax the class of surrogate functions used in MISO [Mairal,
8 2015] and replace that by the respective Monte-Carlo approximations. We provide a general algorithm and global
9 convergence analysis under mild assumptions on the model and show that two examples, MLE for general latent data
10 models and Variational Inference, are its special instances.

11 Working at the crossroads of *Optimization* and *Sampling* constitues, we believe, the novelty and the technicality of our
12 theoretical results.

13 **Reviewer 1:** We thank the reviewer for valuable comments and references. We would like to make the following
14 clarification regarding the difference with MISO:

15 **Originality:** The main contribution of the paper is to extend the MISO algorithm when the surrogate fun ctions are not
16 tractable. We motivate the need for dealing with intractable surrogate functions when nonconvex latent data models
17 are being trained. In this case, the latent structure yields an expected surrogate functions and the nonconvexity yields
18 an intractable expectation to compute. The only option is to build a stochastic surrogate function based on a MC
19 approximation.

20 **Reviewer 2:** We thank the reviewer for the useful comments. Our point-to-point response is as follows:

21 **Numerical Plots:** Due to space constraints, we only presented several dimension for the logistic parameter and the
22 mean of the latent variable. As the reviewer mentioned, we also learn the variance of those latent variables and the
23 convergence plots of those variances will be added to the rebuttal version.

24 **Wallclock Time**:

25 Wallclock time per iteration is comparable for each method. Indeed the methods always only involve first order
26 computation. Yet, we acknowledge that MISSO can present some memory bottlenecks since it requires to store $n$
27 gradients through the run. This has not been a problem for the presented numerical examples

28 **Parameter Tuning:**

29 The baseline methods were tuned and presented to the best of their performances both with regards to their stepsize
30 (grid search) and minibatch size. We believe your remark refers to the first numerical example (logistic regression with
31 missing values): Regarding the stepsize, as MCEM does not have one, we indeed tuned the stepsize of SAEM. Rather
32 than $c/k$, common practice is to tune a parameter $\alpha$ such that $\gamma_k = 1/\gamma^\alpha$. We report results for SAEM with the best $\alpha$
33 ($\alpha = 0.6$). Regarding batch size, for SAEM and MCEM both are full batch methods and the idea here is to compare
34 different values of minibatch size for the MISSO method to see its influence on the performances.

35 **Reviewer 3:** We thank the reviewer for valuable comments and references. Please find the following precisions
36 regarding the numerical examples:

37 **Assumptions and Numerical Examples:**

38 **Reviewer 4:** We thank the reviewer for valuable comments and the numerous related references. Our point-to-point
39 response is as follows:

40 **Comparison to [Murray+, 2012] and [Tran+, 2017]:**

41 **Comparison to [Kang+, 2015]:**

42 **Comparison with MC-ADAM Figure 2:**