



**HAL**  
open science

# Stochastic Variable Metric Proximal Gradient with variance reduction for non-convex composite optimization

Gersende Fort, Eric Moulines

► **To cite this version:**

Gersende Fort, Eric Moulines. Stochastic Variable Metric Proximal Gradient with variance reduction for non-convex composite optimization. 2022. hal-03781216v2

**HAL Id: hal-03781216**

**<https://hal.science/hal-03781216v2>**

Preprint submitted on 31 Dec 2022 (v2), last revised 2 Mar 2023 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Stochastic Variable Metric Proximal Gradient with variance reduction for non-convex composite optimization

Gersende Fort<sup>1\*</sup> and Eric Moulines<sup>2</sup>

<sup>1\*</sup>Institut de Mathématiques de Toulouse, CNRS & Université de Toulouse, 118 route de Narbonne, Toulouse, 31400, France.

<sup>2</sup>CMAP, Ecole Polytechnique, Route de Saclay, Palaiseau, 91128 Cedex, France.

\*Corresponding author(s). E-mail(s): [gersende.fort@math.univ-toulouse.fr](mailto:gersende.fort@math.univ-toulouse.fr);  
Contributing authors: [eric.moulines@polytechnique.edu](mailto:eric.moulines@polytechnique.edu);

## Abstract

This paper introduces a novel algorithm, the Perturbed Proximal Preconditioned SPIDER algorithm (**3P-SPIDER**), designed to solve finite sum non-convex composite optimization. It is a stochastic Variable Metric Forward-Backward algorithm, which allows approximate preconditioned forward operator and uses a variable metric proximity operator as the backward operator; it also proposes a mini-batch strategy with variance reduction to address the finite sum setting. We show that **3P-SPIDER** extends some Stochastic preconditioned Gradient Descent-based algorithms and some Incremental Expectation Maximization algorithms to composite optimization and to the case the forward operator can not be computed in closed form. We also provide an explicit control of convergence in expectation of **3P-SPIDER**, and study its complexity in order to satisfy the epsilon-approximate stationary condition. Our results are the first to combine the non-convex composite optimization setting, a variance reduction technique to tackle the finite sum setting by using a minibatch strategy and, to allow deterministic or random approximations of the preconditioned forward operator. Finally, through an application to inference in a logistic regression model with random effects, we numerically compare **3P-SPIDER** to other stochastic forward-backward algorithms and discuss the role of some design parameters of **3P-SPIDER**.

**Keywords:** Stochastic optimization, Variable Metric Forward-Backward splitting, Preconditioned Stochastic Gradient Descent, Incremental Expectation Maximization, Proximal methods, Variance reduction, Non-asymptotic convergence bounds.

## 1 Introduction

Efficient learning from large data sets require new optimization algorithms designed to be robust to big data and complex models era. In Statistics and Machine Learning, we are often faced with solving

problems of the form

$$\operatorname{argmin}_{s \in \mathbb{R}^q} \left( \frac{1}{n} \sum_{i=1}^n W_i(s) + g(s) \right),$$

where  $n$  is the number of examples in the training data set,  $s$  is an unknown quantity to be learnt from the examples,  $W_i$  is a loss function associated to the example  $\#i$  and  $g$  is a regularization

term encoding a priori knowledge and constraints on  $s$ ;  $g$  may also prevent from overfitting. Quite often, the regularization term  $g : \mathbb{R}^q \rightarrow (0, +\infty]$  is not differentiable, and the data fidelity term  $n^{-1} \sum_{i=1}^n W_i$  is smooth on the domain of  $g$ .

This paper is concerned with stochastic optimization of a non-convex finite sum composite function; more precisely, it addresses the differential inclusion problem

$$0 \in \frac{1}{n} \sum_{i=1}^n G_i(s) + \partial g(s), \quad s \in \mathbb{R}^q, \quad (1)$$

where  $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$  is lower semi-continuous convex with non-empty domain  $\mathcal{S}$  and for all  $i$ ,  $G_i : \mathcal{S} \rightarrow \mathbb{R}^q$  is globally Lipschitz on  $\mathcal{S}$ .

The first goal of this paper is to provide a novel algorithm. Motivated by applications in Statistics and Machine learning, we require this algorithm to satisfy the following three constraints. (c1) The algorithm uses possibly preconditioned operators  $h_i$  instead of the forward operator  $G_i$ :

$$\forall s \in \mathcal{S}, \quad h_i(s, B) \stackrel{\text{def}}{=} -B^{-1} G_i(s), \quad (2)$$

where  $B$  is a  $q \times q$  positive definite matrix. Such a condition encompasses preconditioned gradient methods for example, which also includes gradient methods with adaptive step sizes. It also encompasses Expectation Maximization (EM) algorithms (Dempster et al (1977)) designed for large scale learning. (c2) The algorithm may only have access to approximations of  $h_i(s, B)$ . Such a condition addresses the situations when  $h_i(s, B)$  is not explicit, for example when it is defined by an intractable integral. This occurs at each E-step of EM, when the conditional expectations under the a posteriori distributions can not be computed exactly. (c3) The algorithm addresses the finite sum challenge while keeping the caused variability induced by the algorithmic solution as small as possible. For example, when the solution relies on a random selection of a mini-batch of examples, the algorithm has to propose a variance reduction scheme.

A first class of problems of the form (1) are minimizations of regularized loss functions through gradient-based algorithms. In that case,  $G_i \stackrel{\text{def}}{=} \nabla W_i$ . (2) allows preconditioned gradients; such a variable metric is known to accelerate the

convergence. Variable Metric Forward-Backward (VMFB) algorithms were introduced to solve (1)-(2) in the case  $G$  is a gradient. Nevertheless, as discussed in Section 2.1, to our best knowledge none of the variants of VMFB address the three constraints c1, c2 and c3 simultaneously.

A second application of (1)-(2) is the EM algorithm, an algorithm originally designed to compute the Maximum-Likelihood estimator of an unknown parameter  $\theta$  in latent variable models. When the complete data model is from the curved exponential family, EM is equivalent to an algorithm in the so-called *statistic space* (see e.g. Delyon et al (1999)). This remark is the cornerstone of stochastic EM algorithms including incremental EM ones designed for incremental processing of large data sets. EM only supplies preconditioned forward operators  $-B(s)^{-1} G_i(s)$ . Therefore, stochastic EM algorithms are naturally in the setting (1)-(2) (see Fort et al (2020)). Here again, as discussed in Section 2.2, none of the EM variants in the literature address the constraints (c2) and (c3) simultaneously.

Our first contribution is the design of a novel iterative algorithm, named **Perturbed Proximal Preconditioned SPIDER** (3P-SPIDER), which combines (i) a preconditioned forward operator associated to the smooth part  $n^{-1} \sum_{i=1}^n G_i(s)$ , (ii) a variable metric proximity operator with respect to the non-smooth part  $g$ , (iii) a stochastic approach to address the finite sum setting induced by  $n^{-1} \sum_{i=1}^n G_i(s)$  combined with a variance reduction technique based on the SPIDER algorithm (Nguyen et al (2017); Fang et al (2018)); it also allows (iv) numerical approximations of the preconditioned forward operator when it has no closed form. The algorithm is introduced in Section 3, together with discussions on implementation questions. We also design a stochastic VMFB algorithm which answers the constraints c1 and c2 but does not contain a variance reduction step as required by c3.

The second contribution is to provide a non-asymptotic convergence analysis in expectation of 3P-SPIDER in the case the variable metric  $B$  at iteration  $\#t$  depends on the current value of the iterate, and the  $G_i$ 's are gradient operators; see Section 4. The proof relies on a Lyapunov inequality with an original construction, which

is a consequence of the non-convex optimization setting, and the fact the algorithm uses preconditioned forward operators and variable metric proximity operators (see Section 7.4).

Theorem 4.1 provides a control of convergence in expectation for 3P-SPIDER, which explicitly identifies the impact of non-exact preconditioned forward operators, and the impact of initialization strategies. First, we prove that the learning rate of 3P-SPIDER can be chosen constant over iterations when the preconditioned forward operator is exact or replaced with an unbiased random oracle; and is decreasing along iterations when it is replaced with a biased oracle (deterministic or random). Second, we provide the first convergence result for a stochastic VMFB algorithm addressing c1, c2 and c3 for non-convex finite sum composite optimization. For example, it is the first result for incremental EM with a non-smooth penalty term ( $g \neq 0$ ) and possibly biased Monte Carlo approximations of the E-step.

When the forward operator  $h_i(s, B(s))$  is exact, we study the complexity of 3P-SPIDER (see Corollary 4.3): in order to satisfy the  $\epsilon$ -approximate stationary condition, the number of calls  $\mathcal{K}_{\bar{h}}$  to one of the operator  $h_i(s, B(s))$  is  $O(\sqrt{n}/\epsilon)$ , the number of calls  $\mathcal{K}_{\text{prox}}$  to the backward operator is  $O(1/\epsilon)$  and, the learning rate can be chosen independent of the accuracy  $\epsilon$ . Applied to the Gradient method and applied to the EM method when there are no constraints ( $g = 0$ ), these explicit controls of convergence retrieve previous results in the literature (see e.g. Wang et al (2019); Fort et al (2020)) which are known to be at the state-of-the-art.

Finally we show that this complexity analysis remains valid when the forward operators are approximated. In the difficult case when the approximations are biased random oracles based on Monte Carlo sums, we show that  $\mathcal{K}_{\bar{h}}$  and  $\mathcal{K}_{\text{prox}}$  are not impacted by the approximation and are the same as with exact operators  $h_i(\cdot, B(\cdot))$ , by choosing an adequate number of terms in the Monte Carlo sums. The price to be paid is a Monte Carlo complexity of order  $O(\sqrt{n}/\epsilon^2)$ .

In Section 5, 3P-SPIDER is applied to inference in a logistic regression model with random effects. We show how the problem is of the form (1)-(2). In this example, the preconditioned forward operators are approximated by a Monte Carlo sum computed from a Markov chain Monte Carlo

sampler. Through numerical analyses in the case 3P-SPIDER is a stochastic Expectation Maximization algorithm in the statistic space, we discuss the choice of design parameters. We also show how the SPIDER variance reduction effect can be increased by exploiting the Monte Carlo approximations of the forward operator.

The proofs are given in Section 6 and Section 7; technical details are also provided in Appendix.

**Notations.** We denote by  $\langle \cdot, \cdot \rangle$  the dot product on  $\mathbb{R}^q$ , and by  $\| \cdot \|$  the associated norm. For a  $q \times q$  positive definite matrix  $B$ , we set  $\langle \cdot, \cdot \rangle_B \stackrel{\text{def}}{=} \langle B \cdot, \cdot \rangle$  and  $\| \cdot \|_B$  the associated norm.  $I_q$  denotes the  $q \times q$  identity matrix. For a matrix  $B$ ,  $B^\top$  is its transpose.  $\mathcal{P}_+^q$  denotes the set of the  $q \times q$  positive definite matrices.

$\mathbb{N}$  (resp.  $\mathbb{N}^*$ ) is the set of non negative (resp. positive) integers. For  $n \in \mathbb{N}^*$ , we set  $[n] \stackrel{\text{def}}{=} \{0, \dots, n\}$  and  $[n]^* \stackrel{\text{def}}{=} \{1, \dots, n\}$ .  $\mathbb{R}_+$  is the set of the positive real numbers. For  $q \in \mathbb{R}$ ,  $[q]$  is the upper integer part.

$I$  is the identity function. For a proper function  $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$ ,  $\partial g(s)$  denotes the sub-differential at  $s$ . For a continuously differentiable function  $W$  at  $s \in \mathbb{R}^q$ ,  $\nabla W(s)$  is the gradient of  $W$  at  $s$ .

All the random variables (r.v.) are defined on  $(\Omega, \mathcal{A}, \mathbb{P})$ ; for a r.v.  $U$ ,  $\sigma(U)$  is the sigma-field generated by  $U$ .

## 2 Motivating examples

In this section, we show that the Variable Metric Proximal-Gradient algorithm and the Expectation Maximization algorithm are examples of the general framework described by (1) and (2). In the first case, the preconditioning matrices  $B$  are chosen by the user, while in the second case, they are supplied by the algorithm.

### 2.1 Variable Metric Proximal-Gradient algorithms

Consider the non-convex composite problem with finite sum structure

$$\text{find } s \in \mathbb{R}^q: \quad 0 \in \frac{1}{n} \sum_{i=1}^n \nabla W_i(s) + \partial g(s), \quad (3)$$

where  $g$  is a proper lower semicontinuous convex function with domain  $\mathcal{S}$ ; and for all  $i \in [n]^*$ ,  $W_i$  is continuously differentiable on  $\mathcal{S}$ . It is of the form (1) with  $G_i \stackrel{\text{def}}{=} \nabla W_i$  being a gradient. (3) is an example of the more general problem: finding a zero on  $\mathbb{R}^q$  of the sum of two (set-valued) operators  $0 \in \text{As} + \text{Cs}$ . Here,  $\text{A} \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \nabla W_i$  and  $\text{C} \stackrel{\text{def}}{=} \partial g$  is a maximally monotone operator (see e.g. (Bauschke and Combettes, 2011, Theorem 20.25)). (3) can be solved by Forward-Backward splitting algorithms (see e.g. Combettes and Wajs (2005); Beck (2017)): the forward step uses the gradient of some if not all the functions  $W_i$ 's at each iteration; the backward step uses a proximity operator associated to  $g$ . This yields Proximal-Gradient based algorithms.

In the case  $g = 0$ , which includes unconstrained optimization problems, stochastic gradient methods with variance reduction were proposed in the situation

$$n^{-1} \sum_{i=1}^n G_i(s) = \mathbb{E}[\mathcal{G}(Z, s)] \quad (4)$$

and random oracles  $\mathcal{G}(Z, s)$  are available; in the non-convex setting, let us cite e.g. Ghadimi and Lan (2013); Reddi et al (2016); Allen-Zhu and Hazan (2016); Nguyen et al (2017); Allen-Zhu (2018); Fang et al (2018); Zhou et al (2020). These algorithms address the problem (1)-(2) by choosing  $B$  equal to the identity matrix  $I_q$  and they use unbiased random oracles  $\mathcal{G}(Z, s)$  for the approximation of the forward operator.

For non-convex composite optimization ( $g \neq 0$ ), let us cite Ghadimi et al (2016) and Karimi et al (2016) for stochastic Proximal-Gradient algorithms using unbiased oracles  $\mathcal{G}(Z, s)$  (see (4)). Li and Li (2018), Wang et al (2019), Zhang and Xiao (2019), Nhan et al (2020) and Metel and Takeda (2021) propose stochastic Proximal-Gradient methods with unbiased random oracles and including variance reduction schemes. In Metel and Takeda (2021),  $g$  may be non-convex but admits an efficiently computable proximity operator. Atchadé et al (2017) allow for deterministic or random approximations of the forward operator  $n^{-1} \sum_{i=1}^n G_i(s)$ ; when the perturbation is stochastic, the convergence analysis covers both biased and unbiased oracles, includes Nesterov

acceleration schemes, but is restricted to convex optimization. Here again, all these algorithms address the problem (1)-(2) by choosing  $B = I_q$ .

Forward-Backward suffers from slow convergence, and Variable Metric Forward-Backward (VMFB) methods were proposed by Chen and Rockafellar (1997) in order to accelerate the convergence (see also refs. 11 to 16 in Chouzenoux et al (2014)). VMFB changes the metric at each iteration by using symmetric positive definite scaling matrices multiplying the forward operator. It is an alternative to inertial methods such as Heavy Ball or Nesterov acceleration which use informations from the previous iterates. When solving the inclusion (3), VMFB uses preconditioned gradients with an iteration-dependent preconditioning matrix  $B_t^{-1}$  for the forward step at iteration  $\#t$ , and scales the proximal step consequently. Examples showing that VMFB is more efficient than Forward-Backward and Forward-Backward with inertial schemes, are provided in Chouzenoux et al (2014) and Repetti et al (2014). Different strategies exist for the definition of the variable metric  $B_t$ ; for example, it may be a diagonal matrix depending on the past history of the algorithm (see e.g. Park et al (2019) and references therein for variable scalar metrics; see also Chen et al (2019) in the case  $g = 0$ ), or inherited from Newton-type methods (see e.g. Becker and Fadili (2012); Lee et al (2014); Becker et al (2019) for the composite convex case; see also Kolte et al (2015); Moritz et al (2016); Gower et al (2016) for the smooth convex case with finite sum structure; finally, see Zhang et al (2022) for the smooth non-convex case with finite sum structure), or defined through a Majorize-Minimize strategy to make the backward operator explicit (see e.g. Chouzenoux et al (2014) and Repetti and Wiaux (2021)). Convergence results for VMFB exist in the convex case (see e.g. Combettes and Vũ (2014) and Bonettini et al (2021); and Park et al (2019) for the strongly convex case) and in the non-convex case (see e.g. Chouzenoux et al (2014) and Repetti and Wiaux (2021)). In Yun et al (2021), a stochastic VMFB is studied in the non-convex case; the exact gradient is approximated by a linear combination of random oracles, with exponential forgetting, and the oracles are assumed unbiased and bounded.

3P-SPIDER addresses non-convex composite optimization with a finite sum structure by using Proximal-Gradient algorithms accelerated via the

Variable Metric cuning. It is a stochastic VMFB, which contains a variance reduction technique in order to overcome the finite sum setting; it also allows oracles for the preconditioned forward operators, oracles which can be biased or unbiased when random (see e.g. [Atchadé et al \(2017\)](#) and [Fort et al \(2018\)](#) for examples motivating biased random approximations of the gradient). The combination of these two sources of perturbations is an original setting which, to our best knowledge, is not addressed in the literature.

3P-SPIDER uses preconditioning matrices, which may depend on the current value of the iterate and therefore may be random. The non-asymptotic convergence analysis derived in Section 4 will rely on weaker minorization assumptions on the spectrum of the scaling matrices (see [A 3](#)) than in [Yun et al \(2021\)](#); it will not require ordering assumptions on the sequence of scaling matrices as in [Combettes and Vũ \(2014\)](#) and [Bonettini et al \(2021\)](#), and will not assume a Kurdyka-Lojasiewicz condition on the objective function as in [Chouzenoux et al \(2014\)](#). As a consequence, the construction of the Lyapunov inequality for the convergence analysis of 3P-SPIDER differs from these previous works.

3P-SPIDER requires the backward operator to be explicit, which may be a strong assumption especially for variable metric proximity operators (see Section 3.5); extensions of the convergence analysis to the case of inexact proximity operators is out of the scope of this paper.

## 2.2 Expectation Maximization for curved exponential families

Consider the parametric statistical model: the observations are independent with density

$$y \mapsto \int_{\mathcal{Z}} p(y, z; \theta) \mu_{lv}(dz),$$

with respect to (w.r.t.) a  $\sigma$ -finite positive measure  $\mu_o$  on  $\mathbb{R}^{d_y}$ . In this model,  $z$  acts as a *latent* variable taking values in the measurable set  $(\mathcal{Z}, \mathcal{Z})$  endowed with a  $\sigma$ -finite positive measure  $\mu_{lv}$  (see e.g. [Everitt \(1984\)](#) for examples of latent variable models). The goal is to learn the parameter  $\theta \in \Theta \subseteq \mathbb{R}^d$  from  $n$  observations  $Y_1, \dots, Y_n$ , by

minimizing the negative normalized log-likelihood

$$F(\theta) \stackrel{\text{def}}{=} -\frac{1}{n} \sum_{i=1}^n \log \int_{\mathcal{Z}} p(Y_i, z; \theta) \mu_{lv}(dz) \quad (5)$$

on  $\Theta$ . Unfortunately, this is a non-convex problem and most often, an optimization algorithm for the minimization of (5) can not do better than converging to a critical point of the objective function (see [Wu \(1983\)](#)).

A popular model is the case when the *complete data likelihood*  $(y, z) \mapsto p(y, z; \theta)$  is of the form

$$p(y, z; \theta) \stackrel{\text{def}}{=} H(y, z) \exp(\langle S(y, z), \phi(\theta) \rangle - \psi(\theta))$$

where  $H : \mathbb{R}^{d_y} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ ,  $S : \mathbb{R}^{d_y} \times \mathcal{Z} \rightarrow \mathbb{R}^q$ ,  $\phi : \Theta \rightarrow \mathbb{R}^q$ ,  $\psi : \Theta \rightarrow \mathbb{R}$ ; it corresponds to the so-called *curved exponential family* assumption. It is satisfied by the mixture models as soon as the components of the mixture are from the curved exponential family. See e.g. [Brown \(1986\)](#) for an introduction to curved exponential family of distributions; and [McLachlan and Krishnan \(2008\)](#) for examples of such latent variable models.

EM for curved exponential families defines iteratively a  $\Theta$ -valued sequence  $\{\theta_t, t \geq 0\}$  through the mechanism: given  $\theta_t$ ,

- **(E-step)** Compute  $\bar{s}(\theta_t)$ , the expectation of the *sufficient statistics* w.r.t. the *a posteriori* distributions

$$\bar{s}(\theta_t) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta_t),$$

$$\bar{s}_i(\theta_t) \stackrel{\text{def}}{=} \int_{\mathcal{Z}} S(Y_i, z) \frac{p(Y_i, z; \theta_t) \mu_{lv}(dz)}{\int_{\mathcal{Z}} p(Y_i, u; \theta_t) \mu_{lv}(du)}.$$

- **(M-step)** Update the parameter

$$\theta_{t+1} \stackrel{\text{def}}{=} \operatorname{argmin}_{\theta \in \Theta} (\psi(\theta) - \langle \bar{s}(\theta_t), \phi(\theta) \rangle).$$

The algorithm alternates between a step in the parameter space  $\Theta$  (when computing  $\theta_{t+1} \in \mathbb{R}^d$ ), and a step in the *statistic space* when computing  $\bar{s}(\theta_t) \in \mathbb{R}^q$ . Proposition 2.1 states that the limiting points of EM run in the parameter space  $\Theta$  are the fixed points of an operator onto  $\Theta$ ; finding such a fixed point is equivalent to find a fixed point of an operator onto the statistic space  $\bar{s}(\Theta) \subseteq \mathbb{R}^q$ .

**Proposition 2.1** Assume that for any  $s \in \mathcal{S} \supseteq \bar{s}(\Theta)$ ,

$$\mathsf{T}(s) \stackrel{\text{def}}{=} \operatorname{argmin}_{\theta \in \Theta} (\psi(\theta) - \langle s, \phi(\theta) \rangle)$$

exists and is unique. Set  $\mathcal{L}_\theta \stackrel{\text{def}}{=} \{\theta \in \Theta : \mathsf{T}(\bar{s}(\theta)) = \theta\}$  and  $\mathcal{L}_s \stackrel{\text{def}}{=} \{s \in \bar{s}(\Theta) : \bar{s}(\mathsf{T}(s)) = s\}$ .

$\mathcal{L}_\theta$  is the set of the limiting points of EM. If  $\theta_\star \in \mathcal{L}_\theta$ , then  $s_\star \stackrel{\text{def}}{=} \bar{s}(\theta_\star)$  is in  $\mathcal{L}_s$ . Conversely, if  $s_\star \in \mathcal{L}_s$  then  $\theta_\star \stackrel{\text{def}}{=} \mathsf{T}(s_\star)$  is in  $\mathcal{L}_\theta$ .

See e.g. Delyon et al (1999) for the proof. An algorithmic corollary of Proposition 2.1, is that EM is equivalent to any algorithm run in the statistic space and designed to find the roots of the function

$$s \mapsto \frac{1}{n} \sum_{i=1}^n \bar{h}_i(s), \quad \text{where } \bar{h}_i(s) \stackrel{\text{def}}{=} \bar{s}_i(\mathsf{T}(s)) - s,$$

on the subset  $\bar{s}(\Theta)$  of  $\mathbb{R}^q$ . Under regularity conditions on the statistical model, it is proved in Delyon et al (1999) (see also a statement in (Fort et al, 2020, Proposition 1)) that there exists a  $q \times q$  positive definite matrix  $\mathsf{B}(s)$  such that

$$\nabla(F \circ \mathsf{T})(s) = -\mathsf{B}(s) \left( \frac{1}{n} \sum_{i=1}^n \bar{h}_i(s) \right), \quad (6)$$

where  $F$  is the negative normalized log-likelihood (see (5)). Therefore, the roots of  $n^{-1} \sum_{i=1}^n \bar{h}_i(s)$  on  $\bar{s}(\Theta)$  are the roots of  $n^{-1} \sum_{i=1}^n G_i(s)$  on  $\bar{s}(\Theta)$ , where  $G_i(s) \stackrel{\text{def}}{=} -\mathsf{B}(s) \bar{h}_i(s)$ . It also means that the roots of  $n^{-1} \sum_{i=1}^n \bar{h}_i(s)$  are the roots of the gradient of  $F \circ \mathsf{T}$ , the objective function transferred on the statistic space through the map  $\mathsf{T} : \Theta \rightarrow \mathbb{R}^q$ .

As a conclusion, EM in the statistic space is an example of problem (1)-(2), where the function  $g$  collects the constraint on  $s$  such as  $s \in \mathcal{S} \supseteq \bar{s}(\Theta)$ : (i) it is designed to find a root of  $n^{-1} \sum_{i=1}^n G_i(s) = \nabla(F \circ \mathsf{T})(s)$  under the constraint that  $s \in \mathcal{S}$ ; (ii) it uses the quantities  $\bar{h}_i(s)$  which are preconditioned forward operator since there exists  $\mathsf{B}(s)$  such that  $\bar{h}_i(s) = -\mathsf{B}(s)^{-1} G_i(s)$  (see (6)); (iii) this preconditioned forward operator is intractable for at least two reasons: first, due to the inner integrations on the set  $\mathsf{Z}$  when computing  $\bar{s}_i$ , since  $p(y, z; \theta)$  and  $\mathsf{Z}$  are often too complex to make the integrals explicit; second, due to the outer integration on the  $n$  examples

when computing  $\bar{s}$ , which has a prohibitive computational cost in large scale learning. However, a Monte Carlo approximation of  $\bar{h}_i(s)$  is always possible, whatever  $i$  and  $s$ . This remark is the cornerstone to understand the stochastic versions of EM (see e.g. Celeux and Diebolt (1985); Wei and Tanner (1990); Delyon et al (1999); Fort and Moulines (2003) which address the inner sum intractability; and Neal and Hinton (1998); Ng and McLachlan (2003); Cappé and Moulines (2009); Chen et al (2018); Karimi et al (2019); Fort et al (2020, 2021a) for the outer sum intractability). They consist in running a Stochastic Approximation (SA) algorithm with mean field  $n^{-1} \sum_{i=1}^n \bar{h}_i(s)$  (for an introduction to SA, see e.g. Benveniste et al (1990) or Borkar (2008)); this yields SA within EM procedures. They differ through the construction of the random field used for the approximation of the mean field (see (Fort et al, 2021a, Section 2.2.) for a description of some SA within EM algorithms).

3P-SPIDER is among the SA within EM algorithms. Compared to previous stochastic EM methods, it encompasses the two random approximations (of the sum in  $i$  and of the integrals on  $\mathsf{Z}$ ) and a variance reduction step, and it also allows a more general penalty term  $g$  than the  $\{0, +\infty\}$ -valued indicator function of a set.

### 3 The 3P-SPIDER algorithm

We introduce a novel algorithm named *Perturbed Proximal Preconditioned SPIDER* (3P-SPIDER), solving (1) and satisfying (c1), (c2) and (c3). It requires  $g$  to satisfy the following assumption

**A1.**  $g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$  is proper, lower semicontinuous and convex. Denote by  $\mathcal{S}$  its domain  $\mathcal{S} \stackrel{\text{def}}{=} \{s \in \mathbb{R}^q : g(s) < +\infty\}$ .

Under this assumption, we define a *variable metric* proximity operator. For any  $\gamma > 0$  and  $B \in \mathcal{P}_+^q$ , the proximity operator of the proper lower semicontinuous convex function  $\gamma g : \mathbb{R}^q \rightarrow (-\infty, +\infty]$  relative to the metric induced by  $B$  is defined by (see e.g. (Hiriart-Urruty and Lemaréchal, 1996, Section XV.4))

$$\operatorname{prox}_{\gamma g}^B(s) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbb{R}^q} \left( \gamma g(\cdot) + \frac{1}{2} \|\cdot - s\|_B^2 \right). \quad (7)$$

When  $B = I_q$ , we simply write  $\text{prox}_{\gamma g}(s)$ , which is the proximity operator originally defined by [Moreau \(1965\)](#). Lemma 3.1 shows that under A1,  $\text{prox}_{\gamma g}^B(s)$  exists and is unique for all  $s \in \mathbb{R}^q$ ,  $\gamma > 0$  and  $B \in \mathcal{P}_+^q$ . It also provides characterizations of this point. Its proof is in Section 6.1.

**Lemma 3.1** *Assume A1.*

1. For any  $\gamma > 0$ ,  $B \in \mathcal{P}_+^q$  and  $s \in \mathbb{R}^q$ , the optimization problem (7) has a unique solution, characterized as the unique point  $\mathbf{p} \in \mathcal{S}$  satisfying

$$-\gamma^{-1} B(\mathbf{p} - s) \in \partial g(\mathbf{p}) .$$

2. For any  $\gamma > 0$ ,  $B \in \mathcal{P}_+^q$ ,  $s \in \mathcal{S}$  and  $h \in \mathbb{R}^q$ ,

$$s = \text{prox}_{\gamma g}^B(s + \gamma h) \quad \text{iff} \quad Bh \in \partial g(s). \quad (8)$$

For  $s \in \mathbb{R}^q$  and  $B \in \mathcal{P}_+^q$ , set

$$\begin{aligned} h_i(s, B) &\stackrel{\text{def}}{=} -B^{-1} G_i(s) , \\ h(s, B) &\stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n h_i(s, B) . \end{aligned} \quad (9)$$

By Lemma 3.1-item 2, it holds for any  $B \in \mathcal{P}_+^q$  and  $\gamma > 0$ :  $s = \text{prox}_{\gamma g}^B(s + \gamma h(s, B))$  iff  $-Bh(s, B) \in \partial g(s)$ . By (9), this yields for any  $B \in \mathcal{P}_+^q$ , and  $\gamma > 0$ :

$$\begin{aligned} s_* &= \text{prox}_{\gamma g}^B(s_* + \gamma h(s_*, B)) \\ &\quad \text{iff } s_* \text{ solves (1)}. \end{aligned} \quad (10)$$

### 3.1 Variable Metric Proximal and Preconditioned Gradient

(10) shows that when solving the composite optimization problem (1), as soon as a preconditioned version of the operator  $s \mapsto n^{-1} \sum_{i=1}^n G_i(s)$  is used – with preconditioning matrix  $B^{-1}$ , a proximity operator of  $g$  relative to a metric induced by the matrix  $B$  has to be used.

Based on the characterization (10), a natural splitting algorithm to solve (1) under the condition (c1) is: given  $\sigma_0 \in \mathcal{S}$ , a positive stepsize sequence  $\{\gamma_{k+1}, k \geq 0\}$  and a  $\mathcal{P}_+^q$ -valued sequence

$\{B_{k+1}, k \geq 0\}$ , repeat

$$\sigma_{k+1} = \text{prox}_{\gamma_{k+1} g}^{B_{k+1}}(\sigma_k + \gamma_{k+1} h(\sigma_k, B_{k+1})) . \quad (11)$$

It corresponds to the Variable Metric Forward-Backward algorithm (see e.g. [Chen and Rockafellar \(1997\)](#); [Combettes and Vũ \(2014\)](#)).

In the large scale learning setting, the full sum over the  $n$  functions  $h_i$  (see (9)) can not be computed at each iteration of (11). In addition, it may happen that  $h_i(s)$  is not explicit (see e.g. the case of the incremental EM algorithms, Section 2.2). Therefore, a natural idea is to propose the inexact version of (11) defined by Algorithm 1: the proximal step is unchanged (see line 8); the SA step in line 7 uses a random approximation  $\mathbf{S}_{k+1}$  of the exact mean field  $n^{-1} \sum_{i=1}^n h_i(\hat{S}_k)$ ; this approximation, defined by line 6, combines a mini-batch approximation of a full sum (see line 3) and possibly approximated terms  $\delta_{k+1,i}$  (see line 5).

---

**Algorithm 1** A stochastic Variable Metric Forward-Backward

---

**Require:**  $k^{\text{out}} \in \mathbb{N}^*$ ,  $\gamma_k > 0$  for  $k \in [k^{\text{out}}]^*$ ,  $\mathbf{b} \in \mathbb{N}^*$ ,  $\hat{S}_{\text{init}} \in \mathcal{S}$ .

**Ensure:** The sequence  $\{\hat{S}_k, k \in [k^{\text{out}}]\}$ .

- 1:  $\hat{S}_0 = \hat{S}_{\text{init}}$
  - 2: **for**  $k = 0, \dots, k^{\text{out}} - 1$  **do**
  - 3:     Sample a batch  $\mathcal{B}_{k+1}$  of size  $\mathbf{b}$  in  $[n]^*$
  - 4:     Choose  $B_{k+1} \in \mathcal{P}_+^q$
  - 5:     For  $i \in \mathcal{B}_{k+1}$ , compute an approximation  $\delta_{k+1,i}$  of  $h_i(\hat{S}_k, B_{k+1})$ .
  - 6:      $\mathbf{S}_{k+1} = \mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{k+1}} \delta_{k+1,i}$
  - 7:      $\hat{S}_{k+1/2} = \hat{S}_k + \gamma_{k+1} \mathbf{S}_{k+1}$
  - 8:      $\hat{S}_{k+1} = \text{prox}_{\gamma_{k+1} g}^{B_{k+1}}(\hat{S}_{k+1/2})$ .
  - 9: **end for**
- 

### 3.2 The SPIDER variance reduction technique

3P-SPIDER leverages on Algorithm 1 and on the variance reduction technique SPIDER for the definition of the field  $\mathbf{S}_{k+1}$  that approximates  $n^{-1} \sum_{i=1}^n h_i(\hat{S}_k, B_{k+1})$ . SPIDER stands for *Stochastic Path-Integrated Differential Estimator*, and was originally introduced in the stochastic gradient descent literature by [Fang et al \(2018\)](#) (see also [Nguyen et al \(2017\)](#); [Wang et al \(2019\)](#)).



We give the intuition of SPIDER in the SA setting which encompasses the stochastic gradient one.

SA scheme solves a root finding problem  $\xi(s) = 0$  on  $\mathbb{R}^q$  by: given an initial value  $s_0 \in \mathbb{R}^q$  and a stepsize sequence  $\{\gamma_{k+1}, k \geq 0\}$ , repeat  $s_{k+1} = s_k + \gamma_{k+1} \Xi_{k+1}$ , where at each iteration  $\#(k+1)$ ,  $\Xi_{k+1}$  is a random approximation of  $\xi(s_k)$ . Usually, it is required that conditionally to the past of the algorithm, the expectation of  $\Xi_{k+1}$  is  $\xi(s_k)$ ; in that case,  $\Xi_{k+1}$  can be replaced with  $S_{k+1} \stackrel{\text{def}}{=} \Xi_{k+1} + V_{k+1}$ , where conditionally to the past,  $V_{k+1}$  is centered. SPIDER leverages on this remark and on the *control variate* technique: it proposes a clever construction of a random variable  $V_{k+1}$  approximating zero and correlated to  $\Xi_{k+1}$ .

The recipe is as follows: consider that at iteration  $\#k$ ,  $S_k$  is a random approximation of  $h(\widehat{S}_{k-1}, B_k)$ . Then define  $S_{k+1}$  by  $S_{k+1} \stackrel{\text{def}}{=} H_{k+1} + V_{k+1}$  where

$$H_{k+1} \stackrel{\text{def}}{=} \mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{k+1}} h_i(\widehat{S}_k, B_{k+1}),$$

$$V_{k+1} \stackrel{\text{def}}{=} S_k - \mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{k+1}} h_i(\widehat{S}_{k-1}, B_k),$$

and  $\mathcal{B}_{k+1}$  is sampled at random in  $[n]^*$ . The r.v.  $V_{k+1}$  approximates zero since both  $\mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{k+1}} h_i(\widehat{S}_{k-1}, B_k)$  and  $S_k$  approximate  $n^{-1} \sum_{i=1}^n h_i(\widehat{S}_{k-1}, B_k)$ ;  $V_{k+1}$  and  $H_{k+1}$  are correlated via  $\mathcal{B}_{k+1}$ .

Unfortunately, the r.v.  $S_{k+1}$  is not an unbiased approximation of  $n^{-1} \sum_{i=1}^n h_i(\widehat{S}_k, B_{k+1})$  (see Proposition 7.3 in the case  $B_{k+1}$  is of the form  $B(\widehat{S}_k)$ ). In order to remove the bias, SPIDER restarts the control variate mechanism regularly: every  $k^{\text{in}}$  iterations, compute a full sum over the  $n$  terms and set  $S_{k^{\text{in}}+1} = n^{-1} \sum_{i=1}^n h_i(\widehat{S}_{k^{\text{in}}}, B_{k^{\text{in}}+1})$ .

### 3.3 3P-SPIDER

3P-SPIDER is given by Algorithm 2. The iteration index is  $(t, k)$  where  $t$  is the index of the current *outer* loop and ranges from 1 to  $k^{\text{out}}$ , and  $k$  is the index of the current *inner* loop. At outer loop  $\#t$ , there are  $k_t^{\text{in}}$  inner iterations. The inner iterations are Algorithm 1 (see Lines 8, 9, 12 and 13 of Algorithm 2) combined with the SPIDER variance reduction trick (see Line 11 of

Algorithm 2) adapted to the case when the quantities  $h_i(\widehat{S}_{t,k}, B_{t,k+1}) - h_i(\widehat{S}_{t,k-1}, B_{t,k})$  can not be computed exactly (see Line 10).

When  $G_i$  is a gradient and  $B(s) = \mathbf{I}_q$ , different strategies were proposed for SPIDER for the choice of  $\mathbf{b}'_t$  and  $k_t^{\text{in}}$ . In Fang et al (2018); Nguyen et al (2017); Wang et al (2019), the number of inner loops is constant ( $k_t^{\text{in}} = k^{\text{in}}$  for any  $t \geq 1$ ) and  $\mathbf{b}'_t = n$ ; Nguyen et al (2017) also considers the case when  $k_t^{\text{in}}$  is adapted based on the history of the algorithm while being upper bounded; in Horváth et al (2022),  $\mathbf{b}'_t$  is deterministic and depends on  $t$ ,  $\mathbf{b}$  depends on  $t$ , and  $k_t^{\text{in}}$  is a Geometric random variable with an expectation depending on  $t$ ; in Li et al (2021),  $\mathbf{b}'_t$  does not depend on  $t$  and  $k_t^{\text{in}}$  is random.

For the EM problem (see Section 2.2), Fort et al (2020) introduced SPIDER-EM, a variance reduced stochastic EM designed for large scale learning, in a situation when the computation of  $\bar{h}_i(s)$  is exact for all  $s, i$ . For this algorithm, the benefit of an increasing batch size  $t \mapsto \mathbf{b}'_t$  and a geometric number of inner loops  $k_t^{\text{in}}$  with time-varying expectation, is discussed in Fort et al (2021b). The conclusion is that the best strategy is a deterministic increasing sequence  $\mathbf{b}'_t$  in order to have an increasing accuracy when refreshing the variable  $S$ , and a constant number of inner loops  $k_t^{\text{in}} = k^{\text{in}}$ . This paper allows  $\mathbf{b}'_t$  and  $k_t^{\text{in}}$  to vary with  $t$ : they may be deterministic functions of  $t$  or random ones as well.

The matrices  $\{B_{t,k+1}, t \in [k^{\text{out}}]^*, k \in [k^{\text{in}} - 1]\}$  can be deterministic or random. They could be chosen prior the run of the algorithm; more efficient strategies consist in adapting this matrix along the run of the algorithm, based on its history. In EM (see Section 2.2),  $B_{t,k+1}$  is of the form  $B(\widehat{S}_{t,k})$  where  $B$  is defined by the statistic model.

After  $k_t^{\text{in}}$  inner iterations, the outer loop  $\#(t+1)$  starts: the stochastic mean field  $S_{t+1,0}$  is refreshed (see Line 3 to Line 6). Here again, two approximations of the original SPIDER algorithm are allowed: the first one is when computing  $h_i(\widehat{S}_{t,0}, B_{t,1})$  and the second one avoids the scan of the full data set (one may choose  $\mathbf{b}'_t < n$ ).

The input variables of 3P-SPIDER are the number of outer loops  $k^{\text{out}}$ , the number of inner loops  $k_t^{\text{in}}$ , the stepsize sequence  $\{\gamma_{t,k}, t \geq 1, k \geq 1\}$  for the SA steps, the size of the mini-batches  $\mathbf{b}$  and

---

**Algorithm 2** The Perturbed Proximal Preconditioned SPIDER algorithm (3P-SPIDER)
 

---

**Require:**  $k^{\text{out}} \in \mathbb{N}^*$ ,  $k_t^{\text{in}} \in \mathbb{N}^*$  for  $t \in [k^{\text{out}}]^*$ ,  $\gamma_{t,k+1} > 0$  for  $t \in [k^{\text{out}}]^*$ ,  $k \in [k_t^{\text{in}}]$ ,  $\mathbf{b} \in \mathbb{N}^*$ ,  $\mathbf{b}'_t \in \mathbb{N}^*$  for  $t \in [k^{\text{out}}]^*$ ,  $\widehat{S}_{\text{init}} \in \mathcal{S}$  and  $B_{\text{init}} \in \mathcal{P}_+^q$

**Ensure:** The sequence  $\{\widehat{S}_{t,k}, t \in [k^{\text{out}}]^*, k \in [k_t^{\text{in}}]^*\}$ .

- 1:  $\widehat{S}_{0,k_0^{\text{in}}} = \widehat{S}_{\text{init}}, B_{0,k_0^{\text{in}}} = B_{\text{init}}$
  - 2: **for**  $t = 1, \dots, k^{\text{out}}$  **do**
  - 3:    $\widehat{S}_{t,0} = \widehat{S}_{t-1,k_{t-1}^{\text{in}}}, \widehat{S}_{t,-1} = \widehat{S}_{t-1,k_{t-1}^{\text{in}}}, B_{t,0} = B_{t-1,k_{t-1}^{\text{in}}}$
  - 4:   Sample a batch  $\mathcal{B}_{t,0}$  of size  $\mathbf{b}'_t$  in  $[n]^*$ , with or without replacement.
  - 5:   For all  $i \in \mathcal{B}_{t,0}$ , compute  $\delta_{t,0,i}$  equal to or approximating  $\mathbf{h}_i(\widehat{S}_{t,0}, B_{t,0})$ .
  - 6:    $\mathbf{S}_{t,0} = (\mathbf{b}'_t)^{-1} \sum_{i \in \mathcal{B}_{t,0}} \delta_{t,0,i}$
  - 7:   **for**  $k = 0, \dots, k_t^{\text{in}} - 1$  **do**
  - 8:     Sample a mini batch  $\mathcal{B}_{t,k+1}$  of size  $\mathbf{b}$  in  $[n]^*$ , with or without replacement.
  - 9:     Choose  $B_{t,k+1} \in \mathcal{P}_+^q$ .
  - 10:     For all  $i \in \mathcal{B}_{t,k+1}$ , compute  $\delta_{t,k+1,i}$  equal to or approximating  $\mathbf{h}_i(\widehat{S}_{t,k}, B_{t,k+1}) - \mathbf{h}_i(\widehat{S}_{t,k-1}, B_{t,k})$ .
  - 11:      $\mathbf{S}_{t,k+1} = \mathbf{S}_{t,k} + \mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{t,k+1}} \delta_{t,k+1,i}$
  - 12:      $\widehat{S}_{t,k+1/2} = \widehat{S}_{t,k} + \gamma_{t,k+1} \mathbf{S}_{t,k+1}$
  - 13:      $\widehat{S}_{t,k+1} = \text{prox}_{t,k}(\widehat{S}_{t,k+1/2}),$                    where  $\text{prox}_{t,k} \stackrel{\text{def}}{=} \text{prox}_{\gamma_{t,k+1} g}^{B_{t,k+1}}$ .
  - 14:   **end for**
  - 15: **end for**
- 

$\mathbf{b}'_t$ , and the initial values of the iterate  $\widehat{S}_{\text{init}}$  and the metric  $B_{\text{init}}$  in  $\mathcal{S}$  and  $\mathcal{P}_+^q$  respectively.

### 3.4 Monte Carlo approximation of $\mathbf{h}_i(\mathbf{s}, B)$

Set  $\vartheta \stackrel{\text{def}}{=} (s, i, B) \in \mathcal{S} \times [n]^* \times \mathcal{P}_+^q$ . In some applications, there exist a measurable function  $H_\vartheta$  and a probability measure  $\pi_\vartheta$  defined on the measurable set  $(Z, \mathcal{Z})$  such that

$$\mathbf{h}_i(\mathbf{s}, B) = \int_Z H_\vartheta(z) \pi_\vartheta(dz). \quad (12)$$

This is the case of EM in the statistic space (see Section 2.2) where  $H_\vartheta(z) = S(Y_i, z) - s$  and

$$\pi_\vartheta(dz) \stackrel{\text{def}}{=} \frac{p(Y_i, z; \mathbb{T}(s))}{\int_Z p(Y_i, u; \mathbb{T}(s)) \mu_{l_v}(du)} \mu_{l_v}(dz).$$

When the integral in (12) is intractable, one can resort to Monte Carlo integrations to define the approximations  $\delta_{t,k+1,i}$  and  $\delta_{t,0,i}$  (see e.g. Devroye (1986) for exact sampling methods, and Robert and Casella (2004) for an introduction to Markov chain Monte Carlo methods). If  $\{Z_m^\vartheta, m \geq 0\}$  are independent samples with distribution  $\pi_\vartheta(dz)$  or are a path of an ergodic Markov chain with unique

invariant distribution  $\pi_\vartheta(dz)$ , then we can set

$$\mathbf{h}_i(\widehat{S}_{t,k}, B_{t,k+1}) \approx \frac{1}{M} \sum_{m=1}^M H_{\vartheta_{t,k+1,i}}(Z_m^{\vartheta_{t,k+1,i}}),$$

where  $\vartheta_{t,k+1,i} \stackrel{\text{def}}{=} (\widehat{S}_{t,k}, i, B_{t,k+1})$ . We will show numerically in Section 5 that when approximating the difference  $\mathbf{h}_i(\widehat{S}_{t,k}, B_{t,k+1}) - \mathbf{h}_i(\widehat{S}_{t,k-1}, B_{t,k})$ , there is a gain in correlating the two sequences  $\{Z_m^{\vartheta_{t,k+1,i}}, m \geq 0\}$  and  $\{Z_m^{\vartheta_{t,k,i}}, m \geq 0\}$ ; this makes stronger the effect of the SPIDER control variate (see Section 3.2).

### 3.5 The computation of $\text{prox}_{\gamma g}^B$

When  $g = 0$ ,  $\text{prox}_{\gamma g}^B(s) = s$ . When  $g \neq 0$ ,  $\mathbf{p} \stackrel{\text{def}}{=} \text{prox}_{\gamma g}^B(s)$  solves  $0 \in \mathbf{p} - s + \gamma B^{-1} \partial g(\mathbf{p})$  and there does not always exist an explicit expression of  $\mathbf{p}$ .

When  $B = \mathbf{I}_q$ , (Combettes and Pesquet, 2011, Tables 10.1 and 10.2) provide properties of  $\text{prox}_{\gamma g}$  and expressions of proximity operators for many functions  $g$ .

When  $B$  is the sum of a diagonal matrix and of a rank one matrix, (Becker and Fadili, 2012, Section 3) presents iterative algorithms for the computation of  $\mathbf{p}$ . For a general positive definite matrix  $B$ , we have from (Combettes and Vũ, 2014,

Example 3.9)

$$\text{prox}_{\gamma g}^B(s) = \sqrt{B}^{-1} \text{prox}_{\gamma g(\sqrt{B}^{-1}\cdot)}(\sqrt{B}s),$$

where  $\sqrt{B}$  is the square root of the matrix  $B$ . (Becker and Fadili, 2012, Lemma 5) (see also Combettes and Vũ (2014)) establishes a Moreau identity i.e. an expression of  $\text{prox}_{\gamma g}^B$  as a function of a proximity operator of the Fenchel conjugate of  $g$ .

In the special case  $g$  is the  $\{0, +\infty\}$ -valued indicator function of a closed convex set  $\mathcal{S}$ , the *projected Landweber* method is an iterative algorithm for the computation of  $\mathbf{p}$  (see Eicke (1992), see also (Combettes and Pesquet, 2011, Example 10.10)).

Finally, for applications including a metric selection step, metric selection strategies for the definition of  $B$  can be found in (Park et al, 2019, Section 3) for diagonal variable metrics; and in Repetti and Wiaux (2021) for specific functions  $g$  which circumvent the often challenging computation of  $\text{prox}_{\gamma g}^B$ .

## 4 Non-asymptotic convergence analysis

This section is devoted to explicit non-asymptotic bounds for the convergence in expectation of 3P-SPIDER. We will restrict to the case there exist  $B : \mathcal{S} \rightarrow \mathcal{P}_+^q$  and

$$B_{t,k+1} \stackrel{\text{def}}{=} B(\widehat{S}_{t,k}).$$

This framework encompasses the EM problem (see Section 2.2) and any preconditioned gradient-based algorithms (see Section 2.1) when the preconditioning matrix depends on the past history of the algorithm via the current value of the iterate. We will also use the notation

$$\bar{\mathbf{h}}_i(s) \stackrel{\text{def}}{=} \mathbf{h}_i(s, B(s)), \quad \bar{\mathbf{h}}(s) \stackrel{\text{def}}{=} \mathbf{h}(s, B(s)). \quad (13)$$

3P-SPIDER is designed to solve (1) under the constraints c1 to c3. Therefore, based on (10), we are interested in a control of the quantities  $\text{prox}_{t,k} \left( \widehat{S}_{t,k} + \gamma_{t,k+1} \bar{\mathbf{h}}(\widehat{S}_{t,k}) \right) - \widehat{S}_{t,k}$  where

$$\text{prox}_{t,k}(s) \stackrel{\text{def}}{=} \text{prox}_{\gamma_{t,k+1} g}^{B(\widehat{S}_{t,k})}(s).$$

Roughly speaking, these quantities evaluate how far the algorithm is from the limiting set at iteration  $\#(t, k)$ . More precisely, we will control the cumulative "distances to stationary"  $\sum_{t=1}^{k^{\text{out}}} \sum_{k=0}^{k_t^{\text{in}}-1} \Delta_{t,k+1}^*$  where  $\Delta_{t,k+1}^*$  is equal to

$$\frac{\|\text{prox}_{t,k}(\widehat{S}_{t,k} + \gamma_{t,k+1} \bar{\mathbf{h}}(\widehat{S}_{t,k})) - \widehat{S}_{t,k}\|_{B(\widehat{S}_{t,k})}^2}{\gamma_{t,k+1}^2}; \quad (14)$$

$\mathbf{h}$  is defined by (9).

The controls in expectation of the cumulated distances are obtained under the assumptions A 2 to A 4. A 2 is a smoothness assumption on the functions  $\mathbf{h}_i$ , A 3 assumes that  $n^{-1} \sum_{i=1}^n G_i(s)$  is a gradient operator of some so-called *Lyapunov function*, and the spectrum of the matrices  $B(s)$  are bounded uniformly in  $s$ . A 4 are assumptions on the approximations  $\delta_{t,k+1,i}$ .

**A 2.** For all  $i \in [n]^*$ , the function  $\bar{\mathbf{h}}_i$  is globally Lipschitz on  $\mathcal{S}$ , with constant  $L_i$ : there exists a positive constant  $L_i$  such that  $\forall s, s' \in \mathcal{S}$ ,  $\|\bar{\mathbf{h}}_i(s) - \bar{\mathbf{h}}_i(s')\| \leq L_i \|s - s'\|$ . Set  $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$ .

A 2 only requires a Lipschitz property on  $\mathcal{S}$ ; it is weaker than assuming the Lipschitz property on the full space  $\mathbb{R}^q$  as sometimes assumed in the literature (see e.g. Combettes and Wajs (2005)). A 2 holds for example when  $\mathcal{S}$  is compact and for all  $i \in [n]^*$ , the gradient  $\nabla \mathbf{h}_i$  exists and is continuous on  $\mathcal{S}$ .

**A 3.** 1. There exists a function  $W : \mathbb{R}^q \rightarrow \mathbb{R}$ , continuously differentiable on  $\mathcal{S}$  and such that

$$\forall s \in \mathcal{S}, \quad \nabla W(s) = \frac{1}{n} \sum_{i=1}^n G_i(s);$$

in addition,  $\bar{\mathbf{h}}_i(s) = -B(s)^{-1} G_i(s)$ , where  $B(s) \in \mathcal{P}_+^q$ .

2.  $\nabla W$  is globally  $L_W$ -Lipschitz on  $\mathcal{S}$ .

3. There exist  $0 < v_{\min} \leq v_{\max} < +\infty$  such that for any  $s \in \mathcal{S}$ ,  $v_{\min} \|\cdot\|^2 \leq \|\cdot\|_{B(s)}^2 \leq v_{\max} \|\cdot\|^2$ .

Here again, both the Lipschitz property and the boundedness condition on the spectrum of the matrices  $B(s)$  are required on  $\mathcal{S}$  and not on the full space  $\mathbb{R}^q$ . When  $B(s)$  does not depend on

$s$  ( $\mathbb{B}(s) = B$  for any  $s \in \mathcal{S}$ ), we have  $L_{\hat{W}} \leq v_{\max} n^{-1} \sum_{i=1}^n L_i$ .

The last assumption is on the fluctuations of the errors when approximating  $\bar{h}_i(\hat{S}_{t,k}) - \bar{h}_i(\hat{S}_{t,k-1})$ : set  $\xi_{t,k+1,i} \stackrel{\text{def}}{=} \delta_{t,k+1,i} - \bar{h}_i(\hat{S}_{t,k}) + \bar{h}_i(\hat{S}_{t,k-1})$  and define its conditional bias and variance, conditionally to the  $\sigma$ -field generated by  $\mathcal{B}_{t,k+1}$ ,  $\hat{S}_{t,k}$  and  $\hat{S}_{t,k-1}$ . Set  $\mathcal{P}_{t,k+1/2} \stackrel{\text{def}}{=} \sigma(\mathcal{B}_{t,k+1}, \hat{S}_{t,k}, \hat{S}_{t,k-1})$ .

$$\begin{aligned} \mu_{t,k+1,i} &\stackrel{\text{def}}{=} \mathbb{E}[\xi_{t,k+1,i} | \mathcal{P}_{t,k+1/2}] \\ \sigma_{t,k+1,i}^2 &\stackrel{\text{def}}{=} \mathbb{E}[\|\xi_{t,k+1,i} - \mu_{t,k+1,i}\|^2 | \mathcal{P}_{t,k+1/2}]. \end{aligned}$$

We assume

- A4.** 1. Conditionally to  $\mathcal{B}_{t,k+1}$ ,  $\hat{S}_{t,k}$  and  $\hat{S}_{t,k-1}$ , the approximations  $\{\delta_{t,k+1,i}, i \in \mathcal{B}_{t,k+1}\}$  are independent.  
2. There exists a non negative constant  $C_b$  and for any  $t \in [k^{\text{out}}]^*$ , there exists a non decreasing deterministic sequence  $\{m_{t,k}, k \geq 1\}$  such that for any  $k \in [k_t^{\text{in}} - 1]$ , with probability one,

$$\left\| \frac{1}{n} \sum_{i=1}^n \mu_{t,k+1,i} \right\| \leq \frac{C_b}{m_{t,k+1}}.$$

3. There exist non negative constants  $C_v$  and  $C_{vb}$  and for any  $t \in [k^{\text{out}}]^*$ , there exist non decreasing deterministic sequences  $\{M_{t,k}, k \geq 1\}$  and  $\{\bar{M}_{t,k}, k \geq 1\}$  such that for any  $k \in [k_t^{\text{in}} - 1]$ , with probability one,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \sigma_{t,k+1,i}^2 &\leq \frac{C_v}{M_{t,k+1}}, \\ \frac{1}{n} \sum_{i=1}^n \left\| \mu_{t,k+1,i} - \frac{1}{n} \sum_{j=1}^n \mu_{t,k+1,j} \right\|^2 &\leq \frac{C_{vb}^2}{\bar{M}_{t,k+1}^2}. \end{aligned}$$

We allow the errors  $\xi_{t,k+1,i}$  to be deterministic or random. When there are no errors ( $\xi_{t,k+1,i} = 0$ ) then  $C_b = C_v = C_{vb} = 0$ . When the errors are deterministic, we have  $\xi_{t,k+1,i} = \mu_{t,k+1,i}$  and  $\sigma_{t,k+1,i}^2 = 0$ . When the errors are random and unbiased, then  $\mu_{t,k+1,i} = 0$ . Therefore, some of the constants  $C_b$ ,  $C_v$  or  $C_{vb}$  can be null as summarized in Table 1.

In Section A<sub>2</sub>, we discuss how A 4 is verified in the case  $\bar{h}_i(s') - \bar{h}_i(s)$  is an expectation

	$C_b$	$C_v$	$C_{vb}$
exact	0	0	0
deterministic	$\geq 0$	0	$\geq 0$
random, unbiased	0	$\geq 0$	0
random, biased	$> 0$	$\geq 0$	$\geq 0$

**Table 1** The sign of the constants  $C_b$ ,  $C_v$ ,  $C_{vb}$  when there are no approximations on the  $\bar{h}_i(s)'$  (case *exact*), and when there are approximations.

under a distribution that may depend on  $(s, s', i)$  (see Section 3.4), and  $\delta_{t,k+1,i}$  is a Monte Carlo approximation.

Theorem 4.1 provides an explicit upper bound of the cumulative distance to stationarity as measured by  $\Delta_{t,k+1}^*$  (see (14)) along the  $\sum_{t=1}^{k^{\text{out}}} k_t^{\text{in}}$  iterations of the algorithm. It also provides an upper bound on the cumulative errors  $\mathcal{D}_{t,k+1}^*$  defined by

$$\frac{\|\hat{S}_{t,k+1} - \text{prox}_{t,k}(\hat{S}_{t,k} + \gamma_{t,k+1} \bar{h}(\hat{S}_{t,k}))\|_{\mathbb{B}(\hat{S}_{t,k})}^2}{\gamma_{t,k+1}^2},$$

where  $\bar{h}$  is defined by (13). Given the current iterate  $\hat{S}_{t,k}$ ,  $\mathcal{D}_{t,k+1}^*$  compares two iterations: the ideal one  $\text{prox}_{t,k}(\hat{S}_{t,k} + \gamma_{t,k+1} \bar{h}(\hat{S}_{t,k}))$  and the available one  $\text{prox}_{t,k}(\hat{S}_{t,k} + \gamma_{t,k+1} \mathcal{S}_{t,k+1})$ .

**Theorem 4.1** Assume A 1, A 2, A 3 and A 4. Let  $\{k_t^{\text{in}}, t \in [k^{\text{out}}]^*\}$  be a deterministic positive sequence. For any  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , define  $\Lambda_{t,k+1}$  by

$$\frac{\gamma_{t,k} L_{\hat{W}}}{v_{\min}} + \gamma_{t,k}^2 L^2 \frac{2v_{\max} k_t^{\text{in}}}{v_{\min} \mathfrak{b}} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathfrak{b}} \bar{M}_{t,k+1}} \right). \quad (15)$$

Let  $\{\hat{S}_{t,k}, t \in [k^{\text{out}}]^*, k \in [k_t^{\text{in}}]^*\}$  be the sequence given by Algorithm 2 when the stepsize sequence  $\{\gamma_{t,k+1}, t \in [k^{\text{out}}]^*, k \in [k_t^{\text{in}} - 1]\}$  satisfies

$$\gamma_{t,k+1} \left( 1 + \frac{2C_b}{m_{t,k+1}} \right) \leq \gamma_{t,k}, \quad \Lambda_{t,k+1} \in (0, 1/2). \quad (16)$$

Then,

$$\begin{aligned} &\sum_{t=1}^{k^{\text{out}}} \sum_{k=1}^{k_t^{\text{in}}} \gamma_{t,k} \left( \frac{1}{2} - \Lambda_{t,k+1} \right) \{ \mathbb{E}[\Delta_{t,k}^*] + \mathbb{E}[\mathcal{D}_{t,k}^*] \} \\ &\leq \mathbb{E} \left[ \mathbb{W}(\hat{S}_{1,0}) + g(\hat{S}_{1,0}) \right] - \min_{\mathcal{S}} (\mathbb{W} + g) \\ &+ v_{\max} \sum_{t=1}^{k^{\text{out}}} \gamma_{t,0} k_t^{\text{in}} \mathbb{E}[\|\mathcal{E}_t\|^2] \\ &+ v_{\max} \sum_{t=1}^{k^{\text{out}}} \sum_{k=1}^{k_t^{\text{in}}} (k_t^{\text{in}} - k + 1) \gamma_{t,k} \mathcal{U}_{t,k}, \end{aligned}$$

where  $\mathcal{E}_t \stackrel{\text{def}}{=} \mathcal{S}_{t,0} - \mathfrak{h}(\widehat{\mathcal{S}}_{t,0})$  and

$$\mathcal{U}_{t,k} \stackrel{\text{def}}{=} \frac{2C_b}{m_{t,k}} + \frac{C_b^2}{m_{t,k}^2} + \frac{C_v}{\mathfrak{b}M_{t,k}} + \frac{2C_{vb}}{\sqrt{\mathfrak{b}}M_{t,k}} + \frac{C_{vb}^2}{\mathfrak{b}\bar{M}_{t,k}^2}.$$

The proof of Theorem 4.1 is given in Section 7.5. Note that  $\mathcal{U}_{t,k+1} = 0$  when the algorithm uses exact preconditioned gradients at each iteration:  $\delta_{t,0,i} = \bar{\mathfrak{h}}_i(\widehat{\mathcal{S}}_{t,0})$  and  $\delta_{t,k+1,i} = \bar{\mathfrak{h}}_i(\widehat{\mathcal{S}}_{t,k}) - \bar{\mathfrak{h}}_i(\widehat{\mathcal{S}}_{t,k-1})$  for all  $i, t, k$ .

**Random number of inner loops  $k_t^{\text{in}}$ .** When the number of inner loops  $k_t^{\text{in}}$  at the outer loop  $\#t$  is a random number, we consider it is drawn prior the run of the algorithm. Therefore the expectations in Theorem 4.1 are conditionally to the random sequence  $\{k_t^{\text{in}}, t \in [k^{\text{out}}]^*\}$ . The expectation w.r.t. the randomness of  $k_t^{\text{in}}$  can easily be obtained from Theorem 4.1; details are omitted.

**The step sizes  $\gamma_{t,k}$ .** The conditions on the sequence  $\{\gamma_{t,k+1}, t \in [k^{\text{out}}]^*, k \in [k_t^{\text{in}} - 1]\}$  are satisfied with

$$\gamma_{t,k+1} \stackrel{\text{def}}{=} \prod_{j=0}^k \left(1 + \frac{2C_b}{m_{t,j+1}}\right)^{-1} \gamma_{t,0}$$

where  $\gamma_{t,0}$  is positive and strictly lower than

$$\frac{1}{4Lv_{\max}v} \frac{\mathfrak{b}}{k_t^{\text{in}}} \left( \sqrt{\frac{L_{\mathfrak{W}}^2}{L^2} + 4v_{\min}v_{\max} \frac{k_t^{\text{in}}}{\mathfrak{b}} v} - \frac{L_{\mathfrak{W}}}{L} \right); \quad (17)$$

$v \stackrel{\text{def}}{=} 1 + 2C_{vb}/(\sqrt{\mathfrak{b}} \inf_{t,k} \bar{M}_{t,k+1})$  (see the proof in Section C.1). First, observe that when  $C_b = 0$ , the step size can be a constant function of the inner loop index  $k$  loop ( $\gamma_{t,k+1} = \gamma_{t,0}$  for any  $k$ ). On the contrary, when  $C_b > 0$  i.e. for a deterministic approximation or a biased random approximation (see Table 1), the stepsize sequence is a strictly decreasing function of the inner loop index  $k$ .

Second, the maximal value of  $\gamma_{t,0}$  is larger when  $C_{vb} = 0$  than when  $C_{vb} > 0$ . Here again, deterministic or unbiased random approximations requires more aggressive step sizes.

**The initialization of the outer loops.** Set  $\mathcal{N}_t \stackrel{\text{def}}{=} \|\mathcal{E}_t\|^2$ . When  $\mathcal{B}_{t,0} = \{1, \dots, n\}$  and  $\delta_{t,0,i} = \bar{\mathfrak{h}}_i(\widehat{\mathcal{S}}_{t,0})$  for all  $i$ , then  $\mathcal{N}_t = 0$ ; otherwise,  $\mathcal{N}_t$  is positive.

Let us discuss the behavior of  $\mathcal{N}_t$  when  $\delta_{t,0,i}$  is an unbiased random approximation of  $\bar{\mathfrak{h}}_i(\widehat{\mathcal{S}}_{t,0})$  with variance denoted by  $\sigma_{t,0,i}^2$ . When  $\mathcal{B}_{t,0} = \{1, \dots, n\}$ , then

$$\mathbb{E}[\mathcal{N}_t] = \frac{1}{n^2} \sum_{i=1}^n \sigma_{t,0,i}^2. \quad (18)$$

Nevertheless, the strategy  $\mathcal{B}_{t,0} = \{1, \dots, n\}$  has a large computational cost; sampling a subset of size  $\mathfrak{b}'_t$  reduces the computational cost but increases the squared norm of the error: we have

$$\mathbb{E}[\mathcal{N}_t] \leq \frac{1}{\mathfrak{b}'_t n} \sum_{i=1}^n (\sigma_{t,0,i}^2 + \|\bar{\mathfrak{h}}_i(s) - \bar{\mathfrak{h}}(s)\|^2), \quad (19)$$

with an equality if  $\mathcal{B}_{t,0}$  is sampled with replacement in  $\{1, \dots, n\}$ . See Section C.2 for detailed computations. From a numerical point of view, an efficient strategy consists in increasing the size  $\mathfrak{b}'_t$  with the outer loop index  $t$  (see references in Section 3.3 for 3P-SPIDER applied to EM).

**Random stopping time of the algorithm.** In non-convex optimization, the last iterate  $\widehat{\mathcal{S}}_{k^{\text{out}}, k^{\text{in}}_{k^{\text{out}}}}$  is not necessarily the point which minimizes, over the sequence  $\{\widehat{\mathcal{S}}_{t,k}, t \in [k^{\text{out}}]^*, k \in [k_t^{\text{in}}]^*\}$ , the distance to the set of solutions of (1). The quantity  $\Delta^*$ , motivated by (10), can not be computed exactly in our framework so that the "best" iterate can not be identified thanks to this criterion. It is therefore popular to analyze the algorithm when stopped at a random time (see e.g. (Lan, 2020, Chapter 6)). For sake of simplicity, we consider the case when  $k_t^{\text{in}} = k^{\text{in}}$  for any  $t$  and  $C_b = 0$ . We have the following corollary:

**Corollary 4.2** (of Theorem 4.1) *Assume that  $k_t^{\text{in}} = k^{\text{in}}$ ,  $C_b = 0$  and the stepsize sequence is constant  $\gamma_{t,k} = \gamma_*$ . Let  $(\tau, K)$  be a uniform random variable on  $[k^{\text{out}}]^* \times [k^{\text{in}}]^*$ , independent of the algorithm. Then*

$$\begin{aligned} & \inf_{(t,k) \in [k^{\text{out}}]^* \times [k^{\text{in}}]^*} \left( \frac{1}{2} - \Lambda_{t,k} \right) \mathbb{E}[\Delta_{\tau,K}^* + \mathcal{D}_{\tau,K}^*] \\ & \leq \frac{\mathbb{E}[\mathfrak{W}(\widehat{\mathcal{S}}_{1,0}) + g(\widehat{\mathcal{S}}_{1,0})] - \min_{\mathcal{S}}(\mathfrak{W} + g)}{k^{\text{out}} k^{\text{in}} \gamma_*} \\ & \quad + v_{\max} \mathbb{E}[\|\mathcal{E}_{\tau}\|^2] + v_{\max} \mathbb{E}[(k^{\text{in}} - K + 1) \mathcal{U}_{\tau,K}]. \end{aligned}$$

An upper bound on  $\Lambda_{t,k}$  can easily be obtained from (15) as a function of  $L_{\tilde{W}}, L, v_{\min}, v_{\max}, k^{\text{in}}, \mathbf{b}, \gamma_*, C_{vb}$  and  $\inf_{t,k} \bar{M}_{t,k+1}$ .

Corollary 4.2 shows that, even by stopping 3P-SPIDER with this simple rule, the first term in the RHS is inversely proportional to the maximal number of iterations  $k^{\text{out}} k^{\text{in}}$ .

**Complexity analysis when  $\mathcal{E}_t = 0, \mathcal{U}_{t,k} = 0$  and  $k_t^{\text{in}} = k^{\text{in}}$ .** For smooth first-order optimization, algorithms are compared through their complexity in order to satisfy an  $\epsilon$ -first order stationary condition. In stochastic composite optimization, this criterion is naturally extended to the  $\epsilon$ -approximate stationary condition defined by

$$\mathbb{E} [\Delta_{\tau,K}^*] \leq \epsilon,$$

where  $(\tau, K)$  is a random variable taking values in  $[k^{\text{out}}]^* \times [k^{\text{in}}]^*$ ; see e.g. (Ghadimi et al, 2016, Section 4), (Wang et al, 2019, Section 3) and Fort and Moulines (2021).

Corollary 4.3 studies the proximal complexity  $\mathcal{K}_{\text{prox}}$  defined as the number of calls to the prox operator in order to satisfy the  $\epsilon$ -approximate stationary condition; the stochastic  $\bar{h}$ -complexity  $\mathcal{K}_{\bar{h}}$  defined as the number of calls to one of the  $\bar{h}_i$ 's; and the total number of iterations  $k^{\text{in}} k^{\text{out}}$ . Again for sake of simplicity, and in order to compare our results to the literature, we consider a simplified setting.

**Corollary 4.3** (of Corollary 4.2) *Assume in addition that  $\mathcal{E}_t = 0$  and  $\mathcal{U}_{t,k} = 0$ . The  $\epsilon$ -approximate stationary condition is satisfied with  $\gamma_* = v_{\min}/(4L_{\tilde{W}})$ ,  $k^{\text{in}}/\mathbf{b} = L_{\tilde{W}}^2/(v_{\min}v_{\max}L^2)$ ,  $\mathbf{b} = O(\sqrt{n}\sqrt{v_{\min}v_{\max}}L/L_{\tilde{W}})$  and  $k^{\text{out}}k^{\text{in}} = O(L_{\tilde{W}}/(\epsilon v_{\min}))$ . Moreover,  $\mathcal{K}_{\text{prox}} = O(L_{\tilde{W}}/(v_{\min}\epsilon))$  and  $\mathcal{K}_{\bar{h}} = O(\sqrt{v_{\max}}L\sqrt{n}/(\epsilon\sqrt{v_{\min}}))$ .*

The proof is in Section 7.6. This result shows that the step size  $\gamma_*$  and the number of inner loops  $k^{\text{in}}$  are independent of the accuracy  $\epsilon$ .

When applied to Stochastic Gradient Descent, 3P-SPIDER in the setting of Corollary 4.3 is the Prox-SpiderBoost algorithm studied in Wang et al (2019): Corollary 4.3 and (Wang et al, 2019, Theorem 2) state the same complexity results. (Wang et al, 2019, Table 1) compares Prox-SpiderBoost to other stochastic gradient

algorithms for composite non-convex finite sum optimization. It is shown that the variance reduction based on SPIDER order-level outperforms other variance reduction strategies such as the SVRG one and the SAGA one, introduced respectively by Johnson and Zhang (2013) and Defazio et al (2014). Hence, 3P-SPIDER reaches the state of the art among the proximal stochastic gradient algorithms designed to solve finite sum non-convex composite optimization.

When applied to EM, 3P-SPIDER in the setting of Corollary 4.3 is the extension of the SPIDER-EM algorithm studied in Fort et al (2020) to the case there is a proximal step which manages the constraint  $g$ . Here again, the comparison of Corollary 4.3 and (Fort et al, 2020, Theorem 2) shows that 3P-SPIDER reaches the state of the art among the incremental EM algorithms with variance reduction, including sEM-VR and FIEM introduced respectively in Chen et al (2018) and Karimi et al (2019) (see also Fort et al (2021a)). See the comparison to the literature in Fort et al (2020).

Beyond these two applications, Corollary 4.3 is - to our best knowledge - the first complexity result for an algorithm designed to solve (1) under the constraint (2) and for non-convex finite sum composite optimization.

**$\epsilon$ -approximate stationary condition: the cost of inexact preconditioned forward operators.** Let us discuss the cost of inexact  $\bar{h}_i(s)$ 's when the approximation is unbiased and random (so that  $C_b = C_{vb} = 0$ , see Table 1): does it deteriorate the proximal complexity  $\mathcal{K}_{\text{prox}}$  and the number of calls to an oracle of a preconditioned forward operator  $\bar{h}_i$  (still denoted by  $\mathcal{K}_{\bar{h}}$  below) ? detailed computations of the assertions below can be found in Section 7.7.

If  $\mathbb{E} [\|\mathcal{E}_t\|^2] = O(\epsilon^{1-a'}/(\sqrt{nt})^{a'})$  for some  $a' \in [0, 1)$  and

$$M_{t,k+1} = O\left(\frac{n^{(a-\bar{a})/2}}{\epsilon^{1-a}} t^a (k+1)^{\bar{a}}\right)$$

for some  $a, \bar{a} \in [0, 1)$ , then the  $\epsilon$ -approximate stationary condition is satisfied with  $k_t^{\text{in}} = O(\sqrt{n})$ ,  $\mathbf{b} = O(\sqrt{n})$  and  $k^{\text{out}} = O(1/(\sqrt{n}\epsilon))$ . In addition,  $\mathcal{K}_{\text{prox}} = O(1/\epsilon)$  and  $\mathcal{K}_{\bar{h}} = O(\sqrt{n}/\epsilon)$ . Therefore, the conclusions of Corollary 4.3 remain valid, and the approximations of the  $\bar{h}_i$ 's do not deteriorate

the complexity performances of the algorithms, as soon as the approximation is small enough.

Let us now evaluate the computational cost, in the case the unbiased random approximation is a Monte Carlo approximation computed from independent and identically distributed (i.i.d.) samples. In this case,  $M_{t,\cdot}$  is the number of terms of the Monte Carlo sum (see Section A). The Monte Carlo complexity  $\mathcal{K}_{\text{MC}}$  defined as the total number of Monte Carlo draws required to satisfy the  $\epsilon$ -approximate stationary condition is:  $\mathcal{K}_{\text{MC}} = O(\sqrt{n}/\epsilon^2)$  for any  $\mathbf{a}, \mathbf{a}', \bar{\mathbf{a}} \in [0, 1)$ .

To our best knowledge, it is the first complexity analysis with such a Monte Carlo approximation of the preconditioned forward operators  $\bar{\mathbf{h}}_i$ 's.

## 5 Application: Penalized Logistic Regression with random effects

### 5.1 The model

Motivated by applications in classification, we consider a logistic regression model with random effects.

Let  $n$  pairs of examples  $\{(X_i, Y_i), i \in [n]^*\}$  where  $X_i \in \mathbb{R}^d$  collects the  $d$  explanatory variables, and  $Y_i$  is the binary response variable taking values in  $\{-1, 1\}$ . We assume that given  $\{X_i, i \in [n]^*\}$ , the binary observations  $\{Y_i, i \in [n]^*\}$  are independent with distribution

$$\{-1, 1\} \ni y_i \mapsto \int_{\mathbb{R}^d} (1 + \exp(-y_i \langle X_i, z_i \rangle))^{-1} \times \frac{1}{\sqrt{2\pi}^d \sigma^d} \exp(-(2\sigma^2)^{-1} \|z_i - \theta\|^2) dz_i .$$

In words, each example  $\#i$  has an individual regression vector  $Z_i$  in  $\mathbb{R}^d$  and given  $Z_i$ , the success probability  $\mathbb{P}(Y_i = 1 \mid Z_i)$  is  $(1 + \exp(-\langle X_i, Z_i \rangle))^{-1}$ . The regression vectors  $Z_1, \dots, Z_n$  are independent with a Gaussian distribution  $\mathcal{N}(\theta, \sigma^2 \mathbf{I}_d)$ .  $\theta$  is assumed to be unknown and  $\sigma^2$  is known.

The objective is the estimation of  $\theta$  by maximizing the penalized log-likelihood criterion, with a ridge penalty  $\text{pen}(\theta) \stackrel{\text{def}}{=} n\tau \|\theta\|^2$ , where  $\tau > 0$ . By a change of variable, we obtain that the criterion

to be minimized is (see Lemma B.1)

$$F : \theta \mapsto -\frac{1}{n} \sum_{i=1}^n \log \int_{\mathbb{R}} \frac{\exp(x \langle X_i, \theta \rangle / (\sigma^2 \|X_i\|))}{1 + \exp(-y_i \|X_i\| x)} \times \exp(-x^2 / (2\sigma^2)) dx + \|\theta\|_U^2 ,$$

where

$$U \stackrel{\text{def}}{=} \tau \mathbf{I}_d + \frac{1}{2\sigma^2} \frac{1}{n} \sum_{i=1}^n \frac{X_i X_i^\top}{\|X_i\|^2} .$$

The following lemma shows that the minimizers of  $F$  are in a compact set  $\mathcal{K}$  of  $\mathbb{R}^d$  thus implying that the optimization problem can be constrained to  $\mathcal{K}$ . The proof is given in Section B.2.

**Lemma 5.1** *The minimizers of  $F$  are in the set  $\mathcal{K} \stackrel{\text{def}}{=} \{\theta \in \mathbb{R}^d : \|\theta\|^2 \leq (\ln 4)/\tau\}$ .*

To solve this optimization problem, we propose two approaches: a gradient one, solved in the original space  $\theta \in \mathbb{R}^d$  (see Section 2.1); and an EM one, solved in the statistic space (see Section 2.2). The discussions in Section 5.2 and Section 5.3 show that EM is a gradient approach for finding the critical points of  $s \mapsto F(U^{-1}s/2)$ .

### 5.2 A Gradient approach

We are interested in finding a critical point of  $F$  in  $\mathcal{K}$ . Equivalently, we want to solve

$$0 \in \frac{1}{n} \sum_{i=1}^n G_i(\theta) + \partial g(\theta)$$

where  $g$  is the  $\{0, +\infty\}$ -valued indicator function of the set  $\mathcal{K}$  and

$$G_i(\theta) \stackrel{\text{def}}{=} 2U\theta - \frac{X_i}{\sigma^2 \|X_i\|} \int_{\mathbb{R}} z \pi_{\theta,i}(z) dz ;$$

$\pi_{\theta,i}(z)$  is the probability density proportional to

$$\frac{\exp(z \langle X_i, \theta \rangle / (\sigma^2 \|X_i\|) - z^2 / (2\sigma^2))}{1 + \exp(-y_i \|X_i\| z)} . \quad (20)$$

We apply 3P-SPIDER with  $B \stackrel{\text{def}}{=} \mathbf{I}_q$  and  $\mathbf{h}_i(\theta, \mathbf{I}_q) \stackrel{\text{def}}{=} -G_i(\theta)$ ; note that  $\text{prox}_{\gamma g}(\theta) = \text{argmin}_{x \in \mathcal{K}} \|x -$

$\theta\|^2$ .  $\mathbf{h}_i$  is the sum of an explicit term and an integral with no closed form: it will be approximated by a Monte Carlo method, based on a Markov chain Monte Carlo (MCMC) sampler (see Section 5.4 below). Therefore,  $\delta_{t,k,i}$  will be a biased random approximation.

### 5.3 An EM approach

The criterion  $F$  to be minimized is of the form (5) with  $\mathbf{Z} = \mathbb{R}$ ,  $\mu_{l_v}(dz) = dz$  and  $p(Y_i, z; \theta)$  equal to

$$\frac{\exp(z \langle X_i, \theta \rangle / (\sigma^2 \|X_i\|) - z^2 / (2\sigma^2) - \|\theta\|_U^2)}{1 + \exp(-Y_i \|X_i\| z)}.$$

The curved exponential family assumption on the complete data model is satisfied:  $p(Y_i, z; \theta) = H(Y_i, z) \exp(\langle S(Y_i, z), \phi(\theta) \rangle - \psi(\theta))$  with  $\phi(\theta) \stackrel{\text{def}}{=} \theta$ ,  $\psi(\theta) \stackrel{\text{def}}{=} \|\theta\|_U^2$  and

$$S(Y_i, z) \stackrel{\text{def}}{=} z \frac{X_i}{\sigma^2 \|X_i\|}.$$

From Section 2.2, EM in the statistic space is of the form (1)-(2): it solves  $0 \in n^{-1} \sum_{i=1}^n \bar{G}_i(s) + \partial \bar{g}(s)$  where  $\bar{G}_i(s) = B G_i(Bs)$ ,  $B \stackrel{\text{def}}{=} U^{-1}/2$  and  $\bar{g}(s)$  is the  $\{0, +\infty\}$ -valued indicator function of the set  $\{s \in \mathbb{R}^d : \mathsf{T}(s) \stackrel{\text{def}}{=} Bs\}$ ; it uses

$$\bar{\mathbf{h}}_i(s) \stackrel{\text{def}}{=} \frac{X_i}{\sigma^2 \|X_i\|} \int_{\mathbb{R}} z \pi_{Bs,i}(z) dz - s, \quad (21)$$

and the metric induced by  $\mathsf{B}(s) \stackrel{\text{def}}{=} B$ . See Section B.3 for detailed computations. As in the gradient approach,  $\mathbf{h}_i$  requires the expectation of the distribution  $\pi_{\cdot,i}$  (see (20)) which has no closed form. We will run 3P-SPIDER with  $B(s) \leftarrow B$  and a biased random approximation of the  $\bar{\mathbf{h}}_i(s)$ 's (see Section 5.4); note that  $\text{prox}_{\gamma \bar{g}}^B(s) = B^{-1} \text{argmin}_{x \in \mathcal{K}} ((x - Bs)^\top B^{-1}(x - Bs))$ .

### 5.4 The MCMC approximation of $\bar{\mathbf{h}}_i$

We discuss how to design an efficient MCMC sampler for the approximation of

$$\mathcal{I}_i(\theta) \stackrel{\text{def}}{=} \int_{\mathbb{R}} z \pi_{\theta,i}(z) dz, \quad \theta \in \mathbb{R}^d,$$

where  $\pi_{\theta,i}$  is defined, up to a normalizing constant, by (20). By using an integration by parts and by applying (Polson et al, 2013, Theorem 1), we show that a data augmentation scheme is possible to approximate integrals w.r.t.  $\pi_{\theta,i}(z)$ .

**Lemma 5.2** *For any  $i \in [n]^*$  and  $\theta \in \mathbb{R}^d$ , it holds*

$$\begin{aligned} \mathcal{I}_i(\theta) &= \left\langle \frac{X_i}{\|X_i\|}, \theta \right\rangle \\ &+ y_i \|X_i\| \sigma^2 \int_{\mathbb{R}} \int_0^{+\infty} \frac{\bar{\pi}_{\theta,i}(z, \omega)}{1 + \exp(y_i \|X_i\| z)} dz d\omega, \end{aligned}$$

where  $\bar{\pi}_{\theta,i}(z, \omega)$  is a probability density on  $\mathbb{R} \times (0, +\infty)$ . The conditional distribution of  $z$  given  $\omega$  is a Gaussian distribution with parameters

$$\frac{\langle X_i, \theta \rangle / \|X_i\| + y_i \|X_i\| \sigma^2 / 2}{1 + \omega \sigma^2 \|X_i\|^2}, \quad \frac{\sigma^2}{1 + \omega \sigma^2 \|X_i\|^2};$$

the conditional distribution of  $\omega$  given  $z$  is a Polya-Gamma distribution with parameters  $(1, \|X_i\| z)$ .

The proof is given in Section B.4. Therefore, a Monte Carlo approximation of integrals w.r.t.  $\pi_{\theta,i}$  are obtained from a Gibbs sampler targeting the distribution  $\bar{\pi}_{\theta,i}(z, \omega)$ : it produces a sequence of pairs  $\{(Z_r, \Omega_r), r \geq 0\}$  and only the  $Z_r$ 's are retained for the Monte Carlo approximation. For example,  $\bar{\mathbf{h}}_i(s)$  given by (21) can be approximated by

$$\begin{aligned} \bar{\mathbf{h}}_i(s) &\approx -s + \frac{X_i}{\sigma^2 \|X_i\|^2} \langle X_i, Bs \rangle \\ &+ y_i \|X_i\| \frac{1}{m} \sum_{r=1}^m (1 + \exp(y_i \|X_i\| Z_r^{s,i}))^{-1}. \quad (22) \end{aligned}$$

This Gibbs sampler is uniformly ergodic (see (Choi and Hobert, 2013, Proposition 3.1)); consequently, upon noting that  $z \mapsto H_i(z) \stackrel{\text{def}}{=} (1 + \exp(y_i \|X_i\| z))^{-1}$  is bounded by one uniformly in  $i$  and  $z$ , the conditions A5 in Section A are verified with  $U$  equal to the constant function 1 and with a geometric convergence rate  $\rho(r) \stackrel{\text{def}}{=} v^r$  for some  $v \in (0, 1)$  (remember that  $\mathcal{S}$  is a compact set in our application); details are provided in Section B.5.

Therefore, A4 is verified and the rates  $m_{t,k+1}$ ,  $M_{t,k+1}$  and  $\bar{M}_{t,k+1}$  are equal, and equal to the number of points in the Monte Carlo sum (see Proposition A.1).



## 5.5 Numerical illustrations

Let us run 3P-SPIDER for minimizing the criterion  $F$ ; based on previous results comparing variance reduced Expectation Maximization algorithms and variance reduced Gradient algorithms (see e.g. (Chen et al, 2018, section 4)), we restrict our attention to the EM approach. In this numerical application,  $n = 24989$  and  $d = 21$ ; we choose  $\tau = 1$  and  $\sigma^2 = 0.05$ .

**The data set.** The  $n$  pairs  $(y_i, X_i)$  are built from the MNIST data set. The 13007 examples labeled  $y_i = -1$  are the examples labeled 1 or 7 in the MNIST training data set; the 11982 examples labeled  $y_i = 1$  are the examples labeled 3 or 8 in the MNIST training data set. The covariates  $X_i$  are obtained as follows. Let  $X^{\text{im}}$  be the  $784 \times n$  matrix collecting the 784 pixels for each image. The pixels take values in  $[0, 1]$ . Then the rows of  $X^{\text{im}}$  are centered; by a PCA, each image is reduced to a vector in  $\mathbb{R}^{20}$ . This yields  $X^{\text{red}} \in \mathbb{R}^{20 \times n}$ . Finally,  $X^{\text{red}}$  is augmented with a row of ones, yielding  $X \in \mathbb{R}^{21 \times n}$ . The columns of  $X$  are the  $X_i$ 's.

**The algorithms.** We compare four algorithms. EM denotes the SAEM algorithm (Delyon et al (1999)) combined with a proximal step: each iteration processes the full data set so that there is one iteration of EM per epoch:

$$\widehat{S}_{r+1}^{\text{EM}} \stackrel{\text{def}}{=} \text{prox}_{\gamma g}^B(\widehat{S}_r^{\text{EM}} + \frac{\gamma}{n} \sum_{i=1}^n \widehat{h}_i(\widehat{S}_r^{\text{EM}})).$$

**Online EM** is the algorithm given by Cappé and Moulines (2009) combined with a proximal step; each iteration processes  $\mathbf{b}$  examples and below, we will run  $k^{\text{in}} \stackrel{\text{def}}{=} \lceil n/\mathbf{b} \rceil$  iterations per epoch:

$$\widehat{S}_{r+1}^{\text{OEM}} \stackrel{\text{def}}{=} \text{prox}_{\gamma g}^B(\widehat{S}_r^{\text{OEM}} + \frac{\gamma}{\mathbf{b}} \sum_{i \in \mathcal{B}_{r+1}} \widehat{h}_i(\widehat{S}_r^{\text{OEM}})).$$

For EM and Online EM,  $\widehat{h}_i(\widehat{S}_t^\bullet)$  is a Monte Carlo approximation of  $\bar{h}_i(\widehat{S}_t^\bullet)$  computed with  $m^t$  points. 3P-SPIDER is Algorithm 2; we choose  $k_t^{\text{in}} = k^{\text{in}}$  and  $k^{\text{in}} = \lceil n/\mathbf{b} \rceil$  so that one epoch corresponds to the  $k^{\text{in}}$  inner loops; we choose  $\mathbf{b}'_t = n$  so that the initialization of each outer loop is one epoch; the  $\delta_{t,k,i}$  are computed by Monte Carlo sums (see (22)) with  $m^0$  points for  $\delta_{t,0,i}$  and  $m^t$  points for  $\delta_{t,k+1,i}$ ; since  $\widehat{S}_{t,0} = \widehat{S}_{t,-1}$ , we set  $\delta_{t,1,i} = 0$  for

all  $i$ , so that  $S_{t,1} = S_{t,0} = n^{-1} \sum_{i=1}^n \delta_{t,0,i}$ . Finally, 3P-SPIDER and 3P-SPIDER-corr differ as follows: the Monte Carlo approximation  $\delta_{t,k+1,i}$  necessitates a Monte Carlo approximation of  $\bar{h}_i(\widehat{S}_{t,k})$  and one of  $\bar{h}_i(\widehat{S}_{t,k-1})$ . In 3P-SPIDER, the Monte Carlo approximations are based on two independent chains (see (22)) while in 3P-SPIDER-corr the chains are correlated.

All the algorithms are initialized at the null vector  $\widehat{S}_{\text{init}} = 0 \in \mathbb{R}^d$ . The step size is equal to  $\gamma = 0.4$  during the first six epochs and then equal to  $\gamma = 0.1$ . The length of all the paths is 20 epochs. On all the figures except Figure 3, we report a mean value computed over 25 independent runs of each algorithm; the shadowed area is delimited by the minimal and maximal value of the displayed criterion over these runs.

**Analyses.** Most of the comparisons are based on the evolution of

$$\Delta_{t,k+1} \stackrel{\text{def}}{=} \frac{\|\text{prox}_{\gamma g}^B(\widehat{S}_{t,k} + \gamma S_{t,k+1}) - \widehat{S}_{t,k}\|_B^2}{\gamma^2}$$

as a function of the number of epochs; this criterion is an approximation of  $\Delta_{t,k+1}^*$  (see (14)) which can not be computed here since  $\bar{h}$  has no closed form in this application. The criterion  $\Delta_{t,k+1}$  for 3P-SPIDER and 3P-SPIDER-corr, is compared to  $\Delta_r^{\text{EM}}$  defined by

$$\frac{\|\text{prox}_{\gamma g}^B(\widehat{S}_r^{\text{EM}} + \gamma n^{-1} \sum_{i=1}^n \widehat{h}_i(\widehat{S}_r^{\text{EM}})) - \widehat{S}_r^{\text{EM}}\|_B^2}{\gamma^2};$$

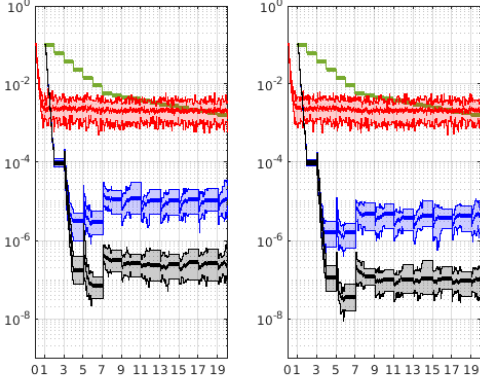
and to  $\Delta_r^{\text{OEM}}$  defined by

$$\frac{\|\text{prox}_{\gamma g}^B(\widehat{S}_r^{\text{OEM}} + \gamma \mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{r+1}} \widehat{h}_i(\widehat{S}_r^{\text{OEM}})) - \widehat{S}_r^{\text{OEM}}\|_B^2}{\gamma^2}.$$

The best algorithm will have the smallest value of  $\Delta_{t,k+1}$ .

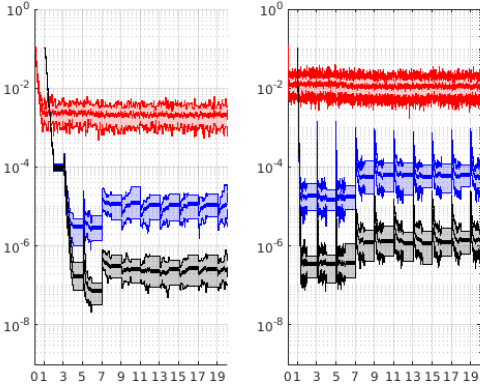
We first study the role of some design parameters of 3P-SPIDER, such as the number of Monte Carlo points when computing  $\delta_{t,0,i}$  (denoted by  $m^0$ ) and  $\delta_{t,k+1,i}$  (denoted by  $m^t$ ) and the balance between  $k^{\text{in}}$  and  $\mathbf{b}$  which satisfy  $k^{\text{in}} \mathbf{b} \approx n$ .

On Figure 1, two strategies are chosen: first,  $m^0 = m^t = 2\lceil \sqrt{n} \rceil$ ; then  $m^0 = m^t = 5\lceil \sqrt{n} \rceil$ ; in all cases,  $k^{\text{in}} = \lceil \sqrt{n}/10 \rceil$  and  $\mathbf{b} = \lceil n/k^{\text{in}} \rceil$ . For comparison, EM and Online EM are also run, with a number of Monte Carlo point equal to  $m^t$  at each iteration.



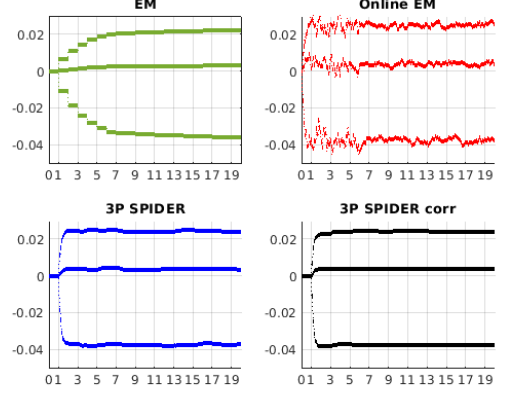
**Fig. 1** Different strategies for the number of Monte Carlo points when approximating  $\bar{h}_i(s)$  - see (22). Evolution of  $\Delta_r^{\text{EM}}$  in green,  $\Delta_r^{\text{OEM}}$  in red,  $\Delta_{t,k+1}$  for 3P-SPIDER in blue and  $\Delta_{t,k+1}$  for 3P-SPIDER corr in black, as a function of the number of epochs. [left]  $m^0 = m^t = 2\lceil\sqrt{n}\rceil$ , [right]  $m^0 = m^t = 5\lceil\sqrt{n}\rceil$ .

On Figure 2, the case when  $k^{\text{in}} = \lceil\sqrt{n}/10\rceil$  is compared to the case  $k^{\text{in}} = \lceil\sqrt{n}/2\rceil$ ; in both cases,  $b = \lceil n/k^{\text{in}}\rceil$  and  $m^0 = m^t = 2\lceil\sqrt{n}\rceil$ .



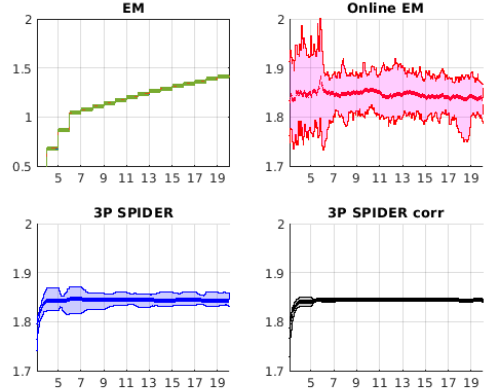
**Fig. 2** Number of inner loops per epoch. Evolution of  $\Delta_r^{\text{OEM}}$  in red,  $\Delta_{t,k+1}$  for 3P-SPIDER in blue and  $\Delta_{t,k+1}$  for 3P-SPIDER corr in black, as a function of the number of epochs. [left]  $k^{\text{in}} = \lceil\sqrt{n}/10\rceil$  and  $b = \lceil n/k^{\text{in}}\rceil$ . [right]  $k^{\text{in}} = \lceil\sqrt{n}/2\rceil$  and  $b = \lceil n/k^{\text{in}}\rceil$ .

Each algorithm returns a sequence of points in the  $s$ -space, from which a sequence of points in the  $\theta$ -space is deduced through the formula  $\theta = \mathsf{T}(s) = Bs \in \mathbb{R}^d$ . On Figure 3, three components of this  $\theta$ -sequence are displayed, versus the number of epochs.



**Fig. 3** Estimation of three parameters. Evolution of the three components of  $\theta$  by EM in green (top, left), OEM in red (top, right), 3P-SPIDER in blue (bottom, left) and 3P-SPIDER corr in black (bottom, right), as a function of the number of epochs.

Finally, we also display on Figure 4 the evolution of the squared norm of the iterates  $\|\hat{S}_{t,k}\|^2$  obtained by 3P-SPIDER and 3P-SPIDER-corr, and  $\|\hat{S}_r^{\text{OM}}\|^2$  and  $\|\hat{S}_r^{\text{OEM}}\|^2$  obtained resp. by EM and Online EM. They are plotted as a function of the epochs.



**Fig. 4** Squared norm of the iterates. Evolution of  $\|\hat{S}_r^{\text{EM}}\|^2$  in green (top, left),  $\|\hat{S}_r^{\text{OEM}}\|^2$  in red (top, right),  $\|\hat{S}_{t,k}\|^2$  for 3P-SPIDER in blue (bottom, left) and  $\|\hat{S}_{t,k}\|^2$  for 3P-SPIDER corr in black (bottom, right), as a function of the number of epochs.

**Conclusions.** EM has a slow convergence rate and even fails to converge before 20 epochs contrary to the other algorithms (see e.g. Figure 4): one update of the iterate per epoch is not enough especially during the first iterations when more

updates even based on part of the data set is a better strategy (see e.g. the behavior of **Online EM**, which contains  $k^{\text{in}}$  updates per epoch).

**Online EM**, **3P-SPIDER** and **3P-SPIDER-corr** process part of the data set at each iteration; compared to **Online EM**, the **3P-SPIDER**'s contain a variance reduction. All the plots illustrate the benefit of this variance reduction, which reduces the variability at convergence.

The choice of  $\gamma$  impacts this variability: see e.g. Figures 1, 2 and Figure 4 where a change occurs at epoch #7 (remember that from epoch  $2\ell$  to  $2\ell+1$ , **Online EM** runs  $k^{\text{in}}$  updates of the iterates while the **3P-SPIDER**'s do not update the iterate since they compute  $S_{t,0}$ ).

**3P-SPIDER-corr** improves on **3P-SPIDER**. The control variate has a larger impact when the correlation is increased, as illustrated by all plots. It decreases the variability introduced by the mini-batches ( $\mathbf{b} < n$ ) and the variability introduced by the Monte Carlo approximation  $\delta_{t,k+1,i}$ .

Given the budget of  $n$  examples processed per outer loops, Figure 2 shows that at convergence, the accuracy is improved by larger mini batch sizes and therefore a smaller number of inner loops. Not surprisingly, a larger number of Monte Carlo points decreases the variability at convergence (see Figure 1).

## 6 Proof of Section 3

### 6.1 Proof of Lemma 3.1

Lemma 6.1 collects the two statements of Lemma 3.1 and a third property.

**Lemma 6.1** *Assume A1.*

1. For any  $\gamma > 0$ ,  $B \in \mathcal{P}_+^q$  and  $s \in \mathbb{R}^q$ , the optimization problem (7) has a unique solution, characterized as the unique point  $\mathbf{p} \in \mathcal{S}$  satisfying  $-\gamma^{-1} B(\mathbf{p} - s) \in \partial g(\mathbf{p})$ .
2. For any  $\gamma > 0$ ,  $B \in \mathcal{P}_+^q$ ,  $s \in \mathcal{S}$  and  $h \in \mathbb{R}^q$ ,

$$s = \text{prox}_{\gamma g}^B(s + \gamma h) \quad \text{iff} \quad Bh \in \partial g(s). \quad (23)$$

3. Let  $\gamma > 0$  and  $B \in \mathcal{P}_+^q$ . The operator  $\text{prox}_{\gamma g}^B$  is firmly nonexpansive; this implies that for any  $s, s' \in \mathbb{R}^q$ ,

$$\|\text{prox}_{\gamma g}^B(s') - \text{prox}_{\gamma g}^B(s)\|_B^2$$

$$\leq \langle \text{prox}_{\gamma g}^B(s') - \text{prox}_{\gamma g}^B(s), s' - s \rangle_B .$$

*Proof* Existence, uniqueness and characterization are established in [Hiriart-Urruty and Lemaréchal \(1996, Chapter XV, Lemma 4.1.1\)](#). The statement (23) follows from the characterization; note that  $\text{prox}_{\gamma g}^B(s) \in \mathcal{S}$  for any  $s \in \mathbb{R}^q$ . The firmly nonexpansive property is a consequence of [Hiriart-Urruty and Lemaréchal \(1996, Chapter XV, Theorem 4.1.4\)](#).  $\square$

## 7 Proof of Section 4

### 7.1 Notations

Define for any  $s \in \mathcal{S}$ ,

$$\bar{h}(s) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \bar{h}_i(s), \quad \bar{h}_{\mathcal{B}} \stackrel{\text{def}}{=} \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}} \bar{h}_i,$$

where  $\mathcal{B}$  is an  $n$ -tuple of elements of  $[n]^*$  (with or without multiplicity) of cardinal  $\mathbf{b}$ .

All the random variables are defined on a probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ . It is endowed with the following filtrations for  $t \geq 0$  and  $k \geq 0$ ,

$$\begin{aligned} \mathcal{F}_{0,k_0^{\text{in}}} &\stackrel{\text{def}}{=} \sigma(\widehat{S}_{\text{init}}), \\ \mathcal{F}_{t,0} &\stackrel{\text{def}}{=} \mathcal{F}_{t-1,k_{t-1}^{\text{in}}} \vee \sigma(\mathcal{B}_{t,0}, \delta_{t,0,i} \text{ for all } i), \\ \mathcal{F}_{t,k+\frac{1}{2}} &\stackrel{\text{def}}{=} \mathcal{F}_{t,k} \vee \sigma(\mathcal{B}_{t,k+1}), \\ \mathcal{F}_{t,k+1} &\stackrel{\text{def}}{=} \mathcal{F}_{t,k+\frac{1}{2}} \vee \sigma(\delta_{t,k+1,i} \text{ for all } i \in \mathcal{B}_{t,k+1}). \end{aligned}$$

For any  $t \in [k^{\text{out}}]^*$ , set

$$\mathcal{E}_t \stackrel{\text{def}}{=} S_{t,0} - \bar{h}(\widehat{S}_{t,0}) = \frac{1}{\mathbf{b}'_t} \sum_{i \in \mathcal{B}_{t,0}} \delta_{t,0,i} - \bar{h}(\widehat{S}_{t,0}).$$

$\mathcal{E}_t$  is the error when replacing the full sum using exact terms  $\bar{h}_i(\widehat{S}_{t,0})$ , with a possibly subsum of size  $\mathbf{b}'_t < n$  using approximations of  $\bar{h}_i(\widehat{S}_{t,0})$ . Remember that

$$\xi_{t,k+1,i} \stackrel{\text{def}}{=} \delta_{t,k+1,i} - \bar{h}_i(\widehat{S}_{t,k}) + \bar{h}_i(\widehat{S}_{t,k-1}),$$

and

$$\begin{aligned} \mu_{t,k+1,i} &\stackrel{\text{def}}{=} \mathbb{E} [\xi_{t,k+1,i} | \mathcal{F}_{t,k+1/2}], \\ \sigma_{t,k+1,i}^2 &\stackrel{\text{def}}{=} \mathbb{E} [\|\xi_{t,k+1,i} - \mu_{t,k+1,i}\|^2 | \mathcal{F}_{t,k+1/2}]. \end{aligned}$$

Finally, set

$$\eta_{t,k+1} \stackrel{\text{def}}{=} \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}_{t,k+1}} \xi_{t,k+1,i}.$$

Throughout the proof, we will use the shorthand notation

$$B_{t,k} \stackrel{\text{def}}{=} \mathbf{B}(\widehat{S}_{t,k}).$$

## 7.2 Preliminary lemmas

**Lemma 7.1** *Let  $\mathcal{B}$  be a batch of  $[n]^*$  of size  $\mathbf{b}$ , sampled at random (with or without replacement).*

1. For any family  $\{f_1, \dots, f_n\}$ ,  
 $\mathbb{E} \left[ \mathbf{b}^{-1} \sum_{i \in \mathcal{B}} f_i \right] = n^{-1} \sum_{i=1}^n f_i$ .
2. For any family  $\{f_1, \dots, f_n\}$ ,

$$\begin{aligned} \mathbb{E} \left[ \left\| \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}} f_i - \frac{1}{n} \sum_{i=1}^n f_i \right\|^2 \right] \\ \leq \frac{1}{\mathbf{b}n} \sum_{i=1}^n \|f_i\|^2 - \frac{1}{n} \sum_{j=1}^n f_j^2. \end{aligned}$$

3. Assume A2. For any  $s, s' \in \mathcal{S}$ , it holds

$$\begin{aligned} \mathbb{E} \left[ \left\| \{ \bar{h}_{\mathcal{B}}(s) - \bar{h}_{\mathcal{B}}(s') \} - \{ \bar{h}(s) - \bar{h}(s') \} \right\|^2 \right] \\ \leq \frac{1}{\mathbf{b}} (L^2 \|s - s'\|^2 - \| \bar{h}(s) - \bar{h}(s') \|^2). \end{aligned}$$

*Proof* The proof is along the same lines as the proof of (Fort et al, 2020, Lemma 4). A detailed proof is provided in Section C.3.  $\square$

**Lemma 7.2** *Assume A4-item 1 and A4-item 3. For any  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , it holds*

$$\begin{aligned} \mathbb{E} \left[ \eta_{t,k+1} | \mathcal{F}_{t,k+1/2} \right] &= \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}_{t,k+1}} \mu_{t,k+1,i}, \\ \mathbb{E} \left[ \eta_{t,k+1} | \mathcal{F}_{t,k} \right] &= \frac{1}{n} \sum_{i=1}^n \mu_{t,k+1,i}, \\ \mathbb{E} \left[ \|\eta_{t,k+1} - \mathbb{E}[\eta_{t,k+1} | \mathcal{F}_{t,k}]\|^2 | \mathcal{F}_{t,k} \right] \\ &\leq \frac{1}{\mathbf{b}} \left( \frac{C_v}{M_{t,k+1}} + \frac{C_{vb}^2}{M_{t,k+1}^2} \right). \end{aligned}$$

*Proof* Let  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ . We have

$$\mathbb{E} \left[ \eta_{t,k+1} | \mathcal{F}_{t,k+1/2} \right] = \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}_{t,k+1}} \mu_{t,k+1,i},$$

since  $\mathcal{B}_{t,k+1} \in \mathcal{F}_{t,k+1/2}$ ; and by Lemma 7.1,

$$\mathbb{E} \left[ \eta_{t,k+1} | \mathcal{F}_{t,k} \right] = \frac{1}{n} \sum_{i=1}^n \mu_{t,k+1,i}.$$

We write

$$\begin{aligned} \eta_{t,k+1} - \mathbb{E} \left[ \eta_{t,k+1} | \mathcal{F}_{t,k} \right] \\ &= \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}_{t,k+1}} \xi_{t,k+1,i} - \frac{1}{n} \sum_{i=1}^n \mu_{t,k+1,i} \\ &= \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}_{t,k+1}} \{ \xi_{t,k+1,i} - \mu_{t,k+1,i} \} \\ &\quad + \frac{1}{\mathbf{b}} \sum_{i \in \mathcal{B}_{t,k+1}} \mu_{t,k+1,i} - \frac{1}{n} \sum_{i=1}^n \mu_{t,k+1,i}. \end{aligned}$$

The RHS is of the form  $U + V$  and we write  $\|U + V\|^2 = \|U\|^2 + \|V\|^2 + 2\langle U, V \rangle$  with  $U \leftarrow \mathbf{b}^{-1} \sum_{i \in \mathcal{B}_{t,k+1}} \{ \xi_{t,k+1,i} - \mu_{t,k+1,i} \}$ . By conditioning and by definition of  $\sigma_{t,k+1,i}^2$ , we have

$$\mathbb{E} \left[ \|U\|^2 | \mathcal{F}_{t,k} \right] = \frac{1}{\mathbf{b}^2} \mathbb{E} \left[ \sum_{i \in \mathcal{B}_{t,k+1}} \sigma_{t,k+1,i}^2 | \mathcal{F}_{t,k} \right].$$

Under A4-item 1, we have by Lemma 7.1

$$\mathbb{E} \left[ \|U\|^2 | \mathcal{F}_{t,k} \right] = \frac{1}{\mathbf{b}n} \sum_{i=1}^n \sigma_{t,k+1,i}^2 \leq \frac{C_v}{\mathbf{b}M_{t,k+1}}.$$

By Lemma 7.1 again, it holds

$$\mathbb{E} \left[ \|V\|^2 | \mathcal{F}_{t,k} \right] \leq \frac{1}{\mathbf{b}n} \sum_{i=1}^n \|\mu_{t,k+1,i} - \frac{1}{n} \sum_{j=1}^n \mu_{t,k+1,j}\|^2,$$

which yields

$$\mathbb{E} \left[ \|V\|^2 | \mathcal{F}_{t,k} \right] \leq \frac{C_{vb}^2}{\mathbf{b}M_{t,k+1}^2}.$$

Finally, upon noting that  $\mathbb{E} \left[ U | \mathcal{F}_{t,k+1/2} \right] = 0$  and  $V \in \mathcal{F}_{t,k+1/2}$ , we have

$$\mathbb{E} \left[ \langle U, V \rangle | \mathcal{F}_{t,k} \right] = \mathbb{E} \left[ \langle \mathbb{E} \left[ U | \mathcal{F}_{t,k+1/2} \right], V \rangle | \mathcal{F}_{t,k} \right] = 0.$$

This concludes the proof.  $\square$

## 7.3 Results on the variables $\mathbf{S}_{t,k}$

Proposition 7.3 studies the bias of the variables  $\mathbf{S}_{t,k+1}$ . It shows that  $\mathbf{S}_{t,k+1}$  is a *biased* approximation of  $\bar{h}(\widehat{S}_{t,k})$ :

$$\mathbb{E} \left[ \mathbf{S}_{t,k+1} | \mathcal{F}_{t,k} \right] \neq \bar{h}(\widehat{S}_{t,k}).$$

When  $k = 0$ , we may have  $\mathbb{E}[\mathbf{S}_{t,1}|\mathcal{F}_{t,0}] = \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,0})$  if  $\delta_{t,0,i} = \mathbf{h}_i(\widehat{\mathbf{S}}_{t,0})$  and  $\mathcal{B}_{t,0} = [n]^*$ . The choice  $\mathcal{B}_{t,0} = [n]^*$  is the strategy proposed in Wang et al (2019) for SpiderBoost; it has an important computational cost but has the advantage to cancel the bias of the variable  $\mathbf{S}$ . at the beginning of each outer loop. Along the inner loops, a (signed) bias appears.

**Proposition 7.3** *For any  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , it holds*

$$\begin{aligned} \mathbb{E}[\mathbf{S}_{t,k+1}|\mathcal{F}_{t,k}] - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k}) \\ = \mathbf{S}_{t,k} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k-1}) + \mathbb{E}[\eta_{t,k+1}|\mathcal{F}_{t,k}], \end{aligned}$$

and

$$\mathbb{E}[\mathbf{S}_{t,k+1} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k})|\mathcal{F}_{t,0}] = \mathcal{E}_t + \sum_{j=1}^{k+1} \mathbb{E}[\eta_{t,j}|\mathcal{F}_{t,0}].$$

*Proof* Let  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ . We write  $\mathbf{S}_{t,k+1} = \mathbf{S}_{t,k} + h_{\mathcal{B}_{t,k+1}}(\widehat{\mathbf{S}}_{t,k}) - h_{\mathcal{B}_{t,k+1}}(\widehat{\mathbf{S}}_{t,k-1}) + \eta_{t,k+1}$ . By Lemma 7.1,

$$\begin{aligned} \mathbb{E}[\mathbf{S}_{t,k+1}|\mathcal{F}_{t,k}] = \mathbf{S}_{t,k} + \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k}) - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k-1}) \\ + \mathbb{E}[\eta_{t,k+1}|\mathcal{F}_{t,k}]. \end{aligned}$$

Since  $\mathcal{F}_{t,0} \subseteq \mathcal{F}_{t,k}$ , we have

$$\begin{aligned} \mathbb{E}[\mathbf{S}_{t,k+1} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k})|\mathcal{F}_{t,0}] = \mathbb{E}[\mathbf{S}_{t,k} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k-1})|\mathcal{F}_{t,0}] \\ + \mathbb{E}[\eta_{t,k+1}|\mathcal{F}_{t,0}]. \end{aligned}$$

Summing from  $j = 0$  to  $j = k$  yields

$$\begin{aligned} \mathbb{E}[\mathbf{S}_{t,k+1} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k})|\mathcal{F}_{t,0}] = \mathbf{S}_{t,0} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,-1}) \\ + \sum_{j=0}^k \mathbb{E}[\eta_{t,j+1}|\mathcal{F}_{t,0}]. \end{aligned}$$

The proof is concluded by using  $\widehat{\mathbf{S}}_{t,0} = \widehat{\mathbf{S}}_{t,-1}$  and the definition of  $\mathcal{E}_t$ ; note that  $\mathcal{E}_t \in \mathcal{F}_{t,0}$ .  $\square$

Proposition 7.4 provides a control of the conditional variance of  $\mathbf{S}_{t,k}$ .

**Proposition 7.4** *Assume A 2, A 4-item 1 and A 4-item 3. For any  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , it holds*

$$\begin{aligned} \mathbb{E} \left[ \left\| \mathbf{S}_{t,k+1} - \mathbb{E}[\mathbf{S}_{t,k+1}|\mathcal{F}_{t,k}] \right\|^2 \middle| \mathcal{F}_{t,k} \right] \\ \leq \frac{L^2}{\mathbf{b}} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}} \right) \|\widehat{\mathbf{S}}_{t,k} - \widehat{\mathbf{S}}_{t,k-1}\|^2 \\ + \frac{C_v}{\mathbf{b} \bar{M}_{t,k+1}} + \frac{C_{vb}^2}{\mathbf{b} \bar{M}_{t,k+1}^2} + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}}. \end{aligned}$$

*Proof* Let  $t \in [k^{\text{out}}]^*$ ,  $k \in [k_t^{\text{in}} - 1]$ . By Lemma 7.1, Proposition 7.3, the definitions of  $\mathbf{S}_{t,k+1}$  and of the filtration  $\mathcal{F}_{t,k}$ ,

$$\begin{aligned} \mathbf{S}_{t,k+1} - \mathbb{E}[\mathbf{S}_{t,k+1}|\mathcal{F}_{t,k}] \\ = \mathbf{S}_{t,k+1} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k}) - \mathbf{S}_{t,k} + \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k-1}) - \mathbb{E}[\eta_{t,k+1}|\mathcal{F}_{t,k}] \\ = \eta_{t,k+1} - \mathbb{E}[\eta_{t,k+1}|\mathcal{F}_{t,k}] \\ + \bar{\mathbf{h}}_{\mathcal{B}_{t,k+1}}(\widehat{\mathbf{S}}_{t,k}) - \bar{\mathbf{h}}_{\mathcal{B}_{t,k+1}}(\widehat{\mathbf{S}}_{t,k-1}) - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k}) + \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k-1}). \end{aligned}$$

The RHS is of the form  $U + V$  with  $U \leftarrow \eta_{t,k+1} - \mathbb{E}[\eta_{t,k+1}|\mathcal{F}_{t,k}]$  and  $V \in \mathcal{F}_{t,k+1/2}$ . Then, we write  $\mathbb{E}[\|U + V\|^2|\mathcal{F}_{t,k}] = \mathbb{E}[\|U\|^2|\mathcal{F}_{t,k}] + \mathbb{E}[\|V\|^2|\mathcal{F}_{t,k}] + 2\mathbb{E}[\langle U, V \rangle|\mathcal{F}_{t,k}]$ .

The term  $\mathbb{E}[\|V\|^2|\mathcal{F}_{t,k}]$  is controlled by Lemma 7.1: an upper bound is  $L^2 \mathbf{b}^{-1} \|\widehat{\mathbf{S}}_{t,k} - \widehat{\mathbf{S}}_{t,k-1}\|^2$ . The term  $\mathbb{E}[\|U\|^2|\mathcal{F}_{t,k}]$  is controlled by Lemma 7.2: an upper bound is  $C_v/(\mathbf{b} M_{t,k+1}) + C_{vb}^2/(\mathbf{b} \bar{M}_{t,k+1}^2)$ . Upon noting that  $V \in \mathcal{F}_{t,k+1/2}$ , and using Lemma 7.2 and Lemma 7.1, the scalar product is upper bounded by

$$\begin{aligned} 2\mathbb{E}[\|V\| \left\| \mathbb{E}[U|\mathcal{F}_{t,k+1/2}] \right\| \middle| \mathcal{F}_{t,k}] \\ \leq 2 \frac{C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}} \left\{ \mathbb{E}[\|V\|^2|\mathcal{F}_{t,k}] \right\}^{1/2} \\ \leq 2 \frac{C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}} \left( 1 + \mathbb{E}[\|V\|^2|\mathcal{F}_{t,k}] \right), \end{aligned}$$

where we used that  $\mathbf{a} \leq 1 + \mathbf{a}^2$ .  $\square$

Proposition 7.5 establishes an upper bound on the conditional expectation of the quadratic error  $\|\mathbf{S}_{t,k+1} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k})\|^2$ .

**Proposition 7.5** *Assume A 2 and A 4. For any  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , it holds*

$$\begin{aligned} \mathbb{E}[\|\mathbf{S}_{t,k+1} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k})\|^2|\mathcal{F}_{t,k}] \\ \leq \left( 1 + \frac{2C_b}{m_{t,k+1}} \right) \|\mathbf{S}_{t,k} - \bar{\mathbf{h}}(\widehat{\mathbf{S}}_{t,k-1})\|^2 \\ + \frac{L^2}{\mathbf{b}} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}} \right) \|\widehat{\mathbf{S}}_{t,k} - \widehat{\mathbf{S}}_{t,k-1}\|^2 \\ + \mathcal{U}_{t,k+1}, \end{aligned}$$

where

$$\mathcal{U}_{t,k} \stackrel{\text{def}}{=} \frac{2C_b}{m_{t,k}} + \frac{C_b^2}{m_{t,k}^2} + \frac{C_v}{\mathbf{b} M_{t,k}} + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k}} + \frac{C_{vb}^2}{\mathbf{b} \bar{M}_{t,k}^2}.$$

*Proof* Let  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ . By definition of the conditional expectation, we have for any r.v.  $U, V$  and any  $\sigma$ -field  $\mathcal{F}$  such that  $V \in \mathcal{F}$ :

$$\begin{aligned} & \mathbb{E} \left[ \|U - V\|^2 | \mathcal{F} \right] \\ &= \mathbb{E} \left[ \|U - \mathbb{E}[U | \mathcal{F}]\|^2 | \mathcal{F} \right] + \|\mathbb{E}[U | \mathcal{F}] - V\|^2. \end{aligned}$$

We apply this equality with  $U \leftarrow \widehat{S}_{t,k+1}$ ,  $V \leftarrow \bar{h}(\widehat{S}_{t,k})$  and  $\mathcal{F} \leftarrow \mathcal{F}_{t,k}$ . Proposition 7.4 controls the first term. For the second one, by Proposition 7.3, Lemma 7.2 and A4-item 2 we have

$$\begin{aligned} & \|\mathbb{E}[\mathbf{S}_{t,k+1} | \mathcal{F}_{t,k}] - \bar{h}(\widehat{S}_{t,k})\|^2 \\ &= \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1}) + \mathbb{E}[\eta_{t,k+1} | \mathcal{F}_{t,k}]\|^2 \\ &\leq \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1})\|^2 + \frac{C_b^2}{m_{t,k+1}^2} \\ &+ 2 \frac{C_b}{m_{t,k+1}} \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1})\|. \end{aligned}$$

We conclude by using  $\|a\| \leq 1 + \|a\|^2$  with  $a \leftarrow \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1})\|$ .  $\square$

**Corollary 7.6** (of Proposition 7.5) *Assume also A3-item 3. For  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , define  $\mathcal{D}_{t,k+1}$  by*

$$\|\widehat{S}_{t,k+1} - \text{prox}_{\gamma_{t,k+1}g}^{B_{t,k}}(\widehat{S}_{t,k} + \gamma_{t,k+1}\bar{h}(\widehat{S}_{t,k}))\|_{B_{t,k}}^2,$$

and  $\mathcal{D}_{t,0} \stackrel{\text{def}}{=} 0$ . For  $t \in [k^{\text{out}}]^*$  and  $k \in [k_t^{\text{in}} - 1]$ , it holds

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{S}_{t,k+1} - \bar{h}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,k} \right] \\ &\leq \left( 1 + \frac{2C_b}{m_{t,k+1}} \right) \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1})\|^2 \\ &+ \gamma_{t,k}^2 \frac{2}{v_{\min}} \frac{L^2}{\mathbf{b}} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}} \right) \Delta_{t,k}^* \\ &+ \frac{2}{v_{\min}} \frac{L^2}{\mathbf{b}} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,k+1}} \right) \mathcal{D}_{t,k} \\ &+ \mathcal{U}_{t,k+1}. \end{aligned}$$

By convention,  $\Delta_{t,0}^* \stackrel{\text{def}}{=} 0$ .

*Proof* The proof consists in an upper bound for  $\|\widehat{S}_{t,k} - \widehat{S}_{t,k-1}\|^2$ . Let  $s \in \mathcal{S}$ ,  $H, h \in \mathbb{R}^q$ ,  $\gamma > 0$  and  $B$  be a  $q \times q$  positive definite matrix. For any  $\beta > 0$ , it holds

$$\begin{aligned} & \|\text{prox}_{\gamma g}^B(s + \gamma H) - s\|_B^2 \leq (1 + \frac{1}{\beta}) \|\text{prox}_{\gamma g}^B(s + \gamma h) - s\|_B^2 \\ &+ (1 + \beta) \|\text{prox}_{\gamma g}^B(s + \gamma H) - \text{prox}_{\gamma g}^B(s + \gamma h)\|_B^2. \end{aligned}$$

We apply these inequalities with  $\gamma \leftarrow \gamma_{t,k}$ ,  $B \leftarrow B_{t,k-1}$ ,  $s \leftarrow \widehat{S}_{t,k-1}$ ,  $H \leftarrow \mathbf{S}_{t,k}$  and  $h \leftarrow \bar{h}(\widehat{S}_{t,k-1})$ . Then, for any  $k > 0$ ,

$$\begin{aligned} & \|\widehat{S}_{t,k} - \widehat{S}_{t,k-1}\|_{B_{t,k-1}}^2 \\ &\leq (1 + \beta^{-1}) \gamma_{t,k}^2 \Delta_{t,k}^* + (1 + \beta) \mathcal{D}_{t,k}. \end{aligned} \quad (24)$$

We choose  $\beta = 1$  and conclude by A3-item 3:  $\|\cdot\|^2 \leq v_{\min}^{-1} \|\cdot\|_{B_{t,k-1}}^2$ .

When  $k = 0$ ,  $\|\widehat{S}_{t,k} - \widehat{S}_{t,k-1}\|_{B_{t,k-1}}^2 = 0$  since by definition,  $\widehat{S}_{t,0} = \widehat{S}_{t,-1}$ . Therefore, (24) remains valid since  $\mathcal{D}_{t,0} = 0$  and  $\Delta_{t,0}^* = 0$  by convention. This concludes the proof.  $\square$

**Corollary 7.7** (of Corollary 7.6) *Let  $\{\rho_{t,k}, t \geq 1, k \geq 0\}$  be a positive sequence satisfying*

$$\rho_{t,k+1} \left( 1 + \frac{2C_b}{m_{t,k+1}} \right) \leq \rho_{t,k}. \quad (25)$$

For any  $t \in [k^{\text{out}}]^*$ ,  $k \in [k_t^{\text{in}} - 1]$ , it holds

$$\begin{aligned} & \rho_{t,k+1} \mathbb{E} \left[ \|\mathbf{S}_{t,k+1} - \bar{h}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,0} \right] \\ &\leq \rho_{t,0} \|\mathcal{E}_t\|^2 + \sum_{\ell=1}^{k+1} \rho_{t,\ell} \mathcal{U}_{t,\ell} + \frac{2}{v_{\min}} \frac{L^2}{\mathbf{b}} \cdots \\ &\times \left\{ \sum_{\ell=1}^k \gamma_{t,\ell}^2 \rho_{t,\ell+1} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,\ell+1}} \right) \mathbb{E}[\Delta_{t,\ell}^* | \mathcal{F}_{t,0}] \right. \\ &\left. + \sum_{\ell=1}^k \rho_{t,\ell+1} \left( 1 + \frac{2C_{vb}}{\sqrt{\mathbf{b}} \bar{M}_{t,\ell+1}} \right) \mathbb{E}[\mathcal{D}_{t,\ell} | \mathcal{F}_{t,0}] \right\}. \end{aligned}$$

*Proof* In Corollary 7.6, the claim is of the form

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{S}_{t,k+1} - \bar{h}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,k} \right] \\ &\leq \left( 1 + \frac{2C_b}{m_{t,k+1}} \right) \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1})\|^2 + A_k. \end{aligned}$$

This yields, by using the condition (25),

$$\begin{aligned} & \rho_{t,k+1} \mathbb{E} \left[ \|\mathbf{S}_{t,k+1} - \bar{h}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,k} \right] \\ &\leq \rho_{t,k} \|\mathbf{S}_{t,k} - \bar{h}(\widehat{S}_{t,k-1})\|^2 + \rho_{t,k+1} A_k. \end{aligned}$$

Using  $\mathbb{E}[U | \mathcal{F}_{t,0}] = \mathbb{E}[\mathbb{E}[U | \mathcal{F}_{t,k}] | \mathcal{F}_{t,0}]$  and summing from  $\ell = 0$  to  $\ell = k$  yields

$$\begin{aligned} & \rho_{t,k+1} \mathbb{E} \left[ \|\mathbf{S}_{t,k+1} - \bar{h}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,0} \right] \\ &\leq \sum_{\ell=0}^k \rho_{t,\ell+1} \mathbb{E} [A_\ell | \mathcal{F}_{t,0}] + \rho_{t,0} \|\mathbf{S}_{t,0} - \bar{h}(\widehat{S}_{t,-1})\|^2; \end{aligned}$$

we then conclude by using the equality  $\widehat{S}_{t,-1} = \widehat{S}_{t,0}$  and the definition of  $\mathcal{E}_t$ . Note also that  $\Delta_{t,0}^* = 0$  and  $\mathcal{D}_{t,0} = 0$ .  $\square$

## 7.4 Lyapunov inequalities for $W$ , $g$ and $W + g$

Lemma 7.8, while being classical in smooth optimization, is provided for a self-content purpose.

**Lemma 7.8** *Assume A3. For any  $s, s' \in \mathcal{S}$  and  $\gamma > 0$ ,*

$$W(s') \leq W(s) - \langle \bar{h}(s), s' - s \rangle_{B(s)} + \frac{L_{\dot{W}}}{2} \|s' - s\|^2.$$

*Proof*  $\mathcal{S}$  is convex since it is the domain of a convex function. By A3,  $W$  is continuously differentiable on  $\mathcal{S}$  with  $L_{\dot{W}}$ -Lipschitz gradient. Then for any  $s, s' \in \mathcal{S}$ ,

$$W(s') - W(s) \leq \langle \nabla W(s), s' - s \rangle + \frac{L_{\dot{W}}}{2} \|s' - s\|^2.$$

We use that  $\nabla W(s) = -B(s)\bar{h}(s)$ , so that

$$\langle \nabla W(s), s' - s \rangle = -\langle \bar{h}(s), s' - s \rangle_{B(s)}.$$

□

**Lemma 7.9** *Assume A1. Let  $B$  be a  $q \times q$  positive definite matrix. For any  $s \in \mathcal{S}$ ,  $\gamma > 0$ ,  $H, h \in \mathbb{R}^q$  and  $\beta > 0$ ,*

$$\begin{aligned} g\left(\text{prox}_{\gamma g}^B(s + \gamma H)\right) &\leq g(s) \\ &- \frac{1}{\gamma} \left(1 - \frac{\beta}{4}\right) \|\text{prox}_{\gamma g}^B(s + \gamma h) - s\|_B^2 \\ &- \frac{1}{\gamma} \left(1 - \frac{1}{\beta}\right) \|\text{prox}_{\gamma g}^B(s + \gamma h) - \text{prox}_{\gamma g}^B(s + \gamma H)\|_B^2 \\ &- \left\langle h, s - \text{prox}_{\gamma g}^B(s + \gamma h) \right\rangle_B \\ &+ \left\langle H, \text{prox}_{\gamma g}^B(s + \gamma H) - \text{prox}_{\gamma g}^B(s + \gamma h) \right\rangle_B. \end{aligned}$$

*Proof* In this proof, we use the shorthand notation  $\mathfrak{p}_H \stackrel{\text{def}}{=} \text{prox}_{\gamma g}^B(s + \gamma H)$  and  $\mathfrak{p}_h \stackrel{\text{def}}{=} \text{prox}_{\gamma g}^B(s + \gamma h)$ . By Lemma 3.1 and the definition of the subdifferential at a point, it holds

$$\begin{aligned} g(\mathfrak{p}_h) &\geq g(\mathfrak{p}_H) - \gamma^{-1} \langle \mathfrak{p}_H - s - \gamma H, \mathfrak{p}_h - \mathfrak{p}_H \rangle_B \\ g(s) &\geq g(\mathfrak{p}_h) - \gamma^{-1} \langle \mathfrak{p}_h - s - \gamma h, s - \mathfrak{p}_h \rangle_B. \end{aligned}$$

This yields

$$\begin{aligned} g(\mathfrak{p}_H) &\leq g(s) - \gamma^{-1} \|\mathfrak{p}_h - s\|_B^2 - \langle h, s - \mathfrak{p}_h \rangle_B \\ &- \langle H, \mathfrak{p}_h - \mathfrak{p}_H \rangle_B + \gamma^{-1} \langle \mathfrak{p}_H - s, \mathfrak{p}_h - \mathfrak{p}_H \rangle_B. \end{aligned}$$

For the last term, we write for any  $\beta > 0$ ,

$$\begin{aligned} &\gamma^{-1} \langle \mathfrak{p}_H - s, \mathfrak{p}_h - \mathfrak{p}_H \rangle_B + \gamma^{-1} \|\mathfrak{p}_h - \mathfrak{p}_H\|_B^2 \\ &= \gamma^{-1} \langle \mathfrak{p}_h - s, \mathfrak{p}_h - \mathfrak{p}_H \rangle_B \end{aligned}$$

$$\begin{aligned} &\leq 2 \left\langle (\mathfrak{p}_h - s) \frac{\sqrt{\beta}}{2\sqrt{\gamma}}, (\mathfrak{p}_h - \mathfrak{p}_H) \frac{1}{\sqrt{\beta\gamma}} \right\rangle_B \\ &\leq \frac{\beta}{4\gamma} \|\mathfrak{p}_h - s\|_B^2 + \frac{1}{\beta\gamma} \|\mathfrak{p}_h - \mathfrak{p}_H\|_B^2. \end{aligned}$$

This concludes the proof. □

**Proposition 7.10** *Assume A1, A2 and A3. For any  $t \in [k^{\text{out}}]^*$ ,  $k \in [k^{\text{in}} - 1]$  and  $\beta > 0$ ,*

$$\begin{aligned} &\mathbb{E} \left[ W(\hat{S}_{t,k+1}) + g(\hat{S}_{t,k+1}) | \mathcal{F}_{t,0} \right] \\ &\leq \mathbb{E} \left[ W(\hat{S}_{t,k}) + g(\hat{S}_{t,k}) | \mathcal{F}_{t,0} \right] \\ &- \gamma_{t,k+1} \left( 1 - \frac{\beta}{4} - \frac{L_{\dot{W}} \gamma_{t,k+1}}{v_{\min}} \right) \mathbb{E} \left[ \Delta_{t,k+1}^* | \mathcal{F}_{t,0} \right] \\ &- \frac{1}{\gamma_{t,k+1}} \left( 1 - \frac{1}{\beta} - \frac{L_{\dot{W}}}{v_{\min}} \gamma_{t,k+1} \right) \mathbb{E} \left[ \mathcal{D}_{t,k+1} | \mathcal{F}_{t,0} \right] \\ &+ \gamma_{t,k+1} \mathbb{E} \left[ \|\mathfrak{S}_{t,k+1} - \bar{h}(\hat{S}_{t,k})\|_{B_{t,k}}^2 | \mathcal{F}_{t,0} \right], \end{aligned}$$

where  $\mathcal{D}_{t,k+1}$  is defined in Corollary 7.6.

*Proof* Let  $\gamma > 0$ ,  $s \in \mathcal{S}$  and  $H \in \mathbb{R}^q$ . Apply Lemma 7.8 with  $s' \leftarrow \text{prox}_{\gamma g}^B(s + \gamma H) \in \mathcal{S}$ ; and Lemma 7.9 with  $h \leftarrow \bar{h}(s)$ . This yields for any  $\beta > 0$ ,

$$\begin{aligned} &W(\text{prox}_{\gamma g}^B(s + \gamma H)) + g(\text{prox}_{\gamma g}^B(s + \gamma H)) \\ &\leq W(s) + g(s) - \frac{1}{\gamma} \left(1 - \frac{\beta}{4}\right) \|\text{prox}_{\gamma g}^B(s + \gamma \bar{h}(s)) - s\|_B^2 \\ &- \frac{1}{\gamma} \left(1 - \frac{1}{\beta}\right) \|\text{prox}_{\gamma g}^B(s + \gamma H) - \text{prox}_{\gamma g}^B(s + \gamma \bar{h}(s))\|_B^2 \\ &- \left\langle \bar{h}(s) - H, \text{prox}_{\gamma g}^B(s + \gamma H) - \text{prox}_{\gamma g}^B(s + \gamma \bar{h}(s)) \right\rangle_B \\ &+ \frac{L_{\dot{W}}}{2} \|\text{prox}_{\gamma g}^B(s + \gamma H) - s\|^2. \end{aligned}$$

Since  $\text{prox}_{\gamma g}^B$  is firmly nonexpansive (see Lemma 6.1), the scalar product is upper bounded by  $\gamma \|H - \bar{h}(s)\|_B^2$ . By A3-item 3, we write

$$\|\text{prox}_{\gamma g}^B(s + \gamma H) - s\|^2 \leq \frac{1}{v_{\min}} \|\text{prox}_{\gamma g}^B(s + \gamma H) - s\|_B^2;$$

then we use  $\|a + b\|_B^2 \leq 2\|a\|_B^2 + 2\|b\|_B^2$  with  $a \leftarrow \text{prox}_{\gamma g}^B(s + \gamma H) - \text{prox}_{\gamma g}^B(s + \gamma \bar{h}(s))$ . This yields

$$\begin{aligned} &\frac{L_{\dot{W}}}{2} \|\text{prox}_{\gamma g}^B(s + \gamma H) - s\|^2 \\ &\leq \frac{L_{\dot{W}}}{v_{\min}} \|\text{prox}_{\gamma g}^B(s + \gamma H) - \text{prox}_{\gamma g}^B(s + \gamma \bar{h}(s))\|_B^2 \\ &+ \frac{L_{\dot{W}}}{v_{\min}} \|\text{prox}_{\gamma g}^B(s + \gamma \bar{h}(s)) - s\|_B^2. \end{aligned}$$

We apply these inequalities with  $s \leftarrow \hat{S}_{t,k}$ ,  $\gamma \leftarrow \gamma_{t,k+1}$ ,  $H \leftarrow \mathfrak{S}_{t,k+1}$ ,  $s' \leftarrow \hat{S}_{t,k+1}$  and  $B \leftarrow B_{t,k}$ . Note that  $\text{prox}_{\gamma g}^B(s + \gamma H) = \hat{S}_{t,k+1}$ . The proof is concluded. □

## 7.5 Proof of Theorem 4.1

Let  $t \in [k^{\text{out}}]^*$ . Let  $\mu \in (0, 1)$ . Throughout the proof, set

$$A_{t,k+1} \stackrel{\text{def}}{=} \left( 1 + \frac{2C_{vb}}{\sqrt{b} \bar{M}_{t,k+1}} \right).$$

From Corollary 7.7 applied with  $\rho_{t,k+1} \leftarrow \gamma_{t,k+1}$  and Proposition 7.10 applied with  $\beta \leftarrow 4\mu$ , it holds for any  $k \in [k_t^{\text{in}} - 1]$ ,

$$\begin{aligned} & \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,k+1}) + g(\widehat{S}_{t,k+1}) | \mathcal{F}_{t,0} \right] \\ & \leq \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,k}) + g(\widehat{S}_{t,k}) | \mathcal{F}_{t,0} \right] \\ & \quad - \gamma_{t,k+1} \left( 1 - \mu - \frac{L_{\dot{W}}}{v_{\min}} \gamma_{t,k+1} \right) \mathbb{E} [\Delta_{t,k+1}^* | \mathcal{F}_{t,0}] \\ & \quad - \frac{1}{\gamma_{t,k+1}} \left( 1 - \frac{1}{4\mu} - \frac{L_{\dot{W}}}{v_{\min}} \gamma_{t,k+1} \right) \mathbb{E} [\mathcal{D}_{t,k+1} | \mathcal{F}_{t,0}] \\ & \quad + \gamma_{t,0} v_{\max} \|\mathcal{E}_t\|^2 + v_{\max} \sum_{\ell=1}^{k+1} \gamma_{t,\ell} \mathcal{U}_{t,\ell} \\ & \quad + \frac{2v_{\max}}{v_{\min}} \frac{L^2}{b} \sum_{\ell=1}^k \gamma_{t,\ell}^3 A_{t,\ell+1} \mathbb{E} [\Delta_{t,\ell}^* | \mathcal{F}_{t,0}] \\ & \quad + \frac{2v_{\max}}{v_{\min}} \frac{L^2}{b} \sum_{\ell=1}^k \gamma_{t,\ell+1} A_{t,\ell+1} \mathbb{E} [\mathcal{D}_{t,\ell} | \mathcal{F}_{t,0}]. \end{aligned}$$

Above, we used that  $\gamma_{t,k+1} \leq \gamma_{t,\ell}$  for any  $\ell \in [k+1]$ . We now sum from  $k=0$  to  $k=k_t^{\text{in}}-1$ . This yields,

$$\begin{aligned} & \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,k_t^{\text{in}}}) + g(\widehat{S}_{t,k_t^{\text{in}}}) | \mathcal{F}_{t,0} \right] \\ & \leq \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,0}) + g(\widehat{S}_{t,0}) | \mathcal{F}_{t,0} \right] \\ & \quad - \sum_{k=1}^{k_t^{\text{in}}} \gamma_{t,k} \left( 1 - \mu - \frac{L_{\dot{W}}}{v_{\min}} \gamma_{t,k} \right) \mathbb{E} [\Delta_{t,k}^* | \mathcal{F}_{t,0}] \\ & \quad - \sum_{k=1}^{k_t^{\text{in}}} \frac{1}{\gamma_{t,k}} \left( 1 - \frac{1}{4\mu} - \frac{L_{\dot{W}}}{v_{\min}} \gamma_{t,k} \right) \mathbb{E} [\mathcal{D}_{t,k} | \mathcal{F}_{t,0}] \\ & \quad + \gamma_{t,0} v_{\max} k_t^{\text{in}} \|\mathcal{E}_t\|^2 \\ & \quad + v_{\max} \sum_{\ell=1}^{k_t^{\text{in}}} (k_t^{\text{in}} - \ell + 1) \gamma_{t,\ell} \mathcal{U}_{t,\ell} \\ & \quad + \frac{2v_{\max}}{v_{\min}} k_t^{\text{in}} \sum_{k=1}^{k_t^{\text{in}}-1} \gamma_{t,k}^3 A_{t,k+1} \mathbb{E} [\Delta_{t,k}^* | \mathcal{F}_{t,0}] \end{aligned}$$

$$+ \frac{2v_{\max}}{v_{\min}} k_t^{\text{in}} \sum_{k=1}^{k_t^{\text{in}}-1} \gamma_{t,k+1} A_{t,k+1} \mathbb{E} [\mathcal{D}_{t,k} | \mathcal{F}_{t,0}].$$

Observe that the coefficient in front of  $\mathbb{E} [\Delta_{t,k}^* | \mathcal{F}_{t,0}]$  is  $\gamma_{t,k}(1 - \mu - \Lambda_{t,k+1})$ ; and the term in front of  $\mathbb{E} [\mathcal{D}_{t,k} | \mathcal{F}_{t,0}]$  is  $\gamma_{t,k}^{-1}(1 - 1/(4\mu) - \Lambda_{t,k+1})$ . By symmetry, we choose  $\mu = 1/2$  so that  $\mu = 1/(4\mu)$ . This yields

$$\begin{aligned} & \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,k_t^{\text{in}}}) + g(\widehat{S}_{t,k_t^{\text{in}}}) | \mathcal{F}_{t,0} \right] \\ & \leq \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,0}) + g(\widehat{S}_{t,0}) | \mathcal{F}_{t,0} \right] \\ & \quad - \sum_{k=1}^{k_t^{\text{in}}} \gamma_{t,k} \left( \frac{1}{2} - \Lambda_{t,k+1} \right) \mathbb{E} [\Delta_{t,k}^* | \mathcal{F}_{t,0}] \\ & \quad - \sum_{k=1}^{k_t^{\text{in}}} \frac{1}{\gamma_{t,k}} \left( \frac{1}{2} - \Lambda_{t,k+1} \right) \mathbb{E} [\mathcal{D}_{t,k} | \mathcal{F}_{t,0}] \\ & \quad + \gamma_{t,0} v_{\max} k_t^{\text{in}} \|\mathcal{E}_t\|^2 \\ & \quad + v_{\max} \sum_{\ell=1}^{k_t^{\text{in}}} (k_t^{\text{in}} - \ell + 1) \gamma_{t,\ell} \mathcal{U}_{t,\ell}. \end{aligned}$$

We now sum for  $t=1$  to  $t=k^{\text{out}}$  and compute the expectation. This yields, by using that  $\widehat{S}_{t+1,0} = \widehat{S}_{t,k_t^{\text{in}}}$ ,

$$\begin{aligned} & \sum_{t=1}^{k^{\text{out}}} \sum_{k=1}^{k_t^{\text{in}}} \gamma_{t,k} \left( \frac{1}{2} - \Lambda_{t,k+1} \right) \mathbb{E} [\Delta_{t,k}^*] \\ & \quad + \sum_{t=1}^{k^{\text{out}}} \sum_{k=1}^{k_t^{\text{in}}} \frac{1}{\gamma_{t,k}} \left( \frac{1}{2} - \Lambda_{t,k+1} \right) \mathbb{E} [\mathcal{D}_{t,k}] \\ & \leq \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{t,0}) + g(\widehat{S}_{t,0}) \right] \\ & \quad - \mathbb{E} \left[ \mathbb{W}(\widehat{S}_{k^{\text{out}},k_{k^{\text{out}}}^{\text{in}}}) + g(\widehat{S}_{k^{\text{out}},k_{k^{\text{out}}}^{\text{in}}}) \right] \\ & \quad + v_{\max} \sum_{t=1}^{k^{\text{out}}} \gamma_{t,0} k_t^{\text{in}} \mathbb{E} [\|\mathcal{E}_t\|^2] \\ & \quad + v_{\max} \sum_{t=1}^{k^{\text{out}}} \sum_{\ell=1}^{k_t^{\text{in}}} (k_t^{\text{in}} - \ell + 1) \gamma_{t,\ell} \mathcal{U}_{t,\ell}. \end{aligned}$$

The proof is concluded upon noting that  $\mathbb{E} \left[ \mathbb{W}(\widehat{S}_{k^{\text{out}},k_{k^{\text{out}}}^{\text{in}}}) + g(\widehat{S}_{k^{\text{out}},k_{k^{\text{out}}}^{\text{in}}}) \right] \geq \min_{\mathcal{S}} (W + g)$ .



## 7.6 Proof of Corollary 4.3

Since  $\mathcal{U}_{t,k} = 0$ , we have  $C_b = C_{vb} = 0$ . In addition,  $k_t^{\text{in}} = k^{\text{in}}$  for any  $t$ . Therefore, we can consider a constant stepsize sequence  $\gamma_{t,k} = \gamma_*$  where  $\gamma_*$  satisfies (see (15) and (16))

$$\gamma_* \frac{L_{\dot{W}}}{v_{\min}} + \gamma_*^2 \frac{2v_{\max}}{v_{\min}} L^2 \frac{k^{\text{in}}}{\mathbf{b}} \in (0, 1/2) .$$

This condition is satisfied by choosing

$$\begin{aligned} \frac{k^{\text{in}}}{\mathbf{b}} &\stackrel{\text{def}}{=} \frac{1}{v_{\min} v_{\max}} \frac{L_{\dot{W}}^2}{L^2} , \\ \gamma_* &\stackrel{\text{def}}{=} \frac{1}{4v_{\max}} \frac{L_{\dot{W}}}{L^2} \frac{\mathbf{b}}{k^{\text{in}}} = \frac{v_{\min}}{4L_{\dot{W}}} . \end{aligned}$$

Such a choice implies that  $\inf_{t,k} (1/2 - \Lambda_{t,k}) = 1/2 - 3/8 = 1/8$ . Since  $\mathcal{E}_t = \mathcal{U}_{t,k} = 0$ , we obtain from Corollary 4.2 that

$$\begin{aligned} &\mathbb{E} [\Delta_{\tau,K}^* + \mathcal{D}_{\tau,K}^*] \\ &\leq \frac{32 L_{\dot{W}}}{v_{\min}} \left( \frac{\mathbb{E} [\mathbb{W}(\widehat{S}_{1,0}) + g(\widehat{S}_{1,0})] - \min_S (\mathbb{W} + g)}{k^{\text{out}} k^{\text{in}}} \right) . \end{aligned}$$

The  $\epsilon$ -approximate stationary condition is satisfied by choosing  $k^{\text{out}} k^{\text{in}} = O(L_{\dot{W}}/(v_{\min} \epsilon))$ . The number of calls to the proximal operator is  $k^{\text{out}} k^{\text{in}}$  so that  $\mathcal{K}_{\text{prox}} = O(L_{\dot{W}}/(v_{\min} \epsilon))$ . Finally, we have  $\mathbf{b}'_t = n$  so that the number of calls to one of the  $\bar{h}_i$ 's is  $k^{\text{out}} n + 2k^{\text{out}} k^{\text{in}} \mathbf{b}$ . We can choose  $\mathbf{b} = O(\sqrt{n} \sqrt{v_{\min} v_{\max}} L / L_{\dot{W}})$ . This yields  $k^{\text{out}} = O(L \sqrt{v_{\max}} / (\sqrt{v_{\min} \epsilon} \sqrt{n}))$ , and  $\mathcal{K}_{\bar{h}} = O(\sqrt{v_{\max}} L \sqrt{n} / (\epsilon \sqrt{v_{\min}}))$ .

## 7.7 Cost of the approximation on the $\bar{h}_i$ 's

Following the rates obtained in Corollary 4.3, let us set  $k_t^{\text{in}} = O(\sqrt{n})$ ,  $\mathbf{b} = O(\sqrt{n})$  and  $k^{\text{out}} = O(1/(\sqrt{n}\epsilon))$  and let us show that we can define random approximations  $\delta_{t,0,i}$  and  $\delta_{t,k+1,i}$  such that the  $\epsilon$ -approximate stationarity condition is satisfied.

**On the term  $\mathbb{E} [\|\mathcal{E}_\tau\|^2]$ .** We write  $\mathbb{E} [\|\mathcal{E}_\tau\|^2] = (k^{\text{out}})^{-1} \sum_{t=1}^{k^{\text{out}}} \mathbb{E} [\|\mathcal{E}_t\|^2]$  and

$$\frac{1}{k^{\text{out}}} \sum_{t=1}^{k^{\text{out}}} \frac{\epsilon^{1-a'}}{\sqrt{n}^{a'} t^{a'}} = \epsilon O(1) .$$

Let us compute the associated Monte Carlo complexity in the case  $\mathcal{E}_t = n^{-1} \sum_{i=1}^n \{\delta_{t,0,i} - \bar{h}_i(\widehat{S}_{t,0})\}$  and  $\delta_{t,0,i}$  is equal to a Monte Carlo sum with  $m_{t,0}$  i.i.d. samples. Then  $\mathbb{E} [\|\mathcal{E}_t\|^2] = n^{-1} m_{t,0}^{-1} O(1)$ . It is equal to  $O(\epsilon^{1-a'} / (\sqrt{nt})^{a'})$  when  $m_{t,0} = O(n^{a'/2-1} t^{a'} / \epsilon^{1-a'})$ . Therefore, the Monte Carlo cost is

$$\sum_{t=1}^{k^{\text{out}}} n m_{t,0} = O\left(\frac{1}{\sqrt{n}\epsilon^2}\right) .$$

**On the term  $\mathbb{E} [(k^{\text{in}} - K + 1) \mathcal{U}_{\tau,K}]$ .** This term is upper bounded by  $k^{\text{in}} \mathbb{E} [\mathcal{U}_{\tau,K}]$  and we write

$$k^{\text{in}} \mathbb{E} [\mathcal{U}_{\tau,K}] \leq \frac{1}{k^{\text{out}}} \sum_{t=1}^{k^{\text{out}}} \sum_{k=1}^{k^{\text{in}}} O\left(\frac{1}{\mathbf{b} M_{t,k+1}}\right) .$$

The RHS is  $O(\epsilon)$ . The associated Monte Carlo complexity is

$$2\mathbf{b} \sum_{t=1}^{k^{\text{out}}} \sum_{k=1}^{k^{\text{in}}} M_{t,k} = O\left(\frac{\sqrt{n}}{\epsilon^2}\right) ,$$

whatever  $\mathbf{a}, \bar{\mathbf{a}} \in [0, 1)$ .

**Acknowledgments.** This work was partly supported by the Fondation Simone et Cino del Duca, under the program OpSiMorE; and by the french Agence Nationale de la Recherche (ANR) under the program ANR-19-CE23 MASDOL.

## Appendix A The condition A 4 in the Monte Carlo case

Following the framework detailed in Section 3.4, let us assume that (i) the intractable quantities  $\mathbf{h}_i(\widehat{S}_{t,k}, B_{t,k+1})$  and  $\mathbf{h}_i(\widehat{S}_{t,k-1}, B_{t,k})$  are of the form

$$\mathbf{h}_i(s, B) = \int_{\mathcal{Z}} H_{\vartheta}(z) \pi_{\vartheta}(dz) , \quad (\text{A1})$$

where  $\vartheta \stackrel{\text{def}}{=} (s, i, B)$ ; and (ii) these integrals are approximated by a Monte Carlo sum: set  $\vartheta_{t,\ell+1,i} \stackrel{\text{def}}{=} (\widehat{S}_{t,\ell}, i, B_{t,\ell+1})$  and

$$\delta_{t,k+1,i} \stackrel{\text{def}}{=} \frac{1}{m_{t,k+1}} \sum_{r=1}^{m_{t,k+1}} \left\{ H_{\vartheta_{t,k+1,i}}(Z_r^{\vartheta_{t,k+1,i}}) - H_{\vartheta_{t,k,i}}(Z_r^{\vartheta_{t,k,i}}) \right\}, \quad (\text{A2})$$

where, conditionally to  $\widehat{S}_{t,k-1}$ ,  $\widehat{S}_{t,k}$ ,  $B_{t,k}$  and  $B_{t,k+1}$ , the samples  $\{Z_r^{\vartheta_{t,\ell,i}}, r \geq 1\}$  are a Markov chain with unique stationary distribution  $\pi_{\vartheta_{t,\ell,i}}(dz)$ ;  $\ell \in \{k, k+1\}$ . Below, we show that **A4** is verified when the Markov chain is ergodic enough. Let us start with introducing few notations from the Markov chain theory (see e.g. [Meyn and Tweedie \(1993\)](#)).

Let  $P$  be a transition kernel onto the measurable set  $(Z, \mathcal{Z})$  and  $\lambda, \pi$  be probability measures on  $(Z, \mathcal{Z})$ . For a measurable function  $\xi : Z \rightarrow [0, +\infty)$ , define

$$\pi(\xi) \stackrel{\text{def}}{=} \int_Z \xi(z) \pi(dz).$$

For any  $r \in \mathbb{N}$ , the  $r$ -iterated transition kernel  $P^r$  is defined by induction:

$$\begin{aligned} P^{r+1}(z, A) &\stackrel{\text{def}}{=} \int_Z P^r(z, dy) P(y, A) \\ &= \int_Z P(z, dy) P^r(y, A), \end{aligned}$$

for all  $z \in Z, A \in \mathcal{Z}$ ; by convention,  $P^0(z, A) \stackrel{\text{def}}{=} \chi_A(z)$  the  $\{0, 1\}$ -valued indicator function and  $P^0(z, A) = \delta_z(A)$ , the Dirac mass at zero. Given a probability measure  $\lambda$  on  $(Z, \mathcal{Z})$ ,  $\lambda P$  stands for the probability measure on  $(Z, \mathcal{Z})$  given by

$$\lambda P(A) \stackrel{\text{def}}{=} \int_Z \lambda(dy) P(y, A), \quad \forall A \in \mathcal{Z}.$$

For a function  $U : Z \rightarrow [1, +\infty)$  such that  $\lambda P^r(U) + \pi(U) < +\infty$ , define the  $U$ -norm of a measurable function  $\xi : Z \rightarrow \mathbb{R}^q$

$$\|\xi\|_U \stackrel{\text{def}}{=} \sup_Z \frac{\|\xi\|}{U};$$

and the  $U$ -norm of the signed measure  $\lambda P^r - \pi$  by

$$\|\lambda P^r - \pi\|_U \stackrel{\text{def}}{=} \sup_{\xi: \|\xi\|_U \leq 1} \|\lambda P^r(\xi) - \pi(\xi)\|.$$

Let us go back to sufficient conditions for verifying **A4**. Denote by  $P_\vartheta$  a Markov transition kernel with

invariant distribution  $\pi_\vartheta(dz)$ : at iteration  $(t, k+1)$ , conditionally to  $(\widehat{S}_{t,k-1}, \widehat{S}_{t,k}, B_{t,k}, B_{t,k+1})$ , the chains  $\{Z_r^{\vartheta_{t,k}}, r \geq 0\}$  and  $\{Z_r^{\vartheta_{t,k+1}}, r \geq 0\}$  are Markov chains with transition kernels  $P_{\vartheta_{t,k}}$  and  $P_{\vartheta_{t,k+1}}$  respectively. They have the same initial value  $\lambda$ . Assume

**A5. 1.** *There exists a measurable function  $U : Z \rightarrow [1, +\infty)$  such that*

$$H_\star \stackrel{\text{def}}{=} \sup_{(s,i,B) \in \mathcal{S} \times [n]^\star \times \mathcal{P}_+^q} \|H_\vartheta\|_U < +\infty,$$

where  $H_\vartheta$  is defined by [\(A1\)](#).

**2.** *There exist a function  $\rho : \mathbb{N} \rightarrow [0, 1]$  and a positive constant  $C_{\text{MC}}$  such that for any  $r \in \mathbb{N}$ ,*

$$\sup_{\vartheta \in \mathcal{S} \times [n]^\star \times \mathcal{P}_+^q} \|\lambda P_\vartheta^r - \pi_\vartheta\|_U \leq C_{\text{MC}} \rho(r).$$

*In addition,  $\sum_{r \geq 1} \rho(r) < +\infty$ .*

**3.** *Let  $\vartheta \in \mathcal{S} \times [n]^\star \times \mathcal{P}_+^q$ . Let  $\{Z_r^\vartheta, r \geq 1\}$  be a Markov chain with transition kernel  $P_\vartheta$  and initial distribution  $\lambda$ . There exists a positive constant  $C'_{\text{MC}}$  such that for any  $\vartheta \in \mathcal{S} \times [n]^\star \times \mathcal{P}_+^q$  and  $m' \in \mathbb{N}_\star$ ,*

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{r=1}^{m'} \{H_\vartheta(Z_r^\vartheta) - \pi_\vartheta(H_\vartheta)\} \right\|^2 \right] \\ \leq H_\star^2 C'_{\text{MC}} m'. \end{aligned}$$

**A5-item 2** is a uniform-in- $s$  ergodicity condition. Sufficient conditions for it are provided in ([Fort et al, 2011](#), Lemma 2.3.) in the case of a geometric rate  $\rho(r) = \kappa^r$  for some  $\kappa \in (0, 1)$ . By adapting ([Andrieu et al, 2015](#), Theorem 1), similar conditions can be obtained in the case of a subgeometric rate  $\rho(r)$ . Sufficient conditions for **A5-item 3** can be obtained from a trivial adaptation of ([Fort and Moulines, 2003](#), Proposition 12).

We prove the following result.

**Proposition A.1** *Assume **A 5**. Let  $\delta_{t,k+1,i}$  be given by [\(A2\)](#), where conditionally to  $(\widehat{S}_{t,k-1}, \widehat{S}_{t,k}, B_{t,k}, B_{t,k+1})$ ,  $\{Z_r^{\vartheta_{t,\ell,i}}, r \geq 0\}$  is a Markov chain with transition kernel  $P_{\vartheta_{t,\ell,i}}$  and initial distribution  $\lambda$ , for  $\ell \in \{k, k+1\}$ . Then **A4** is verified with  $m_{t,k+1} = M_{t,k+1} = \bar{M}_{t,k+1} \leftarrow m_{t,k+1}$ ,*

$$C_b = C_{vb} \stackrel{\text{def}}{=} 2H_\star C_{\text{MC}} \sum_{r \geq 1} \rho(r) \quad \text{and} \\ C_v \stackrel{\text{def}}{=} 2H_\star^2 C'_{\text{MC}}.$$

*Proof* We will use the notations

$$P_{\ell,i} \stackrel{\text{def}}{=} P_{\vartheta_{t,\ell+1,i}}, \quad \pi_{\ell,i} \stackrel{\text{def}}{=} \pi_{\vartheta_{t,\ell+1,i}}, \quad H_{\ell,i} \stackrel{\text{def}}{=} H_{\vartheta_{t,\ell+1,i}}.$$

• **Expression of  $\mu_{t,k+1}$ .** We have

$$\mu_{t,k+1,i} \stackrel{\text{def}}{=} \frac{1}{m_{t,k+1}} \sum_{r=1}^{m_{t,k+1}} (\lambda P_{k,i}^r H_{k,i} - \pi_{k,i}(H_{k,i})) \\ - \frac{1}{m_{t,k+1}} \sum_{r=1}^{m_{t,k+1}} (\lambda P_{k-1,i}^r H_{k-1,i} - \pi_{k-1,i}(H_{k-1,i})).$$

• **The condition A4-Item 2.** By A5 and since  $\widehat{\mathcal{S}}_\bullet \in \mathcal{S}$ , we write

$$\sup_{k,i} \|\lambda P_{k,i}^r H_{k,i} - \pi_{k,i}(H_{k,i})\| \leq H_\star C_{\text{MC}} \rho(r).$$

This implies that

$$\|\mu_{t,k+1,i}\| \leq 2H_\star C_{\text{MC}} \frac{1}{m_{t,k+1}} \sum_{r=1}^{m_{t,k+1}} \rho(r).$$

Since  $\sum_r \rho(r) < \infty$ , the RHS is of the form  $C_b/m_{t,k+1}$  with  $C_b \stackrel{\text{def}}{=} 2H_\star C_{\text{MC}} \sum_r \rho(r)$ .

• **The condition A4-Item 3.** We write

$$\sigma_{t,k+1,i}^2 \leq \mathbb{E} \left[ \|\xi_{t,k+1,i}\|^2 | \mathcal{P}_{t,k+1/2} \right].$$

Then, we have

$$\mathbb{E} \left[ \|\xi_{t,k+1,i}\|^2 | \mathcal{P}_{t,k+1/2} \right] \\ \leq 2 \sup_{s,i} \mathbb{E} \left[ \left\| \frac{1}{m_{t,k+1}} \sum_{r=1}^{m_{t,k+1}} H_{s,i}(Z_r^{s,i}) - \pi_{s,i}(H_{s,i}) \right\|^2 \right]$$

and the RHS is upper bounded by  $2H_\star^2 C'_{\text{MC}}/m_{t,k+1}$  by A5-Item 3.

We also have

$$\frac{1}{n} \sum_{i=1}^n \|\mu_{t,k+1,i}\| - \frac{1}{n} \sum_{j=1}^n \|\mu_{t,k+1,j}\| \leq \frac{1}{n} \sum_{i=1}^n \|\mu_{t,k+1,i}\|^2.$$

From the upper bound on  $\|\mu_{t,k+1,i}\|$  above, we have

$$\|\mu_{t,k+1,i}\|^2 \leq \frac{C_b^2}{m_{t,k+1}^2}.$$

This concludes the proof.  $\square$

## Appendix B Supplementary materials for Section 5

### B.1 The penalized log-likelihood criterion

The observations are assumed independent, so the log-likelihood is given by

$$\theta \mapsto \sum_{i=1}^n \log \int_{\mathbb{R}^d} (1 + \exp(-y_i \langle X_i, z_i \rangle))^{-1} \\ \times \frac{1}{\sqrt{2\pi}^d \sigma^d} \exp(-(2\sigma^2)^{-1} \|z_i - \theta\|^2) dz_i.$$

The penalty term is  $-n\tau\|\theta\|^2$ .

**Lemma B.1** *The sum of the log-likelihood and the penalty term is equal to*

$$- \frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \theta^\top \sum_{i=1}^n \frac{X_i X_i^\top}{\|X_i\|^2} \theta - n\tau\|\theta\|^2 \\ + \sum_{i=1}^n \log \int_{\mathbb{R}} \frac{\exp(x \langle X_i, \theta \rangle / (\sigma^2 \|X_i\|))}{1 + \exp(-y_i \|X_i\| x)} \exp(-\frac{x^2}{2\sigma^2}) dx.$$

*Proof* Let  $i \in [n]^\star$ . Define an orthogonal  $d \times d$  matrix  $Q$  with columns denoted by  $(Q_1, \dots, Q_d)$ , such that  $Q_1 \stackrel{\text{def}}{=} X_i / \|X_i\|$ . We have  $\langle z_i, X_i \rangle = \|X_i\| \langle Q_1, z_i \rangle$ ,  $\langle z_i, \theta \rangle = \langle Q^\top z_i, Q^\top \theta \rangle$  and  $\|z_i\|^2 = \|Q^\top z_i\|^2$ . This implies that

$$-y_i \langle z_i, X_i \rangle = -y_i \|X_i\| \langle Q_1, z_i \rangle \\ \|z_i - \theta\|^2 = \|\theta\|^2 + \|Q^\top z_i\|^2 - 2 \langle Q^\top z_i, Q^\top \theta \rangle$$

so that the log-likelihood of the observation  $Y_i$  is (up to the additive constant  $C_1 \stackrel{\text{def}}{=} -d \ln \sigma - (d/2) \ln(2\pi)$ )

$$y_i \mapsto -\frac{\|\theta\|^2}{2\sigma^2} + \log \int_{\mathbb{R}^d} (1 + \exp(-y_i \|X_i\| \langle Q_1, z_i \rangle))^{-1} \\ \times \exp(-(2\sigma^2)^{-1} (\|Q^\top z_i\|^2 - 2 \langle Q^\top z_i, Q^\top \theta \rangle)) dz_i.$$

By a change of variable  $v = (v_1, \dots, v_q) \leftarrow Q^\top z_i$ , the logarithm of the integral is equal to

$$\log \int_{\mathbb{R}^d} \frac{\exp(-(2\sigma^2)^{-1} (\|v\|^2 - 2 \langle v, Q^\top \theta \rangle))}{1 + \exp(-y_i \|X_i\| v_1)} dv \\ = \log \int_{\mathbb{R}} \frac{\exp(-(2\sigma^2)^{-1} (v_1^2 - 2v_1 \langle Q_1, \theta \rangle))}{1 + \exp(-y_i \|X_i\| v_1)} dv_1$$

$$+ \sum_{u=2}^d \log \int_{\mathbb{R}} \exp \left( -\frac{1}{2\sigma^2} \{v_u^2 - 2v_u \langle Q_u, \theta \rangle\} \right) dv_u .$$

The last  $(d-1)$  integrals have a closed form. Observe indeed that  $v_u^2 - 2v_u \langle Q_u, \theta \rangle = (v_u - \langle Q_u, \theta \rangle)^2 - (\langle Q_u, \theta \rangle)^2$  so that up to the additive constant  $C_2 \stackrel{\text{def}}{=} (d-1)\{\log(2\pi)/2 + \log \sigma\}$

$$\begin{aligned} & \sum_{u=2}^d \log \int_{\mathbb{R}} \exp \left( -\frac{1}{2\sigma^2} \{v_u^2 - 2v_u \langle Q_u, \theta \rangle\} \right) dv_u \\ &= \sum_{u=2}^d \frac{(\langle Q_u, \theta \rangle)^2}{2\sigma^2} = \frac{\|\theta\|^2 - (\langle Q_1, \theta \rangle)^2}{2\sigma^2} \\ &= \frac{\|\theta\|^2 - (\langle X_i, \theta \rangle)^2 / \|X_i\|^2}{2\sigma^2} . \end{aligned}$$

This concludes the proof; the constant (w.r.t. to  $\theta$ ) is equal to  $C_1 + C_2$ .  $\square$

## B.2 Proof of Lemma 5.1

The criterion  $F$  is equal to  $-\mathcal{L}(\theta) - \log(2\pi\sigma^2)/2$ , where  $\mathcal{L}(\theta)$  is the normalized penalized log-likelihood.

The likelihood is the product of probabilities, taking values in  $(0,1)$ ; therefore, its logarithm is negative. The penalized log-likelihood is upper bounded  $-\text{pen}(\theta) = -n\tau\|\theta\|^2$ . The normalized penalized log-likelihood is upper bounded  $-\tau\|\theta\|^2$ . Therefore the criterion is lower bounded by  $\tau\|\theta\|^2 - \ln(2\pi\sigma^2)/2$ .

On the other hand, the minimum of the criterion is smaller than the value of the criterion at  $\theta = 0$ . Let us show that this value is  $(\ln 4)/n - \ln(2\pi\sigma^2)/2$ . This will imply that the minimizers of the criterion are in the set  $\{\theta \in \mathbb{R}^d : \tau\|\theta\|^2 \leq \ln 4\}$  and conclude the proof.

We have  $\text{pen}(0) = 0$ . Let us lower bound the likelihood of an observation  $Y_i = +1$  at  $\theta = 0$ . The likelihood is equal to

$$\frac{1}{\sqrt{2\pi}^d \sigma^d} \int_{\mathbb{R}^d} \frac{\exp(-(2\sigma^2)^{-1}\|z_i\|^2)}{1 + \exp(-\langle X_i, z_i \rangle)} dz_i .$$

By using the same change of variable than in the proof of Lemma B.1, it is equal to

$$\begin{aligned} & \frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} \frac{\exp(-x^2/(2\sigma^2))}{1 + \exp(-\|X_i\|x)} dx \\ & \left( \frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} \exp(-x^2/(2\sigma^2)) dx \right)^{d-1} , \end{aligned}$$

and is lower bounded by (note that the  $(d-1)$  identical integrals are equal to one)

$$\frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}^+} \frac{\exp(-x^2/(2\sigma^2))}{1 + \exp(-\|X_i\|x)} dx ,$$

which is in turn lower bounded by  $1/4$  since  $1 + \exp(-\|X_i\|x) \leq 2$  for all  $x \geq 0$ .

The proof for the case  $Y_i = -1$  is on the same lines and is omitted.

This implies that the likelihood of the  $n$  variables is lower bounded by  $1/4^n$ ; the normalized log-likelihood is lower bounded by  $-\ln 4$ ; the criterion is upper bounded by  $\ln 4 - \ln(2\pi\sigma^2)/2$ .

## B.3 The optimization problem seen as an EM

We established that for any  $\theta \in \mathbb{R}^d$ ,  $\nabla F(\theta) = n^{-1} \sum_{i=1}^n G_i(\theta)$  where

$$G_i(\theta) \stackrel{\text{def}}{=} 2U\theta - \frac{X_i}{\sigma^2 \|X_i\|} \int_{\mathbb{R}} z \pi_{\theta,i}(z) dz ,$$

and  $\pi_{\theta,i}(z)$  is the probability density proportional to (20).

From the expressions of  $\phi, \psi$  and  $S(Y_i, z)$ , we obtain that  $\mathbb{T}$ , defined by Proposition 2.1, is given by  $\mathbb{T}(s) \stackrel{\text{def}}{=} U^{-1}s/2$  for any  $s \in \mathbb{R}^d$ . This implies that for any  $s \in \mathbb{R}^d$ ,

$$\begin{aligned} h_i(s, B) & \stackrel{\text{def}}{=} \int_{\mathbb{R}} S(Y_i, z) \pi_{\mathbb{T}(s),i}(z) dz - s \\ & = \frac{X_i}{\sigma^2 \|X_i\|} \int_{\mathbb{R}} z \pi_{B_s,i}(z) dz - s . \end{aligned}$$

For any  $s \in \mathbb{R}^d$ , let us find the matrix  $\mathbb{B}(s)$  satisfying  $\nabla(F \circ \mathbb{T})(s) = -n^{-1} \sum_{i=1}^n \mathbb{B}(s) h_i(s)$  (see (6)). We have  $\nabla(F \circ \mathbb{T}) = \nabla F(B \cdot) = B^\top (\nabla F)(B \cdot) = B (\nabla F)(B \cdot)$ . This yields

$$\begin{aligned} \nabla(F \circ \mathbb{T})(s) &= B \frac{1}{n} \sum_{i=1}^n G_i(Bs) \\ &= B \frac{1}{n} \sum_{i=1}^n \left( 2U\theta - \frac{X_i}{\sigma^2 \|X_i\|} \int_{\mathbb{R}} z \pi_{B_s,i}(z) dz \right) \\ &= Bs - B \frac{1}{n} \sum_{i=1}^n \frac{X_i}{\sigma^2 \|X_i\|} \int_{\mathbb{R}} z \pi_{B_s,i}(z) dz \end{aligned}$$

$$= -B \frac{1}{n} \sum_{i=1}^n h_i(s) .$$

This yields  $B(s) \stackrel{\text{def}}{=} B$  for any  $s \in \mathbb{R}^d$ .

## B.4 Proof of Lemma 5.2

Let  $i \in [n]^*$  and  $\theta \in \mathbb{R}^d$ . **Step 1.** By using

$$\begin{aligned} & -z^2/(2\sigma^2) + z \langle X_i, \theta \rangle / (\sigma^2 \|X_i\|) \\ &= -(z - \langle X_i, \theta \rangle / \|X_i\|)^2 / (2\sigma^2) \\ & \quad + (\langle X_i, \theta \rangle)^2 / (2\sigma^2 \|X_i\|^2) , \end{aligned}$$

we write

$$\pi_{\theta,i}(z) = \frac{\exp\left(-\frac{(z - \langle X_i, \theta \rangle / \|X_i\|)^2}{2\sigma^2}\right)}{Z_{\theta,i} \left(1 + \exp(-y_i \|X_i\| z)\right)}$$

where  $Z_{\theta,i}$  is the normalizing constant. Second, we use  $z \pi_{\theta,i}(z) = (z - a_i) \pi_{\theta,i}(z) + a_i \pi_{\theta,i}(z)$  with  $a_i \leftarrow \langle X_i, \theta \rangle / \|X_i\|$  and since  $\int_{\mathbb{R}} \pi_{\theta,i}(z) dz = 1$ , we obtain

$$\mathcal{I}_i(\theta) = \frac{\langle X_i, \theta \rangle}{\|X_i\|} + \int_{\mathbb{R}} \left( z - \frac{\langle X_i, \theta \rangle}{\|X_i\|} \right) \pi_{\theta,i}(z) dz .$$

Finally, the integral in the RHS being of the form

$$\sigma^2 \int_{\mathbb{R}} \frac{f'(z)}{Z_{\theta,i} \left(1 + \exp(-y_i \|X_i\| z)\right)} dz$$

with

$$f(z) \stackrel{\text{def}}{=} -\exp\left(-\frac{(z - \langle X_i, \theta \rangle / \|X_i\|)^2}{2\sigma^2}\right) ,$$

we use an integration by parts. Upon noting that the derivative of  $z \mapsto 1/(1 + \exp(-y_i \|X_i\| z))$  is

$$y_i \|X_i\| \frac{\exp(-y_i \|X_i\| z)}{(1 + \exp(-y_i \|X_i\| z))^2} ,$$

we write

$$\begin{aligned} & \int_{\mathbb{R}} \frac{f'(z)}{1 + \exp(-y_i \|X_i\| z)} dz \\ &= -y_i \|X_i\| \int_{\mathbb{R}} f(z) \frac{\exp(-y_i \|X_i\| z)}{(1 + \exp(-y_i \|X_i\| z))^2} dz . \end{aligned}$$

Therefore, the conclusion of this first step is

$$\begin{aligned} \mathcal{I}_i(\theta) &= \left\langle \frac{X_i}{\|X_i\|}, \theta \right\rangle \\ & \quad + y_i \|X_i\| \sigma^2 \int_{\mathbb{R}} \frac{\pi_{\theta,i}(z)}{1 + \exp(y_i \|X_i\| z)} dz . \end{aligned}$$

**Step 2.** This step is classical in the MCMC literature (see e.g. [Choi and Hobert \(2013\)](#) and references therein). We prove that for any  $z \in \mathbb{R}$ ,

$$\pi_{\theta,i}(z) = \int_0^{+\infty} \bar{\pi}_{\theta,i}(z, \omega) d\omega .$$

By ([Polson et al, 2013](#), Theorem 1), it holds

$$\begin{aligned} \frac{1}{1 + \exp(-y_i \|X_i\| z)} &= \frac{1}{2} \exp(y_i \|X_i\| z/2) \\ & \quad \times \int_0^{+\infty} \exp(-\omega \|X_i\|^2 z^2/2) \mathbf{p}(\omega; 1) d\omega , \end{aligned}$$

where  $\mathbf{p}(\omega; b)d\omega$  is a Polya-Gamma distribution with parameter  $b$ . This implies that  $\pi_{\theta,i}(z)$  is equal to

$$\begin{aligned} & \exp\left(\frac{y_i \|X_i\| z}{2} - \frac{(z - \langle X_i, \theta \rangle / \|X_i\|)^2}{2\sigma^2}\right) \\ & \quad \times \frac{1}{2Z_{\theta,i}} \int_0^{+\infty} \exp(-\omega \|X_i\|^2 z^2/2) \mathbf{p}(\omega; 1) d\omega . \end{aligned}$$

This concludes the proof.

## B.5 The assumption A4 is verified.

Define the Markov kernel with density

$$P_{s,i}(z; z') \stackrel{\text{def}}{=} \left( \int_0^{\infty} \pi_2(\omega|z; i) \pi_1(z'|\omega; s, i) d\omega \right)$$

w.r.t. the Lebesgue measure on  $\mathbb{R}$ ; here,  $\pi_1(z'|\omega; s, i)$  is the density of a Gaussian distribution with expectation  $\mathbf{m}_{s,i}(\omega)$  and variance  $v_i(\omega)$  given by

$$\begin{aligned} v_i(\omega) &\stackrel{\text{def}}{=} \frac{\sigma^2}{1 + \omega \sigma^2 \|X_i\|^2} , \\ \mathbf{m}_{s,i}(\omega) &\stackrel{\text{def}}{=} v_i(\omega) \left( \frac{1}{\sigma^2} \left\langle \frac{X_i}{\|X_i\|}, Bs \right\rangle + \frac{1}{2} y_i \|X_i\| \right) ; \end{aligned}$$

and  $\pi_2(\omega|z; i)$  is a Polya-Gamma distribution with parameter  $(1, \|X_i\| z)$ . The Gibbs kernel described by Lemma 5.2 and targeting the density

distribution  $\bar{\pi}_{B_{s,i}}(z, \omega) dz d\omega$ , produces a Markov chain  $\{(Z_r^{s,i}, \Omega_r^{s,i}), r \geq 0\}$  such that the marginal  $\{Z_r^{s,i}, r \geq 0\}$  is a Markov chain with transition kernel  $P_{s,i}(z; z') dz'$ . We apply the results of (Choi and Hobert, 2013, Proposition 3.1) with

$y \in \{0, 1\}$	$y_i \in \{-1, 1\}$
$n$	1
$\Omega(\omega)$	$\omega$
$X$	$\ X_i\ $
$B$	$\sigma^2$
$b$	$\langle X_i, B_s \rangle / \ X_i\ $

**Table B1** [left] The notations of Choi and Hobert (2013). [right] the notations in this paper.

This yields

$$\begin{aligned} & \int_A P_{s,i}(z; z') dz' \\ & \geq \varepsilon \int_A \exp(-0.5(x - m_\star)^2 / v_\star^2) dx \quad (\text{B3}) \end{aligned}$$

where

$$\varepsilon \stackrel{\text{def}}{=} \inf_{s \in \mathcal{S}, i \in [n]^\star} \frac{\exp\left(-\frac{1}{4} - \frac{\{m_{s,i}(1/2)\}^2 \sigma^2 \|X_i\|^2}{4 v_i(1/2)}\right)}{2\sqrt{1 + \sigma^2 \|X_i\|^2 / 2}}$$

and  $(m_\star, v_\star)$  satisfy for any  $x \in \mathbb{R}$ ,

$$\begin{aligned} & \inf_{s \in \mathcal{S}, i \in [n]^\star} \exp(-0.5(x - m_{s,i}(1/2))^2 / v_i(1/2)) \\ & \geq \exp(-0.5(x - m_\star)^2 / v_\star^2). \end{aligned}$$

**Lemma B.2** *Since  $\mathcal{S}$  is bounded, then  $\varepsilon > 0$  and  $m_\star, v_\star$  exist in  $\mathbb{R} \times (0, +\infty)$ .*

The minorization condition (B3) implies that the kernel  $P_{s,i}(z, z') dz'$  is uniformly ergodic, uniformly in  $s, i$  and  $z$ . By (Meyn and Tweedie, 1993, Theorem 16.0.2.) and (Fort and Moulines, 2003, Proposition 1), A 5-Item 2 and A 5-Item 3 are satisfied.

## Appendix C Detailed proofs

### C.1 Proof of (17)

Let  $t \in [k^{\text{out}}]^\star$ . The sequence given by  $\gamma_{t,k+1} \stackrel{\text{def}}{=} \prod_{j=0}^k \left(1 + \frac{2C_b}{m_{t,j+1}}\right)^{-1} \gamma_{t,0}$  for any  $k \geq 0$ , satisfies

$$\gamma_{t,k+1} \left(1 + \frac{2C_b}{m_{t,k+1}}\right) \leq \gamma_{t,k}.$$

A sufficient condition for the property  $\Lambda_{t,k+1} \in (0, 1/2)$  to hold is  $\mathfrak{a}\gamma_{t,0}^2 + \bar{\mathfrak{a}}\gamma_{t,0} - 1/2 < 0$  where

$$\begin{aligned} \bar{\mathfrak{a}} & \stackrel{\text{def}}{=} \frac{L_{\dot{W}}}{v_{\min}}, \\ \mathfrak{a} & \stackrel{\text{def}}{=} L^2 \frac{2v_{\max} k_t^{\text{in}}}{v_{\min} \mathfrak{b}} \left(1 + \frac{2C_{vb}}{\sqrt{\mathfrak{b}} \bar{M}_{t,k+1}}\right). \end{aligned}$$

The function  $x \mapsto \mathfrak{a}x^2 + \bar{\mathfrak{a}}x - 1/2$  possesses two roots: one is positive and one is negative. The positive one is given by  $(-\bar{\mathfrak{a}} + \sqrt{\bar{\mathfrak{a}}^2 + 2\mathfrak{a}})/(2\mathfrak{a})$ ; it is equal to (17).

### C.2 Proof of (18) and (19)

We write  $\mathbf{S}_{t,0} - \mathfrak{h}(\widehat{\mathbf{S}}_{t,0}) = U + V$  where

$$\begin{aligned} U & \stackrel{\text{def}}{=} \frac{1}{\mathfrak{b}'_t} \sum_{i \in \mathcal{B}_{t,0}} \left( \delta_{t,0,i} - \mathbb{E} \left[ \delta_{t,0,i} | \widehat{\mathbf{S}}_{t,0}, \mathcal{B}_{t,0} \right] \right), \\ V & \stackrel{\text{def}}{=} \frac{1}{\mathfrak{b}'_t} \sum_{i \in \mathcal{B}_{t,0}} \mathbb{E} \left[ \delta_{t,0,i} | \widehat{\mathbf{S}}_{t,0}, \mathcal{B}_{t,0} \right] - \mathfrak{h}(\widehat{\mathbf{S}}_{t,0}). \end{aligned}$$

We have  $\mathbb{E} [\|U + V\|^2] = \mathbb{E} [\|U\|^2] + \mathbb{E} [\|V\|^2]$  by definition of the conditional expectation. Since  $\delta_{t,0,i}$  is an unbiased random approximation of  $\mathfrak{h}_i(\widehat{\mathbf{S}}_{t,0})$ , we have  $\mathbb{E} \left[ \delta_{t,0,i} | \widehat{\mathbf{S}}_{t,0}, \mathcal{B}_{t,0} \right] = \mathfrak{h}_i(\widehat{\mathbf{S}}_{t,0})$ .

In the case  $\mathcal{B}_{t,0} = \{1, \dots, n\}$ , then  $V = 0$ . Therefore,

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{S}_{t,0} - \mathfrak{h}(\widehat{\mathbf{S}}_{t,0})\|^2 \right] \\ & = \frac{1}{n^2} \mathbb{E} \left[ \left\| \sum_{i=1}^n \{ \delta_{t,0,i} - \mathbb{E} [ \delta_{t,0,i} | \widehat{\mathbf{S}}_{t,0} ] \} \right\|^2 \right] \\ & = \frac{1}{n^2} \sum_{i=1}^n \sigma_{t,0,i}^2, \end{aligned}$$

where we used that the variables  $\{\delta_{t,0,i}, i \in [n]^\star\}$  are independent conditionally to  $\widehat{\mathbf{S}}_{t,0}$ , and with variance  $\sigma_{t,0,i}^2$ .

In the case  $\mathcal{B}_{t,0}$  is a subset of  $[n]^*$  of cardinality  $b'_t$ , then we write

$$\begin{aligned}\mathbb{E} [\|U\|^2] &= \frac{1}{(b'_t)^2} \mathbb{E} \left[ \mathbb{E} [\|U\|^2 | \mathcal{B}_{t,0}, \widehat{S}_{t,0}] \right] \\ &= \frac{1}{b'_t} \mathbb{E} \left[ \frac{1}{b'_t} \sum_{i \in \mathcal{B}_{t,0}} \sigma_{t,0,i}^2 \right]\end{aligned}$$

and we conclude by Lemma 7.1. Again from Lemma 7.1, we have

$$\mathbb{E} [\|V\|^2] \leq \frac{1}{b'_t n} \sum_{i=1}^n \|\mathbf{h}_i(\widehat{S}_{t,0}) - \mathbf{h}(\widehat{S}_{t,0})\|^2,$$

with an equality when  $\mathcal{B}_{t,0}$  is sampled with replacement. This concludes the proof.

### C.3 Proof of Lemma 7.1

Set, for ease of notations,

$$\mathcal{B} \stackrel{\text{def}}{=} \mathcal{B}_{t,k}, \quad \mathbf{h}_{\mathcal{B}} \stackrel{\text{def}}{=} \frac{1}{b} \sum_{i \in \mathcal{B}} \mathbf{h}_i.$$

#### C.3.1 Case with replacement

We write  $\mathcal{B} = \{I_1, \dots, I_b\}$  where the r.v.  $I_i$ 's are independent, and uniformly distributed on  $[n]^*$ .

- Then

$$\mathbb{E} [f_{\mathcal{B}}] = \frac{1}{b} \sum_{\ell=1}^b \mathbb{E} [f_{I_{\ell}}] = \mathbb{E} [f_{I_1}] = n^{-1} \sum_{i=1}^n f_i(u).$$

- Set  $\bar{f} \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n f_i$ . We have, by using that the r.v.  $\{I_1, \dots, I_b\}$  are independent,

$$\begin{aligned}\mathbb{E} [\|f_{\mathcal{B}} - \bar{f}\|^2] &= \mathbb{E} \left[ \left\| \frac{1}{b} \sum_{\ell=1}^b (f_{I_{\ell}} - \bar{f}) \right\|^2 \right] \\ &= \frac{1}{b^2} \sum_{\ell=1}^b \mathbb{E} [\|f_{I_{\ell}} - \bar{f}\|^2] \\ &= \frac{1}{b} \mathbb{E} [\|f_{I_1} - \bar{f}\|^2] = \frac{1}{b n} \sum_{i=1}^n \|f_i - \bar{f}\|^2.\end{aligned}$$

- Since the variance of the sum is the sum of the variance for independent r.v.

$$\mathbb{E} [\|\mathbf{h}_{\mathcal{B}}(u) - \mathbf{h}_{\mathcal{B}}(u') - \{\mathbf{h}(u) - \mathbf{h}(u')\}\|^2]$$

$$= \frac{1}{b^2} \sum_{\ell=1}^b \mathbb{E} [\|\mathbf{h}_{I_{\ell}}(u) - \mathbf{h}_{I_{\ell}}(u') - \mathbf{h}(u) + \mathbf{h}(u')\|^2].$$

Then, since  $I_{\ell}$  is uniformly distributed on  $[n]^*$ ,

$$\begin{aligned}\mathbb{E} [\|\mathbf{h}_{I_{\ell}}(u) - \mathbf{h}_{I_{\ell}}(u') - \mathbf{h}(u) + \mathbf{h}(u')\|^2] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\mathbf{h}_i(u) - \mathbf{h}_i(u')\|^2] - \|\mathbf{h}(u) - \mathbf{h}(u')\|^2 \\ &\leq \|u - u'\|^2 \frac{1}{n} \sum_{i=1}^n L_i^2 - \|\mathbf{h}(u) - \mathbf{h}(u')\|^2.\end{aligned}\tag{C4}$$

#### C.3.2 Case without replacement

Set  $\mathcal{B} = \{I_1, \dots, I_b\}$  and  $\bar{f} \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n f_i$ .  $I_1$  is a uniform random variable on  $[n]^*$  so that  $\mathbb{E} [f_{I_1}] = \bar{f}$ .

- Conditionally to  $I_1$ ,  $I_2$  is a uniform random variable on  $[n]^* \setminus \{I_1\}$ . Therefore

$$\begin{aligned}\mathbb{E} [f_{I_2}] &= \frac{1}{n-1} \left( \sum_{j=1}^n f_j - \mathbb{E} [f_{I_1}] \right) \\ &= \frac{n}{n-1} \bar{f} - \frac{1}{n-1} \bar{f} = \bar{f}.\end{aligned}$$

By induction, for any  $\ell \geq 2$ ,

$$\begin{aligned}\mathbb{E} [f_{I_{\ell}}] \\ &= \frac{1}{n-\ell+1} \left( \sum_{j=1}^n f_j - \sum_{r=1}^{\ell-1} \mathbb{E} [f_{I_r}] \right) \\ &= \frac{n}{n-\ell+1} \bar{f} - \frac{\ell-1}{n-\ell+1} \bar{f} = \bar{f}.\end{aligned}$$

As a conclusion,  $b^{-1} \sum_{\ell=1}^b \mathbb{E} [f_{I_{\ell}}] = \bar{f}$ .

- Let  $u, u' \in \mathcal{S}$ ; set  $\phi(I_{\ell}) \stackrel{\text{def}}{=} \mathbf{h}_{I_{\ell}}(u) - \mathbf{h}(u) - \mathbf{h}_{I_{\ell}}(u') + \mathbf{h}(u')$ . Then  $\mathbb{E} [\phi(I_{\ell})] = 0$ . First, we prove by induction that  $\mathbb{E} [\|\phi(I_{\ell})\|^2] = \mathbb{E} [\|\phi(I_1)\|^2]$ . Upon noting that  $I_1$  is a uniform random variable on  $[n]^*$  and by using the induction assumption,

$$\begin{aligned}(n-\ell+1) \mathbb{E} [\|\phi(I_{\ell})\|^2] \\ &= \left( \sum_{i=1}^n \|\phi(i)\|^2 - \mathbb{E} \left[ \sum_{r=1}^{\ell-1} \|\phi(I_r)\|^2 \right] \right) \\ &= n \mathbb{E} [\|\phi(I_1)\|^2] - (\ell-1) \mathbb{E} [\|\phi(I_1)\|^2],\end{aligned}$$

which concludes the induction. Second, let us prove that for any  $\ell \geq 0$ ,

$$\mathbb{E} \left[ \left\| \sum_{r=1}^{\ell+1} \phi(I_r) \right\|^2 \right] \leq (\ell + 1) \mathbb{E} [\|\phi(I_1)\|^2] . \quad (\text{C5})$$

Since  $\sum_{i=1}^n \phi(i) = n \mathbb{E} [\phi(I_1)] = 0$ ,

$$\begin{aligned} & \mathbb{E} \left[ \left\langle \sum_{p=1}^{\ell} \phi(I_p), \phi(I_{\ell+1}) \right\rangle \right] \\ &= \mathbb{E} \left[ \left\langle \sum_{p=1}^{\ell} \phi(I_p), \mathbb{E} [\phi(I_{\ell+1}) | I_1, \dots, I_{\ell}] \right\rangle \right] \\ &= \frac{1}{n - \ell} \mathbb{E} \left[ \left\langle \sum_{p=1}^{\ell} \phi(I_p), \sum_{i=1}^n \phi(i) - \sum_{p=1}^{\ell} \phi(I_p) \right\rangle \right] \\ &= -\frac{1}{n - \ell} \mathbb{E} \left[ \left\| \sum_{p=1}^{\ell} \phi(I_p) \right\|^2 \right] , \end{aligned}$$

so that

$$\begin{aligned} & \mathbb{E} \left[ \left\| \sum_{p=1}^{\ell+1} \phi(I_p) \right\|^2 \right] \\ &= \left( 1 - \frac{2}{n - \ell} \right) \mathbb{E} \left[ \left\| \sum_{p=1}^{\ell} \phi(I_p) \right\|^2 \right] + \mathbb{E} [\|\phi(I_{\ell+1})\|^2] \\ &\leq (\ell + 1) \mathbb{E} [\|\phi(I_1)\|^2] . \end{aligned}$$

The proof of the first bound follows from (C5) and (C4) since here again,  $I_1$  is uniformly distributed on  $[n]^*$ .

• The proof of the second bound is similar (change the definition of  $\phi(I_{\ell})$ ); it is omitted.

## References

Allen-Zhu Z (2018) Natasha 2: Faster Non-Convex Optimization Than SGD. In: Bengio S, Wallach H, Larochelle H, et al (eds) *Advances in Neural Information Processing Systems*, vol 31. Curran Associates, Inc.

Allen-Zhu Z, Hazan E (2016) Variance reduction for faster non-convex optimization. In: Balcan M, Weinberger K (eds) *33rd International Conference on Machine Learning, ICML*

2016. International Machine Learning Society (IMLS), 33rd International Conference on Machine Learning, ICML 2016, pp 1093–1101

Andrieu C, Fort G, Vihola M (2015) Quantitative convergence rates for subgeometric markov chains. *J Appl Probab* 52(2):391–404. <https://doi.org/10.1239/jap/1437658605>

Atchadé Y, Fort G, Moulines E (2017) On Perturbed Proximal Gradient Algorithms. *Journal of Machine Learning Research* 18(10):1–33

Bauschke HH, Combettes PL (2011) *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, 1st edn. Springer Publishing Company, Incorporated, <https://doi.org/10.1007/978-1-4419-9467-7>

Beck A (2017) *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, <https://doi.org/10.1137/1.9781611974997>

Becker S, Fadili J (2012) A quasi-Newton proximal splitting method. In: Pereira F, Burges C, Bottou L, et al (eds) *Advances in Neural Information Processing Systems*, vol 25. Curran Associates, Inc.

Becker S, Fadili J, Ochs P (2019) On quasi-newton forward-backward splitting: Proximal calculus and convergence. *SIAM Journal on Optimization* 29(4):2445–2481. <https://doi.org/10.1137/18M1167152>

Benveniste A, Métivier M, Priouret P (1990) *Adaptive Algorithms and Stochastic Approximations*. Springer Verlag, <https://doi.org/https://doi.org/10.1007/978-3-642-75894-2>

Bonettini S, Porta F, Ruggiero V, et al (2021) Variable metric techniques for forward-backward methods in imaging. *Journal of Computational and Applied Mathematics* 385:113,192. <https://doi.org/https://doi.org/10.1016/j.cam.2020.113192>

Borkar VS (2008) *Stochastic approximation*. Cambridge University Press, Cambridge; Hindustan Book Agency, New Delhi, <https://doi.org/https://doi.org/10.1007/978-93-86279-38-5>, a



dynamical systems viewpoint

- Brown L (1986) Fundamentals of statistical exponential families : with applications in statistical decision theory. Lecture notes-monograph series Fundamentals of statistical exponential families, Institute of Mathematical Statistics, <https://doi.org/10.1214/lnms/1215466757>
- Cappé O, Moulines E (2009) On-line Expectation Maximization algorithm for latent data models. *J Roy Stat Soc B Met* 71(3):593–613. <https://doi.org/https://doi.org/10.1111/j.1467-9868.2009.00698.x>
- Celex G, Diebolt J (1985) The SEM algorithm: a probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. *Computational Statistics Quarterly* 2:73–82
- Chen HG, Rockafellar R (1997) Convergence rates in forward-backward splitting. *SIAM J Optim* 7:421–444. <https://doi.org/https://doi.org/10.1137/S1052623495290179>
- Chen J, Zhu J, Teh Y, et al (2018) Stochastic Expectation Maximization with Variance Reduction. In: Bengio S, Wallach H, Larochelle H, et al (eds) *Advances in Neural Information Processing Systems* 31. Curran Associates, Inc., p 7967–7977, <https://doi.org/10.5555/3327757.3327893>
- Chen X, Liu S, Sun R, et al (2019) On the convergence of a class of adam-type algorithms for non-convex optimization. In: *International Conference on Learning Representations*
- Choi H, Hobert J (2013) The Polya-Gamma Gibbs sampler for Bayesian logistic regression is uniformly ergodic. *Electronic Journal of Statistics* 7(none):2054 – 2064. <https://doi.org/10.1214/13-EJS837>
- Chouzenoux E, Pesquet JC, Repetti A (2014) Variable Metric Forward-Backward Algorithm for Minimizing the Sum of a Differentiable Function and a Convex Function. *Journal of Optimization Theory and Applications* 162(1):107–132. <https://doi.org/10.1007/s10957-013-0465-7>
- Combettes P, Pesquet J (2011) Proximal Splitting Methods in Signal Processing. In: Bauschke HH, Burachik RS, Combettes PL, et al (eds) *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*. Springer Optimization and Its Applications, Springer, p 185–212, <https://doi.org/10.1007/978-1-4419-9569-8>
- Combettes P, Vũ B (2014) Variable metric forward-backward splitting with applications to monotone inclusions in duality. *Optimization* 63(9):1289–1318. <https://doi.org/10.1080/02331934.2012.733883>
- Combettes PL, Wajs VR (2005) Signal Recovery by Proximal Forward-Backward Splitting. *Multiscale Modeling & Simulation* 4(4):1168–1200. <https://doi.org/10.1137/050626090>
- Defazio A, Bach F, Lacoste-Julien S (2014) Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In: *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1*. MIT Press, Cambridge, MA, USA, NIPS'14, p 1646–1654
- Delyon B, Lavielle M, Moulines E (1999) Convergence of a Stochastic Approximation version of the EM algorithm. *Ann Statist* 27(1):94–128. <https://doi.org/10.1214/aos/1018031103>
- Dempster A, Laird N, Rubin D (1977) Maximum Likelihood from Incomplete Data via the EM Algorithm. *J Roy Stat Soc B Met* 39(1):1–38
- Devroye L (1986) Non-Uniform Random Variate Generation (originally published with. Springer-Verlag, <https://doi.org/https://doi.org/10.1007/978-1-4613-8643-8>
- Eicke B (1992) Iteration methods for convexly constrained ill-posed problems in hilbert space. *Numerical Functional Analysis and Optimization* 13(5-6):413–429. <https://doi.org/10.1080/01630569208816489>
- Everitt B (1984) An introduction to latent variable models. Chapman and Hall London ; New York, <https://doi.org/https://doi.org/10.1007/978-94-009-5564-6>

- Fang C, Li C, Lin Z, et al (2018) Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In: Bengio S, Wallach H, Larochelle H, et al (eds) *Advances in Neural Information Processing Systems*, vol 31. Curran Associates, Inc.
- Fort G, Moulines E (2003) Convergence of the Monte Carlo Expectation Maximization for curved exponential families. *Ann Statist* 31(4):1220–1259
- Fort G, Moulines E (2021) The Perturbed Prox-Preconditioned Spider Algorithm: Non-Asymptotic Convergence Bounds. In: 2021 IEEE Statistical Signal Processing Workshop (SSP), pp 96–100, <https://doi.org/10.1109/SSP49050.2021.9513846>
- Fort G, Moulines E, Priouret P (2011) Convergence of adaptive and interacting markov chain monte carlo algorithms. *Ann Statist* 39(6):3262–3289
- Fort G, Risser L, Atchadé Y, et al (2018) Stochastic fista algorithms: So fast ? In: 2018 IEEE Statistical Signal Processing Workshop (SSP), pp 796–800, <https://doi.org/10.1109/SSP.2018.8450740>
- Fort G, Moulines E, Wai HT (2020) A Stochastic Path-Integrated Differential Estimator Expectation Maximization Algorithm. In: *Proceedings of the 34th International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA, NIPS'20
- Fort G, Gach P, Moulines E (2021a) Fast Incremental Expectation Maximization for finite-sum optimization: nonasymptotic convergence. *Stat Comput* 31(4):48. <https://doi.org/10.1007/s11222-021-10023-9>
- Fort G, Moulines E, Wai HT (2021b) Geom-Spider-EM: Faster Variance Reduced Stochastic Expectation Maximization for Nonconvex Finite-Sum Optimization. In: *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp 3135–3139, <https://doi.org/10.1109/ICASSP39728.2021.9414271>
- Ghadimi S, Lan G (2013) Stochastic First- and Zeroth-Order Methods for Nonconvex Stochastic Programming. *SIAM Journal on Optimization* 23(4):2341–2368. <https://doi.org/10.1137/120880811>
- Ghadimi S, Lan G, Zhang H (2016) Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Math Program* 155(1-2):267–305. <https://doi.org/10.1007/s10107-014-0846-1>
- Gower R, Goldfarb D, Richtarik P (2016) Stochastic Block BFGS: Squeezing More Curvature out of Data. In: Balcan MF, Weinberger KQ (eds) *Proceedings of The 33rd International Conference on Machine Learning, Proceedings of Machine Learning Research*, vol 48. PMLR, New York, New York, USA, pp 1869–1878
- Hiriart-Urruty JB, Lemaréchal C (1996) *Convex Analysis and Minimization Algorithms*. Springer Verlag, Heidelberg, <https://doi.org/https://doi.org/10.1007/978-3-662-02796-7>, two volumes - 2nd printing
- Horváth S, Lei L, Richtárik P, et al (2022) Adaptivity of Stochastic Gradient Methods for Nonconvex Optimization. *SIAM Journal on Mathematics of Data Science* 4(2):634–648. <https://doi.org/10.1137/21M1394308>
- Johnson R, Zhang T (2013) Accelerating stochastic gradient descent using predictive variance reduction. In: Burges C, Bottou L, Welling M, et al (eds) *Advances in Neural Information Processing Systems*, vol 26. Curran Associates, Inc.
- Karimi B, Wai HT, Moulines E, et al (2019) On the Global Convergence of (Fast) Incremental Expectation Maximization Methods. In: Wallach H, Larochelle H, Beygelzimer A, et al (eds) *Advances in Neural Information Processing Systems* 32. Curran Associates, Inc., p 2837–2847
- Karimi H, Nutini J, Schmidt M (2016) Linear Convergence of Gradient and Proximal-Gradient Methods Under the Polyak-Łojasiewicz Condition. In: Frascconi P, Landwehr N, Manco G, et al (eds) *Machine Learning and Knowledge Discovery in*

- Databases. Springer International Publishing, pp 795–811
- Kolte R, Erdogdu M, Ozgur A (2015) Accelerating svrg via second-order information. In: *Advances in Neural Information Processing Systems - Workshop OptML*, pp 1–5
- Lan G (2020) *First-order and Stochastic Optimization Methods for Machine Learning*. Springer Series in the Data Sciences, Springer International Publishing, <https://doi.org/https://doi.org/10.1007/978-3-030-39568-1>
- Lee JD, Sun Y, Saunders MA (2014) Proximal Newton-Type Methods for Minimizing Composite Functions. *SIAM Journal on Optimization* 24(3):1420–1443. <https://doi.org/10.1137/130921428>
- Li Z, Li J (2018) A Simple Proximal Stochastic Gradient Method for Nonsmooth Nonconvex Optimization. In: Bengio S, Wallach H, Larochelle H, et al (eds) *Advances in Neural Information Processing Systems*, vol 31. Curran Associates, Inc.
- Li Z, Bao H, Zhang X, et al (2021) PAGE: A Simple and Optimal Probabilistic Gradient Estimator for Nonconvex Optimization. In: Meila M, Zhang T (eds) *Proceedings of the 38th International Conference on Machine Learning*, *Proceedings of Machine Learning Research*, vol 139. PMLR, pp 6286–6295
- McLachlan G, Krishnan T (2008) *The EM algorithm and extensions*, 2nd edn. Wiley series in probability and statistics, Wiley, <https://doi.org/10.1002/9780470191613>
- Metel M, Takeda A (2021) Stochastic proximal methods for non-smooth non-convex constrained sparse optimization. *Journal of Machine Learning Research* 22(115):1–36
- Meyn S, Tweedie R (1993) *Markov Chains and Stochastic Stability*. Springer-Verlag, London, <https://doi.org/https://doi.org/10.1007/978-1-4471-3267-7>
- Moreau J (1965) Proximité et dualité dans un espace hilbertien. *Bulletin de la Société Mathématique de France* 93:273–299. <https://doi.org/10.24033/bsmf.1625>
- Moritz P, Nishihara R, Jordan M (2016) A Linearly-Convergent Stochastic L-BFGS Algorithm. In: Gretton A, Robert CC (eds) *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, *Proceedings of Machine Learning Research*, vol 51. PMLR, pp 249–258
- Neal RM, Hinton GE (1998) A View of the EM Algorithm that Justifies Incremental, Sparse, and other Variants. In: Jordan MI (ed) *Learning in Graphical Models*. Springer Netherlands, Dordrecht, p 355–368, [https://doi.org/10.1007/978-94-011-5014-9\\_12](https://doi.org/10.1007/978-94-011-5014-9_12)
- Ng SK, McLachlan GJ (2003) On the choice of the number of blocks with the incremental EM algorithm for the fitting of normal mixtures. *Stat Comput* 13(1):45–55. <https://doi.org/10.1023/A:1021987710829>
- Nguyen L, Liu J, Scheinberg K, et al (2017) SARAH: A novel method for machine learning problems using stochastic recursive gradient. In: Precup D, Teh YW (eds) *Proceedings of the 34th International Conference on Machine Learning*, pp 2613–2621
- Nhan HP, Lam MN, Dzung TP, et al (2020) Prox-SARAH: An Efficient Algorithmic Framework for Stochastic Composite Nonconvex Optimization. *Journal of Machine Learning Research* 21(110):1–48
- Park Y, Dhar S, Boyd S, et al (2019) Variable Metric Proximal Gradient Method with Diagonal Barzilai-Borwein Stepsize. <https://doi.org/10.48550/ARXIV.1910.07056>
- Polson NG, Scott J, Windle J (2013) Bayesian Inference for Logistic Models Using P’olya–Gamma Latent Variables. *Journal of the American Statistical Association* 108(504):1339–1349. <https://doi.org/https://doi.org/10.1080/01621459.2013.829001>

- Reddi SJ, Hefny A, Sra S, et al (2016) Stochastic Variance Reduction for Nonconvex Optimization. In: Balcan MF, Weinberger KQ (eds) Proceedings of The 33rd International Conference on Machine Learning, Proceedings of Machine Learning Research, vol 48. PMLR, New York, New York, USA, pp 314–323
- Repetti A, Wiaux Y (2021) Variable metric forward-backward algorithm for composite minimization problems. *SIAM J Optim* 31:1215–1241. <https://doi.org/https://doi.org/10.1137/19M1277552>
- Repetti A, Chouzenoux E, Pesquet JC (2014) A preconditioned forward-backward approach with application to large-scale nonconvex spectral unmixing problems. In: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 1498–1502, <https://doi.org/10.1109/ICASSP.2014.6853847>
- Robert C, Casella G (2004) Monte Carlo statistical methods. Springer Verlag, <https://doi.org/https://doi.org/10.1007/978-1-4757-4145-2>
- Wang Z, Ji K, Zhou Y, et al (2019) Spiderboost and momentum: Faster variance reduction algorithms. In: Wallach HM, Larochelle H, Beygelzimer A, et al (eds) Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp 2403–2413
- Wei G, Tanner M (1990) A Monte Carlo Implementation of the EM Algorithm and the Poor Man’s Data Augmentation Algorithms. *J Am Stat Assoc* 85(411):699–704. <https://doi.org/10.1080/01621459.1990.10474930>
- Wu C (1983) On the Convergence Properties of the EM Algorithm. *Ann Statist* 11(1):95–103. <https://doi.org/10.1214/aos/1176346060>
- Yun J, Lozano AC, Yang E (2021) Adaptive proximal gradient methods for structured neural networks. In: Ranzato M, Beygelzimer A, Dauphin Y, et al (eds) Advances in Neural Information Processing Systems, vol 34. Curran Associates, Inc., pp 24,365–24,378
- Zhang J, Xiao L (2019) A stochastic composite gradient method with incremental variance reduction. In: Wallach H, Larochelle H, Beygelzimer A, et al (eds) Advances in Neural Information Processing Systems, vol 32. Curran Associates, Inc.
- Zhang Q, Huang F, Deng C, et al (2022) Faster stochastic quasi-newton methods. *IEEE Transactions on Neural Networks and Learning Systems* 33(9):4388–4397. <https://doi.org/10.1109/TNNLS.2021.3056947>
- Zhou D, Xu P, Gu Q (2020) Stochastic Nested Variance Reduction for Nonconvex Optimization. *Journal of Machine Learning Research* 21(103):1–63