

# Parallel Random Block-Coordinate Forward-Backward Algorithm: A Unified Convergence Analysis

Saverio Salzo\* and Silvia Villa†

## Abstract

We study the block-coordinate forward-backward algorithm in which the blocks are updated in a random and possibly parallel manner, according to arbitrary probabilities. The algorithm allows different stepsizes along the block-coordinates to fully exploit the smoothness properties of the objective function. In the convex case and in an infinite dimensional setting, we establish almost sure weak convergence of the iterates and the asymptotic rate  $o(1/n)$  for the mean of the function values. We derive linear rates under strong convexity and error bound conditions. Our analysis is based on an abstract convergence principle for stochastic descent algorithms which allows to extend and simplify existing results.

**Keywords.** Convex optimization, parallel algorithms, random block-coordinate descent, arbitrary sampling, error bounds, stochastic quasi-Fejér sequences, forward-backward algorithm, convergence rates.

**AMS Mathematics Subject Classification:** 65K05, 90C25, 90C06, 49M27

## 1 Introduction and problem setting

Random block-coordinate descent algorithms are nowadays among the methods of choice for solving large scale optimization problems [28, 34, 41]. Indeed, they have low complexity and low memory requirements and, additionally, they are amenable for distributed and parallel implementations [32, 34]. In the last decade a number of works have appeared on the topic which address several aspects, that is: the way the block sampling is performed, the composite structure, the partial separability, and the smoothness/geometrical properties of the objective function, accelerations, and iteration complexity [4, 5, 13, 21, 23, 24, 28, 30, 31, 33, 34, 38].

In this work we consider the following optimization problem

$$\underset{x \in H}{\text{minimize}} \quad f(x) + h(x), \quad h(x) = \sum_{i=1}^m h_i(x_i), \quad (1.1)$$

where  $H$  is the direct sum of  $m$  separable real Hilbert spaces  $(H_i)_{1 \leq i \leq m}$ , that is,

$$H = \bigoplus_{i=1}^m H_i, \quad (\forall x = (x_i)_{1 \leq i \leq m}, y = (y_i)_{1 \leq i \leq m} \in H) \quad \langle x, y \rangle = \sum_{i=1}^m \langle x_i, y_i \rangle,$$

and the following assumptions hold:

\*Istituto Italiano di Tecnologia, Via Melen, 83, 16152 Genova, Italy (saverio.salzo@iit.it).

†Università degli Studi di Genova, Via Dodecaneso, 35, 16146 Genova, Italy (silvia.villa@unige.it). Supported by the H2020-MSCA-RISE project NoMADS-GA No. 777826 and by Gruppo Nazionale per l'Analisi Matematica, la Probabilità e le loro Applicazioni (GNAMPA) of the Istituto Nazionale di Alta Matematica (INdAM).

H1  $f: H \rightarrow \mathbb{R}$  is convex and differentiable,

H2 for every  $i = 1, \dots, m$ ,  $h_i: H_i \rightarrow ]-\infty, +\infty]$  is proper, convex, and lower semicontinuous.

The objective of this study is a stochastic algorithm, called *parallel random block-coordinate forward-backward algorithm*, that depends on a random variable  $\varepsilon$  satisfying the following hypothesis

H3  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$  is a random variable with values in  $\{0, 1\}^m$  such that, for every  $i \in \{1, \dots, m\}$ ,  $p_i := P(\varepsilon_i = 1) > 0$  and  $P(\varepsilon = (0, \dots, 0)) = 0$ .

**Algorithm 1.1.** Let  $(\varepsilon^n)_{n \in \mathbb{N}} = (\varepsilon_1^n, \dots, \varepsilon_m^n)_{n \in \mathbb{N}}$  be a sequence of independent copies of  $\varepsilon$ . Let  $(\gamma_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  and  $x^0 = (x_1^0, \dots, x_m^0) \equiv x^0 \in \text{dom } h$  be a constant random variable. Iterate

$$\begin{array}{l} \text{for } n = 0, 1, \dots \\ \quad \left[ \begin{array}{l} \text{for } i = 1, \dots, m \\ \quad \left[ x_i^{n+1} = x_i^n + \varepsilon_i^n (\text{prox}_{\gamma_i h_i}(x_i^n - \gamma_i \nabla_i f(x^n)) - x_i^n). \end{array} \right. \end{array} \right. \quad (1.2)$$

For every  $n \in \mathbb{N}$ , we denote by  $\mathfrak{E}_n$  the sigma-algebra generated by  $\varepsilon^0, \dots, \varepsilon^n$ .

In Algorithm 1.1, the role of the random variable  $\varepsilon^n$  is to select, at iteration  $n$ , the blocks to update in parallel (those indexed in  $\{i \in \{1, \dots, m\} \mid \varepsilon_i^n = 1\}$ ). When all block-coordinates are simultaneously updated at each iteration, Algorithm 1.1 reduces to the (deterministic) forward-backward algorithm, which converges only if the stepsizes are appropriately set. More specifically, if  $\nabla f$  is  $L$ -Lipschitz continuous, then convergence is ensured if the stepsizes  $\gamma_i$  are all equal and *strictly less* than  $2/L$  [6, 7]. This fact is proved by using the so called *descent lemma*, i.e.,

$$(\forall x \in H)(\forall v \in H) \quad f(x + v) \leq f(x) + \langle \nabla f(x), v \rangle + \frac{L}{2} \|v\|^2. \quad (1.3)$$

Indeed, (1.3) is itself an assumption concerning the smoothness of  $f$ , since it is well-known to be equivalent to the Lipschitz continuity of the gradient of  $f$  [1, Theorem 18.15]. By contrast, when the block-coordinates are updated one by one in a serial manner, it is desirable to allow moving along the block-coordinates with different stepsizes, depending on the Lipschitz constants of the partial gradients of  $f$  across the block-coordinates [2, 28]. So, in this case it is more appropriate to assume that a descent lemma holds on each block-coordinate subspace individually, that is,

$$(\forall i = 1, \dots, m)(\forall x \in H)(\forall v_i \in H_i) \quad f(x + J_i v_i) \leq f(x) + \langle \nabla_i f(x), v_i \rangle + \frac{L_i}{2} \|v_i\|^2, \quad (1.4)$$

where  $J_i v_i = (0, \dots, 0, v_i, 0, \dots, 0)$  ( $v_i$  occurring at the  $i$ -th position), for some positive constants  $L_i$ 's. In the setting of Algorithm 1.1, multiple block-coordinates may be updated in parallel at each iteration, according to the random sampling  $\varepsilon$ . Therefore, it is reasonable to assume that there exists  $(\nu_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  so that one of the *generalized smoothness conditions* below holds

$$\text{S1 } (\forall x, v \in H) \quad \mathbb{E}[f(x + \varepsilon \odot v)] \leq f(x) + \mathbb{E}[\langle \nabla f(x), \varepsilon \odot v \rangle] + \frac{1}{2} \sum_{i=1}^m p_i \nu_i \|v_i\|^2,$$

$$\text{S2 } (\forall x, v \in H) \quad f(x + \varepsilon \odot v) \leq f(x) + \langle \nabla f(x), \varepsilon \odot v \rangle + \frac{1}{2} \sum_{i=1}^m \nu_i \varepsilon_i \|v_i\|^2 \quad \text{P a.s.,}$$

where  $\varepsilon \odot \mathbf{v} = (\varepsilon_i \mathbf{v}_i)_{1 \leq i \leq m} \in \mathbb{H}$ . Conditions **S1** and **S2** can be interpreted as descent lemmas on random block-coordinate subspaces, depending on the chosen random sampling of the block-coordinates. They reduce to (1.4), with  $\nu_i = L_i$ , if the sampling  $\varepsilon$  selects only one block at a time almost surely (see Section 3.2). We call  $(\nu_i)_{1 \leq i \leq m}$  the *smoothness parameters* of  $f$ . Then, similarly to the deterministic case, we will adopt the following stepsize rule

$$(\forall i \in \{1, \dots, m\}) \quad \gamma_i < \frac{2}{\nu_i}. \quad (1.5)$$

Another smoothness condition suitable for Algorithm 1.1, which was considered in [23], is

$$\mathbf{S3} \quad (\forall \mathbf{x}, \mathbf{v} \in \mathbb{H}) \quad f(\mathbf{x} + \mathbf{v}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \sum_{i=1}^m \nu_i \|\mathbf{v}_i\|^2.$$

Note that, **S3**  $\Rightarrow$  **S2**  $\Rightarrow$  **S1**, which in turn implies the Lipschitz continuity of the gradient of  $f$  (see Theorem 3.1(iv)). So, possibly with different values of the  $\nu_i$ 's, the above conditions are all equivalent. The point is that in the parallel setting (where multiple blocks are updated in parallel at each iteration), **S1** may be fulfilled with values of  $\nu_i$  that are much smaller than those related to the other two conditions, ultimately allowing to significantly increase the stepsizes and hence speeding up the convergence. Moreover, **S1** makes parallelization particularly effective on problems with a sparse structure and superior to the serial strategy (which updates a single block per iteration). See the discussion after Theorem 4.9. The critical role played by assumption **S1** in the analysis of parallel randomized block-coordinate descent methods was pointed out in [31, 34, 35, 38]. There, it was called *expected separable overapproximation* (ESO) inequality. Condition **S2** is new and serves to guarantee that Algorithm 1.1 is almost surely descending (Proposition 4.7), which is a property that is especially relevant when error bound conditions hold (see Section 4.3). Note that in [34] the issue of monotonicity of the algorithm was addressed for each sampling separately without any general guidance. Finally, we stress that, except for [4, 5] (which study the convergence of the iterates only), in all previous works the stepsizes  $\gamma_i$ 's are set equal to  $1/\nu_i$ . This is an unnecessary limitation that we remove, so to match the standard stepsize rule of the forward-backward algorithm [6, 7].

**Remark 1.2.** For every  $i = 1, \dots, m$ , the canonical embedding of  $\mathbb{H}_i$  into  $\mathbb{H}$  is the operator  $J_i: \mathbb{H}_i \rightarrow \mathbb{H}$ ,  $\mathbf{x} \mapsto (0, \dots, 0, \mathbf{x}, 0, \dots, 0)$ , where  $\mathbf{x}$  occurs in the  $i$ -th position. Then Algorithm 1.1 can be written as

$$x^{n+1} = x^n + \sum_{i=1}^m \varepsilon_i^n J_i (\text{prox}_{\gamma_i h_i}(x_i^n - \gamma_i \nabla_i f(x^n)) - x_i^n).$$

## 1.1 Main contributions and comparison to previous work

In the following we summarize the main contributions of this paper, where, for the sake of brevity, we set  $F = f + h$ . We assume that **H1–H3** are satisfied and that **S1** is met. Then, the following hold.

- Algorithm 1.1 is descending in expectation and  $E[F(x^n)] - \inf F \rightarrow 0$ , even if the infimum is not attained. If  $\text{argmin} F \neq \emptyset$ , then  $E[F(x^n)] - \inf F = o(1/n)$ . In addition, a nonasymptotic bound for  $E[F(x^n)] - \inf F$  of order  $O(1/n)$  holds. Finally, there exists a random variable  $x_*$  with values in  $\text{argmin} F$  such that  $x^n \rightarrow x_*$  P-a.s. See Theorem 4.9.
- If  $F$  is strongly convex or satisfies an error bound condition of Luo-Tseng type (see condition **EB**), then the iterates as well as the corresponding function values generated by Algorithm 1.1, converge linearly in expectation. See Theorem 4.10, Theorem 4.16, and Theorem 4.19.

Our results advance the state-of-the-art in the study of random block-coordinate descent methods under several aspects. We comment on this below. 1) While convergence of the function values has been intensively studied in the related literature (see e.g., [16, 21, 23, 28, 29, 34, 35, 38]), surprisingly, in a convex setting, convergence of the iterates has been investigated only recently in [4], but with stepsizes set according to the global Lipschitz constant of  $\nabla f$ . See also [12] which addresses the convergence of the iterates in the framework of primal-dual algorithms with a serial and uniform block sampling. We improve the existing results, since we show *convergence of the iterates* for Algorithm 1.1 in an *infinite dimensional setting* even when the stepsizes are chosen according to the condition S1, which can incorporate the block Lipschitz constants of the gradient of  $f$  and is at the basis of the effectiveness of the parallel block-coordinatewise approach. 2) The worst case *asymptotic rate*  $o(1/n)$  for the mean of the function values is new in the setting of stochastic algorithms. 3) Our analysis spotlights an abstract convergence principle for stochastic descent algorithms (Theorem 4.1) which is essentially a special form of the stochastic quasi-Fejér monotonicity property, involving also the values of the objective functions. This principle, previously investigated in a deterministic setting in [36], allows to prove in a unified way both the almost sure convergence of the iterates and rates of convergence for the mean of the function values. 4) As a by-product of the above analysis we single out an inequality (Proposition 4.4) which is pivotal for studying *the convergence under error bound conditions*, improving the results and simplifying the analysis in [23]. 5) We allow for *parallel and arbitrary sampling of the blocks* in a composite setting. The benefit of such sampling in terms of convergence rate have been first investigated in [35] for a strongly convex and smooth objective function. In [30] a composite objective optimization problem was analyzed but for a slightly different algorithm. The rest of the studies deal either with parallel uniform sampling of the blocks [34], or with the case where a single block is updated at each iteration [21, 28]. 6) We also allow for *stepsizes larger than those considered in literature* [16, 21, 23, 28, 29, 34, 35, 38], since we can let the stepsizes go beyond  $1/\nu_i$  and be arbitrarily close to  $2/\nu_i$ , matching the standard rule for the forward-backward algorithm. This provides additional flexibility to the algorithm. Indeed, in the strongly convex case we show that the optimal stepsizes are strictly larger than  $1/\nu_i$ .

The rest of the paper is organized as follows. In Section 2 we give notation and basic facts. Section 3 shows how to determine the smoothness parameters  $\nu_i$  when  $f$  features a partially separable structure. In Section 4 we carry out the convergence analysis and give the related theorems. Finally, Section 5 shows three applications and Section 6 provides some numerical experiments.

## 2 Notation and background

**Notation.** We define  $\mathbb{R}_+ = [0, +\infty[$ ,  $\mathbb{R}_{++} = ]0, +\infty[$ , for every integer  $s \geq 1$ ,  $[s] = \{1, \dots, s\}$ , and for every  $a \in \mathbb{R}^s$ ,  $\text{spt}(a) = \{i \in [s] \mid a_i \neq 0\}$ . Scalar products and norms in Hilbert spaces are denoted by  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|$  respectively. If  $U: H \rightarrow G$  is a bounded linear operator between real Hilbert spaces,  $U^\top: G \rightarrow H$  is its transpose operator, that is, the one satisfying  $\langle Ux, y \rangle = \langle x, U^\top y \rangle$ , for every  $(x, y) \in H \times G$ . Let  $(H_i)_{1 \leq i \leq m}$  be  $m$  separable real Hilbert spaces and let  $H = \bigoplus_{i=1}^m H_i$  be their direct sum. For every  $v \in H$  and  $\epsilon \in \{0, 1\}^m$  we set  $\epsilon \odot v = (\epsilon_i v_i)_{1 \leq i \leq m} \in H$ . We will consider random variables with underlying probability space  $(\Omega, \mathfrak{A}, P)$  taking values in  $H_i$  or  $H$ . We use the default font for random variables and sans serif font for their realizations. The expected value operator is denoted by  $E$ . A copy of a random variable is random variable having the same distribution of the given one. Let  $(w_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$ . The direct sum operator  $W = \bigoplus_{i=1}^m w_i \text{Id}_i$ , where  $\text{Id}_i$  is the identity operator on  $H_i$ , is the positive bounded linear operator on  $H$  acting as  $x = (x_i)_{1 \leq i \leq m} \mapsto (w_i x_i)_{1 \leq i \leq m}$ .

$W$  defines an equivalent inner product on  $H$

$$(\forall x \in H)(\forall y \in H) \quad \langle x, y \rangle_W = \langle Wx, y \rangle = \sum_{i=1}^m w_i \langle x_i, y_i \rangle,$$

which gives the norm  $\|x\|_W^2 = \sum_{i=1}^m w_i \|x_i\|^2$ . If  $S \subset H$  and  $x \in H$ , we set  $\text{dist}_W(x, S) = \inf_{z \in S} \|x - z\|_W$ . Let  $\varphi: H \rightarrow ]-\infty, +\infty]$  be proper, convex, and lower semicontinuous. The domain of  $\varphi$  is  $\text{dom } \varphi = \{x \in H \mid \varphi(x) < +\infty\}$  and the set of minimizers of  $\varphi$  is  $\text{argmin } \varphi = \{x \in H \mid \varphi(x) = \inf \varphi\}$ . The subdifferential of  $\varphi$  in the metric  $\langle \cdot, \cdot \rangle_W$  is the multivalued operator

$$\partial^W \varphi: H \rightarrow 2^H, \quad x \mapsto \partial^W \varphi(x) = \{u \in H \mid (\forall y \in H) \varphi(y) \geq \varphi(x) + \langle u, y - x \rangle_W\}.$$

In case  $W = \text{Id}$ , it is simply denoted by  $\partial \varphi$ . Clearly  $\partial^W \varphi = W^{-1} \partial \varphi$ . If the function  $\varphi: H \rightarrow \mathbb{R}$  is differentiable, then, for every  $x \in H$ ,  $\partial^W \varphi(x) = \{\nabla^W \varphi(x)\}$  and for all  $v \in H$ ,  $\langle \nabla^W \varphi(x), v \rangle_W = \langle \nabla \varphi(x), v \rangle$ . The proximity operator of  $\varphi$  in the metric  $\langle \cdot, \cdot \rangle_W$  is defined as

$$\text{prox}_\varphi^W: H \rightarrow H, \quad \text{prox}_\varphi^W(x) = \underset{z \in H}{\text{argmin}} \varphi(z) + \frac{1}{2} \|x - z\|_W^2.$$

Referring to the functions in (1.1), we denote by  $\mu_{\Gamma^{-1}}$  and  $\sigma_{\Gamma^{-1}}$  the moduli of strong convexity of  $f$  and  $h$  respectively, in the norm  $\|\cdot\|_{\Gamma^{-1}}$ , where  $\Gamma = \bigoplus_{i=1}^m \gamma_i \text{Id}_i$  and the  $\gamma_i$ 's are the stepsizes occurring in Algorithm 1.1. This means that  $\mu_{\Gamma^{-1}}, \sigma_{\Gamma^{-1}} \in \mathbb{R}_+$  and that, for every  $x, y \in H$ ,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu_{\Gamma^{-1}}}{2} \sum_{i=1}^m \frac{1}{\gamma_i} \|y_i - x_i\|^2, \quad (2.1)$$

$$(\forall v \in \partial h(x)) \quad h(y) \geq h(x) + \langle v, y - x \rangle + \frac{\sigma_{\Gamma^{-1}}}{2} \sum_{i=1}^m \frac{1}{\gamma_i} \|y_i - x_i\|^2. \quad (2.2)$$

Note that, since  $h$  is separable, by taking  $y = x + J_i(y_i - x_i)$  in (2.2), we have

$$(\forall i \in [m]) \quad h_i(y_i) \geq h_i(x_i) + \langle v_i, y_i - x_i \rangle + \frac{\sigma_{\Gamma^{-1}}}{2} \frac{1}{\gamma_i} \|y_i - x_i\|^2. \quad (2.3)$$

**Remark 2.1.** If S1 is satisfied, the  $\gamma_i$ 's are chosen as in (1.5), and  $\delta = \max_{1 \leq i \leq m} \gamma_i \nu_i$  (according to the convergence theorems), then we have

$$\mu_{\Gamma^{-1}} \leq \min_{1 \leq i \leq m} \gamma_i \nu_i \leq \delta < 2. \quad (2.4)$$

Indeed, let  $x \in H$ ,  $i \in [m]$  and  $v_i \in H_i$ ,  $v_i \neq 0$ . It follows from (2.1) with  $y = x + \varepsilon \odot J_i v_i$  (where  $J_i$  is defined in Remark 1.2) and S1 that

$$\begin{aligned} f(x) + E[\langle \nabla f(x), \varepsilon \odot J_i v_i \rangle] + \frac{\mu_{\Gamma^{-1}}}{2} \frac{p_i}{\gamma_i} \|v_i\|^2 &\leq E[f(x + \varepsilon \odot J_i v_i)] \\ &\leq f(x) + E[\langle \nabla f(x), \varepsilon \odot J_i v_i \rangle] + \frac{1}{2} p_i \nu_i \|v_i\|^2. \end{aligned}$$

Thus, (2.4) follows.

**Fact 2.2** ([10, Example 5.1.5]). *Let  $\zeta_1$  and  $\zeta_2$  be independent random variables with values in the measurable spaces  $\mathcal{Z}_1$  and  $\mathcal{Z}_2$  respectively. Let  $\varphi: \mathcal{Z}_1 \times \mathcal{Z}_2 \rightarrow \mathbb{R}$  be measurable and suppose that  $E[|\varphi(\zeta_1, \zeta_2)|] < +\infty$ . Then  $E[\varphi(\zeta_1, \zeta_2) \mid \zeta_1] = \psi(\zeta_1)$ , where for all  $z_1 \in \mathcal{Z}_1$ ,  $\psi(z_1) = E[\varphi(z_1, \zeta_2)]$ .*

**Fact 2.3.** *Let  $\varepsilon$  be a random variable with values in  $\{0, 1\}^m$  and, for all  $i \in [m]$ ,  $p_i = P(\varepsilon_i = 1)$ . Then  $E[\varepsilon_i] = p_i$  and, for every  $v = (v_i)_{1 \leq i \leq m} \in \mathbb{R}^m$ ,  $E[\langle \varepsilon, v \rangle] = \sum_{i=1}^m p_i v_i$ .*

**Fact 2.4** ([17]). *Let  $(a_n)_{n \in \mathbb{N}}$  be a decreasing sequence in  $\mathbb{R}_+$ . If  $\sum_{n=0}^{+\infty} a_n < +\infty$ , then, for every  $n \in \mathbb{N}$ ,  $a_n \leq (1/(n+1)) \sum_{n=0}^{+\infty} a_n$  and  $a_n = o(1/(n+1))$ .*

### 3 Determining the smoothness parameters

In this section we provide few scenarios for which the relaxed smoothness conditions **S1** and **S2** can be fully exploited, attaining tight values for the  $\nu_i$ 's. This ultimately allows to take larger stepsizes and improves rates of convergence. In [31, 38] an extensive analysis of cases in which **S1** is satisfied is presented.

#### 3.1 General estimates.

We consider the following setting.

H4 The function  $f: H \rightarrow \mathbb{R}$  is such that

$$(\forall x \in H) \quad f(x) = \sum_{k=1}^p g_k \left( \sum_{i=1}^m U_{k,i} x_i \right), \quad (3.1)$$

where, for every  $k = 1, \dots, p$ ,  $g_k: G_k \rightarrow \mathbb{R}$  is a convex differentiable function defined on a real Hilbert space  $G_k$  and, for every  $i \in [m]$ ,  $U_{k,i}: H_i \rightarrow G_k$  is a bounded linear operator. Moreover,  $\bigcup_{k=1}^p I_k \neq \emptyset$ , where, for all  $k = 1, \dots, p$ ,  $I_k = \{i \in [m] \mid U_{k,i} \neq 0\}$ , and  $\eta = \max_{1 \leq k \leq p} \text{card}(I_k)$ .

We will also consider one of the following conditions.

L1 For every  $i = 1, \dots, m$  there exists  $L_i > 0$  such that, for every  $x \in H$ , the function  $\nabla_i f(x_1, \dots, x_{i-1}, \cdot, x_{i+1}, \dots, x_m): H_i \rightarrow H_i$  is  $L_i$ -Lipschitz continuous.

L2 For every  $k = 1, \dots, p$ ,  $\nabla g_k: G_k \rightarrow G_k$  is  $L^{(k)}$ -Lipschitz continuous and for every  $i, j \in [m]$ ,  $i \neq j$ , the ranges of  $U_{k,i}$  and  $U_{k,j}$  are orthogonal.

Assumption **H4** concerns the *partial separability* of the function  $f$ . Depending on the number of the nonzero operators  $U_{k,i}$ ,  $g_k$  might depend only on few block-variables  $x_i$ 's: if  $\eta = 1$ ,  $f$  is fully separable, whereas if  $\eta = m$ ,  $f$  is not separable. Note that **L1** is equivalent to (1.4) and, since  $f$  is convex, implies the global Lipschitz continuity of the gradient of  $f$  (Corollary A.2). So either **L1** or **L2** implies the global Lipschitz smoothness of  $f$ . However, considering the constants  $L_i$ 's or  $L^{(k)}$ 's leads in general to a finer analysis of the smoothness properties of  $f$ , eventually determining parameters  $\nu_i$  that are smaller than the global Lipschitz constant of  $\nabla f$ . Instances of problem (1.1) where  $f$  has the structure shown in **H4**, occur very often in applications. In particular, a prominent example is that of the *Lasso problem* which will be discussed in Section 5.1. The following theorem, which is proved in Appendix B, relates the smoothness parameters  $(\nu_i)_{1 \leq i \leq m}$  to the block Lipschitz constants of the partial gradients of  $f$  and to the Lipschitz constants of the gradients of its components  $g_k$ 's in (3.1), as well as to the distribution of the random variable  $\varepsilon$ .

**Theorem 3.1.** Assume **H3** and **H4** and let  $(\nu_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^{\mathbb{N}}$ . Then the following hold.

(i) **L1**  $\Rightarrow$  **S1** provided that

$$(\forall i \in [m]) \quad \nu_i \geq \beta_{1,i} L_i, \quad \text{where} \quad \beta_{1,i} := \mathbb{E} \left[ \max_{1 \leq k \leq p} \left( \sum_{j \in I_k} \varepsilon_j \right) \mid \varepsilon_i = 1 \right].$$

(ii) **L1**  $\Rightarrow$  **S2** provided that

$$(\forall i \in [m]) \quad \nu_i \geq \beta_2 L_i, \quad \text{where} \quad \beta_2 := \text{ess sup} \left( \max_{1 \leq k \leq p} \left( \sum_{j \in I_k} \varepsilon_j \right) \right).$$

(iii) **L2**  $\Rightarrow$  **S3** provided that

$$(\forall i \in [m]) \quad \nu_i \geq \tilde{L}_i := \left\| \sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \right\|. \quad (3.2)$$

(iv) **S1**  $\Rightarrow$  **L1** with, for all  $i \in [m]$ ,  $L_i = \nu_i$ . In particular, **S1** implies that  $f$  is Lipschitz smooth.

**Remark 3.2.**

- (i) Suppose that in **H4**, for all  $k \in [p]$ ,  $\mathbf{G}_k = \mathbf{H}$ ,  $\mathbf{g}_k$  is  $L^{(k)}$ -Lipschitz smooth, and, for all  $i \in I_k$ ,  $\mathbf{U}_{k,i} = \mathbf{J}_i$ , the canonical embedding of  $\mathbf{H}_i$  into  $\mathbf{H}$  (see Remark 1.2). Then, **L2** holds and, for every  $i \in [m]$ ,  $\tilde{L}_i = \sum_{k|i \in I_k} L^{(k)}$ . Hence, in view of Theorem 3.1(iii), **S3** is met with  $\nu_i = \tilde{L}_i$ . This setting was studied in [23].
- (ii) If  $\nabla f$  is  $L$ -Lipschitz continuous, then **S3** is satisfied with, for every  $i \in [m]$ ,  $\nu_i = L$ . Therefore, we cover the analysis of the random block-coordinate forward-backward algorithm given in [4, 5] which set the stepsizes as  $\gamma_i < 2/L$ .
- (iii) Let, for every  $k \in [p]$ ,  $f_k(x) = \mathbf{g}_k(\sum_{i=1}^m \mathbf{U}_{k,i} x_i)$ . If, for every  $k \in [p]$ , **S1** (resp. **S2**) holds for  $f_k$  with  $(\nu_i^{(k)})_{1 \leq i \leq m}$ , then **S1** (resp. **S2**) holds for  $f$  with  $\nu_i = \sum_{k=1}^p \nu_i^{(k)}$ .
- (iv) Using similar ideas as in the proof of [34, Theorem 12] we show in Appendix B that item (i) in Theorem 3.1 remains true with

$$\beta_{1,i} := \sum_{t=1}^{\eta} t \max_{\substack{1 \leq k \leq p \\ i \in I_k}} \mathbb{P} \left( \sum_{j \in I_k} \varepsilon_j = t \mid \varepsilon_i = 1 \right). \quad (3.3)$$

**Remark 3.3.** Referring to Theorem 3.1, for all  $i \in [m]$ , we have  $1 \leq \beta_{1,i} \leq \beta_2 \leq \min\{\eta, \tau_{\max}\}$ , where

$$\tau_{\max} := \text{ess sup} \left( \sum_{i=1}^m \varepsilon_i \right) = \min \left\{ \tau' \in \mathbb{N} \mid \mathbb{P} \left( \sum_{i=1}^m \varepsilon_i \leq \tau' \right) = 1 \right\}$$

is the maximum number of blocks processed in parallel. Indeed, since  $\mathbb{P}(\varepsilon \equiv 0) = 0$  we have  $\mathbb{P}(\max_{1 \leq i \leq m} \varepsilon_i \geq 1) = 1$ . Moreover, since  $\max_{1 \leq k \leq p} (\sum_{i \in I_k} \varepsilon_i) \geq \max_{1 \leq i \leq m} \varepsilon_i$ , we have  $1 \leq \beta_{1,i}$ . The inequality  $\beta_{1,i} \leq \beta_2$  is immediate, while the last one derives from the following

$$(\forall k \in [p]) \quad \sum_{i \in I_k} \varepsilon_i \leq \min \left\{ \text{card}(I_k), \sum_{i=1}^m \varepsilon_i \right\} \leq \min \left\{ \eta, \sum_{i=1}^m \varepsilon_i \right\} \leq \min\{\eta, \tau_{\max}\}.$$

### 3.2 The smoothness parameters for some special block samplings.

Here we show how to compute (or estimate) the constants  $(\beta_{1,i})_{1 \leq i \leq m}$  and  $\beta_2$  in Theorem 3.1 and Remark 3.2(iv), and the related  $(\nu_i)_{1 \leq i \leq m}$ , in some relevant scenarios, when **H4** and **L1** are satisfied.

*Arbitrary parallel sampling.* It follows from Theorem 3.1(ii) and Remark 3.3 that for an arbitrary (possibly nonuniform) block sampling  $\varepsilon$ , **S2** is satisfied provided that  $\nu_i = \min\{\eta, \tau_{\max}\} L_i$ , for every  $i \in [m]$ . Additionally, if we denote by  $L_i^{(k)}$  the blockwise Lipschitz constants of the gradient of the function  $x \mapsto \mathbf{g}_k(\sum_{i=1}^m \mathbf{U}_{k,i} x_i)$ , then we derive from Remark 3.2(iii) and the above discussion

that **S2** holds with<sup>1</sup>  $\nu_i = \sum_{k|i \in I_k} \min\{\text{card}(I_k), \tau_{\max}\} L_i^{(k)}$ . However, the above estimates are rather conservative and can be improved for special choices of the block sampling as we will show below. We refer to [31, 35] for further results on nonuniform samplings.

*Serial sampling or full separability.* Suppose that  $\tau_{\max} = 1$  or  $\eta = 1$ . Then, Remark 3.3 yields  $\beta_{1,i} = \beta_2 = \min\{\eta, \tau_{\max}\} = 1$ . Moreover, recalling Theorem 3.1(ii)-(iv), this also shows that **S1**, **S2**, and **L1** are indeed equivalent with the same smoothness parameters  $\nu_i = L_i$ . So, conditions **S1** or **S2** find their justification only in the parallel case ( $\tau_{\max} > 1$ ) and when  $f$  is not fully separable ( $\eta > 1$ ).

*Fully Parallel.* If  $\mathbb{P}(\sum_{i=1}^m \varepsilon_i = m) = 1$ , then for every  $i \in [m]$   $p_i = 1$ . This yields a fully parallel (deterministic) algorithm. Moreover, since  $\mathbb{P}(\varepsilon = (1, \dots, 1)) = 1$ , we have  $\beta_{1,i} = \eta = \beta_2$  and hence **S2** holds with  $\nu_i = \eta L_i$ . Actually, also **S3** holds with  $\nu_i = \eta L_i$  (see Corollary A.2(iv)).

*Uniform samplings.* Suppose that  $m > 1$ . The sampling is *uniform* if  $p_i = p_j$ , with  $i \neq j$ . In this case if we denote by  $\bar{\tau}$  the average number of block updates per iteration, we have  $\bar{\tau} = \mathbb{E}[\sum_{i=1}^m \varepsilon_i] = \sum_{i=1}^m p_i$  and hence  $p_i = \bar{\tau}/m$ , for every  $i \in [m]$ . In [34] several types of uniform samplings are studied. In the following we single out two of them. The sampling is said to be *doubly uniform* if any two sets of blocks with the same number of blocks have the same probability to be chosen. In formula, this means that for every  $J_1, J_2 \subset [m]$  such that  $\text{card}(J_1) = \text{card}(J_2)$ ,  $\mathbb{P}(\cap_{i \in J_1} \{\varepsilon_i = 1\}) = \mathbb{P}(\cap_{i \in J_2} \{\varepsilon_i = 1\})$ . For such sampling one directly derives from (3.3) in Remark 3.2(iv) (see Appendix B) that

$$\beta_{1,i} = \beta_1 := 1 + \frac{\eta - 1}{m - 1} \left( \frac{\mathbb{E}[(\sum_{i=1}^m \varepsilon_i)^2]}{\mathbb{E}[\sum_{i=1}^m \varepsilon_i]} - 1 \right). \quad (3.4)$$

A special type of doubly uniform sampling is the  $\tau$ -nice sampling in which  $\sum_{i=1}^m \varepsilon_i = \tau$  P-a.s. for some  $\tau \in [m]$ . In this case (3.4) reduces to

$$\beta_{1,i} = \beta_1 := 1 + \frac{(\eta - 1)(\tau - 1)}{m - 1}. \quad (3.5)$$

Now, according to Remark 3.2(iv), if we set, for every  $i \in [m]$ ,  $\nu_i = \beta_1 L_i$ , then condition **S1** holds. Additionally, if we denote by  $L_i^{(k)}$  the blockwise Lipschitz constants of the gradient of the function  $x \mapsto \mathbf{g}_k(\sum_{i=1}^m \mathbf{U}_{k,i} x_i)$ , then we derive from Remark 3.2(iii) and (3.5) that **S1** is satisfied with  $\nu_i = \sum_{k|i \in I_k} (1 + (\tau - 1)(\text{card}(I_k) - 1)/(m - 1)) L_i^{(k)}$ . This result provides possibly even smaller values for the parameters  $(\nu_i)_{1 \leq i \leq m}$  and was given, in the special setting of Remark 3.2(i), in [13, 38].

## 4 Convergence analysis

In the rest of the paper, referring to Algorithm 1.1, we set

$$\Gamma^{-1} = \bigoplus_{i=1}^m \frac{1}{\gamma_i} \text{Id}_i, \quad (w_i)_{1 \leq i \leq m} = \left( \frac{1}{\gamma_i p_i} \right)_{1 \leq i \leq m}, \quad \mathbf{W} = \bigoplus_{i=1}^m w_i \text{Id}_i, \quad (4.1)$$

where  $\text{Id}_i$  is the identity operator on  $H_i$ , and

$$\bar{x}^{n+1} = (\text{prox}_{\gamma_i h_i}(x_i^n - \gamma_i \nabla_i f(x^n)))_{1 \leq i \leq m}, \quad \Delta^n = x^n - \bar{x}^{n+1}. \quad (4.2)$$

---

<sup>1</sup>If  $i \notin I_k$ , then  $L_i^{(k)} = 0$ .



Then, we have

$$\bar{x}^{n+1} = \text{prox}_h^{\Gamma^{-1}}(x^n - \nabla^{\Gamma^{-1}} f(x^n)), \quad x^{n+1} = x^n + \varepsilon^n \odot (\bar{x}^{n+1} - x^n), \quad (4.3)$$

and, recalling (1.2), that for every  $i \in [m]$  such that  $\varepsilon_i^n = 1$ ,

$$\bar{x}_i^{n+1} = \text{prox}_{\gamma_i h_i}(x_i^n - \gamma_i \nabla_i f(x^n)) = x_i^{n+1}, \quad \Delta_i^n = x_i^n - x_i^{n+1}. \quad (4.4)$$

Note that  $x^n$  and  $\bar{x}^{n+1}$  are functions of the random variables  $\varepsilon^0, \dots, \varepsilon^{n-1}$  only, hence they are both discrete random variables, which are measurable with respect to  $\mathfrak{E}_{n-1}$ .

#### 4.1 An abstract principle for stochastic convergence

We provide an abstract convergence principle for stochastic descent algorithms in the same spirit of [36, Theorem 3.10]. It simultaneously addresses the convergence of the iterates and that of the function values.

**Theorem 4.1.** *Let  $H$  be a separable real Hilbert space with norm  $\|\cdot\|$ . Let  $\Phi: H \rightarrow ]-\infty, +\infty]$  be a proper, lower semicontinuous, and convex function and set  $S_* = \text{argmin } \Phi$  and  $\Phi_* = \inf \Phi$ . Let  $(x^n)_{n \in \mathbb{N}}$  be a sequence of  $H$ -valued random variables such that  $x^0 \equiv x^0 \in \text{dom } \Phi$  and, for every  $n \in \mathbb{N}$ ,  $\Phi(x^n)$  is P-summable. Consider the following conditions*

P1  $(E[\Phi(x^n)])_{n \in \mathbb{N}}$  is decreasing.

P2 There exist a sequence  $(\mathfrak{X}_n)_{n \in \mathbb{N}}$  of sub-sigma algebras of  $\mathfrak{A}$  such that,  $(\forall n \in \mathbb{N}) \mathfrak{X}_n \subset \mathfrak{X}_{n+1}$  and  $x^n$  is  $\mathfrak{X}_n$ -measurable, a sequence  $(\xi_n)_{n \in \mathbb{N}}$  of  $\mathfrak{X}_n$ -measurable real-valued positive random variables such that  $\sum_{n \in \mathbb{N}} E[\xi_n] \leq b < +\infty$ , and  $a > 0$  such that, for every  $x \in \text{dom } \Phi$  and  $n \in \mathbb{N}$ ,

$$E[\|x^{n+1} - x\|^2 | \mathfrak{X}_n] \leq \|x^n - x\|^2 + aE[\Phi(x) - \Phi(x^{n+1}) | \mathfrak{X}_n] + \xi_n \quad \text{P-a.s.} \quad (4.5)$$

P3 There exist  $(y^n)_{n \in \mathbb{N}}$  and  $(v^n)_{n \in \mathbb{N}}$ , sequences of  $H$ -valued random variables, such that  $(\forall n \in \mathbb{N}) v^n \in \partial\Phi(y^n)$ ,  $y^n - x^n \rightarrow 0$ , and  $v^n \rightarrow 0$  P-a.s.

Assume P1 and that  $(\inf_{n \in \mathbb{N}} E[\Phi(x^n)] > -\infty) \Rightarrow$  P2. Then, the following hold.

(i)  $E[\Phi(x^n)] \rightarrow \Phi_*$ .

(ii) Suppose that  $S_* \neq \emptyset$ . Then  $E[\Phi(x^n)] - \Phi_* = o(1/n)$  and,

$$(\forall n \in \mathbb{N}, n \geq 1) \quad E[\Phi(x^n)] - \Phi_* \leq \left[ \frac{\text{dist}^2(x^0, S_*)}{a} + \frac{b}{a} \right] \frac{1}{n}.$$

(iii) Suppose that P3 holds and  $S_* \neq \emptyset$ . Then, there exists a random variable  $x_*$  taking values in  $S_*$  such that  $x^n \rightarrow x_*$  P-a.s.

*Proof.* Taking the expectation in (4.5), we obtain

$$a(E[\Phi(x^{n+1})] - \Phi(x)) \leq E[\|x^n - x\|^2] - E[\|x^{n+1} - x\|^2] + E[\xi_n]. \quad (4.6)$$

(i): Since  $(E[\Phi(x^n)])_{n \in \mathbb{N}}$  is decreasing,  $E[\Phi(x^n)] \rightarrow \inf_{n \in \mathbb{N}} E[\Phi(x^n)] \geq \Phi_*$ . Thus, the statement is true if  $\inf_{n \in \mathbb{N}} E[\Phi(x^n)] = -\infty$ . Suppose that  $\inf_{n \in \mathbb{N}} E[\Phi(x^n)] > -\infty$  and let  $x \in \text{dom } \Phi$ . Then,

**P2** holds and the right hand side of (4.6), being summable, converges to zero. Therefore,  $\Phi_* \leq \lim_{n \rightarrow +\infty} \mathbb{E}[\Phi(x^{n+1})] \leq \Phi(x)$ . Since  $x$  is arbitrary in  $\text{dom } \Phi$ ,  $\mathbb{E}[\Phi(x^n)] \rightarrow \Phi_*$ .

(ii): Let  $x \in S_*$ . Then,  $\inf_{n \in \mathbb{N}} \mathbb{E}[\Phi(x^n)] \geq \Phi(x) > -\infty$ . Hence **P2** holds and (4.6) yields

$$a \sum_{n \in \mathbb{N}} (\mathbb{E}[\Phi(x^{n+1})] - \Phi_*) \leq \mathbb{E}[\|x^0 - x\|^2] + \sum_{n \in \mathbb{N}} \mathbb{E}[\xi_n] \leq \|x^0 - x\|^2 + b.$$

Therefore,  $\sum_{n \in \mathbb{N}} (\mathbb{E}[\Phi(x^{n+1})] - \Phi_*) \leq (1/a)\|x^0 - x\|^2 + b/a$ . Since  $(\mathbb{E}[\Phi(x^{n+1})] - \Phi_*)_{n \in \mathbb{N}}$  is decreasing, the statement follows from Fact 2.4.

(iii): Let  $x \in S_*$ . Then **P2** holds and, since  $\Phi(x) \leq \Phi(x^{n+1})$ , we derive from (4.5) that,

$$(\forall n \in \mathbb{N}) \quad \mathbb{E}[\|x^{n+1} - x\|^2 | \mathfrak{X}_n] \leq \|x^n - x\|^2 + \xi_n \quad \text{P-a.s.} \quad (4.7)$$

Note that  $\xi_n$  and  $\|x^n - x\|^2$  are  $\mathfrak{X}_n$ -measurable. Moreover  $\mathbb{E}[\sum_{n \in \mathbb{N}} \xi_n] = \sum_{n \in \mathbb{N}} \mathbb{E}[\xi_n] < +\infty$  and hence  $\sum_{n \in \mathbb{N}} \xi_n < +\infty$  P-a.s. Therefore  $(x^n)_{n \in \mathbb{N}}$  is a stochastic quasi-Fejér sequence with respect to  $S_*$  [11]. Then, in view of [4, Proposition 2.3(iv)] it is sufficient to prove that the weak limit points of  $(x^n)_{n \in \mathbb{N}}$  are contained in  $S_*$  P-a.s. By assumption **P3** there exist two sequences of  $H$ -valued random variables  $(y^n)_{n \in \mathbb{N}}$  and  $(v^n)_{n \in \mathbb{N}}$  and  $\tilde{\Omega} \subset \Omega$ ,  $\mathbb{P}(\tilde{\Omega}) = 1$  such that, for every  $\omega \in \tilde{\Omega}$ ,  $v^n(\omega) \in \partial\Phi(y^n(\omega))$ ,  $y^n(\omega) - x^n(\omega) \rightarrow 0$ ,  $v^n(\omega) \rightarrow 0$ . Let  $\omega \in \tilde{\Omega}$  and let  $(x^{n_k}(\omega))_{n \in \mathbb{N}}$  be a subsequence of  $(x^n(\omega))_{n \in \mathbb{N}}$  such that  $x^{n_k}(\omega) \rightarrow \bar{x}$ , for some  $\bar{x} \in H$ . Then,

$$y^{n_k}(\omega) \rightarrow \bar{x}, v^{n_k}(\omega) \rightarrow 0, v^{n_k}(\omega) \in \partial\Phi(y^{n_k}(\omega)).$$

Since  $\partial\Phi$  is weakly-strongly closed [1], we have  $0 \in \partial\Phi(\bar{x})$ , so  $\bar{x} \in S_*$ .  $\square$

**Remark 4.2.** Inequalities similar to (4.5) appear implicitly in the analysis of several deterministic and stochastic algorithms [3, 18, 26], to get rate of convergence for the function values. Moreover, (4.5) is related also to the concept introduced in [20], in a deterministic setting.

## 4.2 Convergence under convexity and strong convexity assumptions

In this section we address the convergence of Algorithm 1.1 in the convex and strongly convex case. The main results consist in the  $o(1/n)$  rate of convergence for the mean of the function values and in the almost sure weak convergence of the iterates. We start by recalling a standard result (see [36, Lemma 3.12(iii)]). Here we give a slightly more general version, including the moduli of strong convexity. The proof is given in Appendix B for reader's convenience.

**Lemma 4.3.** *Let  $H$  be a real Hilbert space. Let  $\varphi: H \rightarrow \mathbb{R}$  be differentiable and convex with modulus of strong convexity  $\mu_\varphi \geq 0$  and  $\psi: H \rightarrow ]-\infty, +\infty]$  be proper, lower semicontinuous, and convex with modulus of strong convexity  $\mu_\psi \geq 0$ . Let  $x \in H$  and set  $x^+ = \text{prox}_\psi(x - \nabla\varphi(x))$ . Then, for every  $z \in H$ ,*

$$\begin{aligned} (1 + \mu_\psi)\langle x - x^+, z - x \rangle &\leq \left( (\varphi + \psi)(z) - (\varphi + \psi)(x) - \frac{\mu_\varphi + \mu_\psi}{2} \|z - x\|^2 \right) \\ &\quad + (\psi(x) - \psi(x^+) + \langle \nabla\varphi(x), x - x^+ \rangle) - \left( 1 + \frac{\mu_\psi}{2} \right) \|x - x^+\|^2. \end{aligned}$$

**Proposition 4.4.** *Let **H1–H3** be satisfied. Let  $(\nu_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  and suppose that **S1** holds. Let  $(x^n)_{n \in \mathbb{N}}$  be generated by Algorithm 1.1 with, for every  $i \in [m]$ ,  $\gamma_i < 2/\nu_i$ . Set  $\delta = \max_{1 \leq i \leq m} \gamma_i \nu_i$  and  $\rho_{\min} =$*

$\min_{1 \leq i \leq m} \mathfrak{p}_i$ . Let  $\Gamma^{-1}$  be as in (4.1) and  $\mu_{\Gamma^{-1}}$  and  $\sigma_{\Gamma^{-1}}$  be the moduli of strong convexity of  $f$  and  $h$  respectively, in the norm  $\|\cdot\|_{\Gamma^{-1}}$ . Set  $F = f + h$ . Then,

$$\begin{aligned} (1 + \sigma_{\Gamma^{-1}}) \langle x^n - \bar{x}^{n+1}, x - x^n \rangle_{\Gamma^{-1}} &\leq \frac{1}{\mathfrak{p}_{\min}} \mathbb{E} [F(x^n) - F(x^{n+1}) \mid \mathfrak{E}_{n-1}] \\ &\quad + \left( F(x) - F(x^n) - \frac{\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}}}{2} \|x^n - x\|_{\Gamma^{-1}}^2 \right) \\ &\quad + \frac{\delta - 2 - \sigma_{\Gamma^{-1}}}{2} \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}}^2. \end{aligned} \quad (4.8)$$

*Proof.* Let  $x \in \text{dom } F$  and  $n \in \mathbb{N}$ . Since for all  $v \in H$ ,  $\langle \nabla^{\Gamma^{-1}} f(x^n), v \rangle_{\Gamma^{-1}} = \langle \nabla f(x^n), v \rangle$ , we derive from Lemma 4.3, written in the norm  $\|\cdot\|_{\Gamma^{-1}}$ , and (4.3) that

$$\begin{aligned} (1 + \sigma_{\Gamma^{-1}}) \langle x^n - \bar{x}^{n+1}, x - x^n \rangle_{\Gamma^{-1}} &\leq (h(x^n) - h(\bar{x}^{n+1}) + \langle \nabla f(x^n), x^n - \bar{x}^{n+1} \rangle) \\ &\quad + \left( F(x) - F(x^n) - \frac{\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}}}{2} \|x^n - x\|_{\Gamma^{-1}}^2 \right) \\ &\quad - \left( 1 + \frac{\sigma_{\Gamma^{-1}}}{2} \right) \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}}^2. \end{aligned} \quad (4.9)$$

Now, we majorize  $h(x^n) - h(\bar{x}^{n+1}) + \langle \nabla f(x^n), x^n - \bar{x}^{n+1} \rangle$ . By Fact 2.2 and Fact 2.3, we have

$$\begin{aligned} h(x^n) - h(\bar{x}^{n+1}) + \langle \nabla f(x^n), x^n - \bar{x}^{n+1} \rangle \\ = \mathbb{E} \left[ \sum_{i=1}^m \frac{\varepsilon_i^n}{\mathfrak{p}_i} \left( h_i(x_i^n) - h_i(\bar{x}_i^{n+1}) + \langle \nabla_i f(x^n), x_i^n - \bar{x}_i^{n+1} \rangle \right) \mid \mathfrak{E}_{n-1} \right]. \end{aligned}$$

Moreover,

$$\begin{aligned} &\sum_{i=1}^m \frac{\varepsilon_i^n}{\mathfrak{p}_i} \left( h_i(x_i^n) - h_i(\bar{x}_i^{n+1}) + \langle \nabla_i f(x^n), x_i^n - \bar{x}_i^{n+1} \rangle \right) \\ &= \sum_{i=1}^m \frac{1}{\mathfrak{p}_i} \left( h_i(x_i^n) - h_i(x_i^{n+1}) + \langle \nabla_i f(x^n), x_i^n - x_i^{n+1} \rangle \right) \\ &= \frac{1}{\mathfrak{p}_{\min}} \left( h(x^n) - h(x^{n+1}) + \langle \nabla f(x^n), x^n - x^{n+1} \rangle \right) \\ &\quad - \sum_{i=1}^m \underbrace{\left( \frac{1}{\mathfrak{p}_{\min}} - \frac{1}{\mathfrak{p}_i} \right)}_{\geq 0} \left( h_i(x_i^n) - h_i(x_i^{n+1}) + \langle \nabla_i f(x^n), x_i^n - x_i^{n+1} \rangle \right) \\ &\leq \frac{1}{\mathfrak{p}_{\min}} \left( h(x^n) - h(x^{n+1}) + \langle \nabla f(x^n), x^n - x^{n+1} \rangle \right) \\ &\quad - \left( 1 + \frac{\sigma_{\Gamma^{-1}}}{2} \right) \sum_{i=1}^m \left( \frac{1}{\mathfrak{p}_{\min}} - \frac{1}{\mathfrak{p}_i} \right) \frac{\varepsilon_i^n}{\gamma_i} \|\Delta_i^n\|^2, \end{aligned}$$

where in the last inequality we used that

$$-\left( h_i(x_i^n) - h_i(x_i^{n+1}) + \langle \nabla_i f(x^n), x_i^n - x_i^{n+1} \rangle \right) \leq -\frac{\varepsilon_i^n}{\gamma_i} \left( 1 + \frac{\sigma_{\Gamma^{-1}}}{2} \right) \|\Delta_i^n\|^2, \quad (4.10)$$

which was obtained from (2.3) with

$$x_i = x_i^n, \quad y_i = x_i^{n+1}, \quad v_i = \frac{x_i^n - x_i^{n+1}}{\gamma_i} - \nabla_i f(x^n) \in \partial h_i(x_i^{n+1}), \quad \text{for } \varepsilon_i^n = 1.$$

Therefore,

$$\begin{aligned} & h(x^n) - h(\bar{x}^{n+1}) + \langle \nabla f(x^n), x^n - \bar{x}^{n+1} \rangle \\ & \leq \frac{1}{\rho_{\min}} \mathbb{E}[h(x^n) - h(x^{n+1}) + \langle \nabla f(x^n), x^n - x^{n+1} \rangle \mid \mathfrak{E}_{n-1}] \\ & \quad - \frac{1}{\rho_{\min}} \left(1 + \frac{\sigma_{\Gamma-1}}{2}\right) \sum_{i=1}^m \frac{\rho_i}{\gamma_i} \|\Delta_i^n\|^2 + \left(1 + \frac{\sigma_{\Gamma-1}}{2}\right) \|\bar{x}^{n+1} - x^n\|_{\Gamma-1}^2. \end{aligned} \quad (4.11)$$

Next, it follows from (4.3), S1, and Fact 2.2 that

$$\mathbb{E}[\langle \nabla f(x^n), x^n - x^{n+1} \rangle \mid \mathfrak{E}_{n-1}] \leq \mathbb{E}[f(x^n) - f(x^{n+1}) \mid \mathfrak{E}_{n-1}] + \frac{1}{2} \sum_{i=1}^m \rho_i \nu_i \|\Delta_i^n\|^2.$$

Then, we derive from (4.11) that

$$\begin{aligned} & h(x^n) - h(\bar{x}^{n+1}) + \langle \nabla f(x^n), x^n - \bar{x}^{n+1} \rangle \\ & \leq \frac{1}{\rho_{\min}} \mathbb{E}[F(x^n) - F(x^{n+1}) \mid \mathfrak{E}_{n-1}] \\ & \quad - \frac{1}{2\rho_{\min}} \sum_{i=1}^m \left(2 + \sigma_{\Gamma-1} - \gamma_i \nu_i\right) \frac{\rho_i}{\gamma_i} \|\Delta_i^n\|^2 + \left(1 + \frac{\sigma_{\Gamma-1}}{2}\right) \|\bar{x}^{n+1} - x^n\|_{\Gamma-1}^2. \end{aligned}$$

The statement follows from (4.9), considering that

$$\begin{aligned} \frac{1}{\rho_{\min}} \sum_{i=1}^m \left(\gamma_i \nu_i - 2 - \sigma_{\Gamma-1}\right) \frac{\rho_i}{\gamma_i} \|\Delta_i^n\|^2 & \leq \underbrace{\frac{\delta - 2 - \sigma_{\Gamma-1}}{\rho_{\min}} \sum_{i=1}^m \frac{\rho_i}{\gamma_i} \|\Delta_i^n\|^2}_{\leq 0} \\ & \leq \frac{\delta - 2 - \sigma_{\Gamma-1}}{\rho_{\min}} \sum_{i=1}^m \frac{\rho_{\min}}{\gamma_i} \|\Delta_i^n\|^2 \\ & = (\delta - 2 - \sigma_{\Gamma-1}) \|x^n - \bar{x}^{n+1}\|_{\Gamma-1}^2. \end{aligned} \quad \square$$

**Proposition 4.5.** *Let H1–H3 be satisfied. Let  $\Gamma^{-1}$  and  $W$  be as in (4.1) and  $(x^n)_{n \in \mathbb{N}}$  be generated by Algorithm 1.1. Let  $n \in \mathbb{N}$  and  $x$  be an  $H$ -valued random variable which is measurable w.r.t.  $\mathfrak{E}_{n-1}$ . Then*

$$\mathbb{E}[\|x^{n+1} - x\|_W^2 \mid \mathfrak{E}_{n-1}] - \|x^n - x\|_W^2 = \|\bar{x}^{n+1} - x\|_{\Gamma-1}^2 - \|x^n - x\|_{\Gamma-1}^2 \quad (4.12)$$

and  $\mathbb{E}[\|x^{n+1} - x^n\|_W^2 \mid \mathfrak{E}_{n-1}] = \|\bar{x}^{n+1} - x^n\|_{\Gamma-1}^2$ .

*Proof.* It follows from (4.3), Fact 2.2, and Fact 2.3 that

$$\begin{aligned} \mathbb{E}[\|x^{n+1} - x\|_W^2 \mid \mathfrak{E}_{n-1}] & = \mathbb{E}\left[\sum_{i=1}^m \frac{1}{\gamma_i \rho_i} \|x_i^{n+1} - x_i\|^2 \mid \mathfrak{E}_{n-1}\right] \\ & = \mathbb{E}\left[\sum_{i=1}^m \frac{\varepsilon_i^n}{\gamma_i \rho_i} \|\bar{x}_i^{n+1} - x_i\|^2 \mid \mathfrak{E}_{n-1}\right] + \mathbb{E}\left[\sum_{i=1}^m \frac{1 - \varepsilon_i^n}{\gamma_i \rho_i} \|x_i^n - x_i\|^2 \mid \mathfrak{E}_{n-1}\right] \\ & = \|\bar{x}^{n+1} - x\|_{\Gamma-1}^2 + \|x^n - x\|_W^2 - \|x^n - x\|_{\Gamma-1}^2. \end{aligned}$$

The second equation follows from (4.12), by choosing  $x = x^n$ .  $\square$

The following result is a stochastic version of [36, Proposition 3.15].

**Proposition 4.6.** *Let H1–H3 be satisfied. Let  $(\nu_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  and suppose that S1 holds. Let  $(x^n)_{n \in \mathbb{N}}$  be generated by Algorithm 1.1 with, for every  $i \in [m]$ ,  $\gamma_i < 2/\nu_i$ . Set  $\delta = \max_{1 \leq i \leq m} \gamma_i \nu_i$  and  $\rho_{\min} = \min_{1 \leq i \leq m} \rho_i$ . Let  $\Gamma^{-1}$  and  $W$  be as in (4.1) and  $\mu_{\Gamma^{-1}}$  and  $\sigma_{\Gamma^{-1}}$  be the moduli of strong convexity of  $f$  and  $h$  respectively, in the norm  $\|\cdot\|_{\Gamma^{-1}}$ . Set  $F = f + h$ . Then, the following hold.*

- (i)  $(\mathbb{E}[F(x^n)])_{n \in \mathbb{N}}$  is decreasing.
- (ii) Suppose that  $\inf_{n \in \mathbb{N}} \mathbb{E}[F(x^n)] > 0$ . Then,

$$\sum_{n \in \mathbb{N}} \|\bar{x}^{n+1} - x^n\|_{\Gamma^{-1}}^2 = \sum_{n \in \mathbb{N}} \mathbb{E}[\|x^n - x^{n+1}\|_W^2 | \mathfrak{E}_{n-1}] < +\infty \quad \text{P a.s.}$$

- (iii) For every  $n \in \mathbb{N}$  and every  $x \in \text{dom } F$

$$\begin{aligned} (1 + \sigma_{\Gamma^{-1}}) \mathbb{E}[\|x^{n+1} - x\|_W^2 | \mathfrak{E}_{n-1}] &\leq (1 + \sigma_{\Gamma^{-1}}) \|x^n - x\|_W^2 \\ &\quad - 2 \left( F(x^n) - F(x) + \frac{\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}}}{2} \|x^n - x\|_{\Gamma^{-1}}^2 \right) \\ &\quad + \frac{2}{\rho_{\min}} \left( \frac{(\delta - 1)_+}{2 + \sigma_{\Gamma^{-1}} - \delta} + 1 \right) \mathbb{E}[F(x^n) - F(x^{n+1}) | \mathfrak{E}_{n-1}]. \end{aligned}$$

*Proof.* Let  $n \in \mathbb{N}$  and  $x \in \text{dom } F$ . Since

$$\|x^n - x\|_{\Gamma^{-1}}^2 - \|\bar{x}^{n+1} - x\|_{\Gamma^{-1}}^2 = -\|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}}^2 + 2\langle x^n - \bar{x}^{n+1}, x^n - x \rangle_{\Gamma^{-1}},$$

we derive from (4.8), multiplied by 2, that

$$\begin{aligned} (1 + \sigma_{\Gamma^{-1}}) \|\bar{x}^{n+1} - x\|_{\Gamma^{-1}}^2 &\leq (1 + \sigma_{\Gamma^{-1}}) \|x^n - x\|_{\Gamma^{-1}}^2 + (\delta - 1) \|\bar{x}^{n+1} - x^n\|_{\Gamma^{-1}}^2 \\ &\quad + \frac{2}{\rho_{\min}} \mathbb{E}[F(x^n) - F(x^{n+1}) | \mathfrak{E}_{n-1}] \\ &\quad - 2 \left( F(x^n) - F(x) + \frac{\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}}}{2} \|x^n - x\|_{\Gamma^{-1}}^2 \right). \end{aligned} \quad (4.13)$$

Then for an  $H$ -valued  $\mathfrak{E}_{n-1}$ -measurable random variable  $x$ , Proposition 4.5 yields

$$\begin{aligned} (1 + \sigma_{\Gamma^{-1}}) \mathbb{E}[\|x^{n+1} - x\|_W^2 | \mathfrak{E}_{n-1}] &\leq (1 + \sigma_{\Gamma^{-1}}) \|x^n - x\|_W^2 + (\delta - 1) \mathbb{E}[\|x^{n+1} - x^n\|_W^2 | \mathfrak{E}_{n-1}] \\ &\quad + \frac{2}{\rho_{\min}} \mathbb{E}[F(x^n) - F(x^{n+1}) | \mathfrak{E}_{n-1}] \\ &\quad - 2 \left( F(x^n) - F(x) + \frac{\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}}}{2} \|x^n - x\|_{\Gamma^{-1}}^2 \right). \end{aligned} \quad (4.14)$$

Taking  $x = x^n$  in (4.14), we have

$$\frac{\rho_{\min}}{2} (2 + \sigma_{\Gamma^{-1}} - \delta) \mathbb{E}[\|x^{n+1} - x^n\|_W^2 | \mathfrak{E}_{n-1}] \leq \mathbb{E}[F(x^n) - F(x^{n+1}) | \mathfrak{E}_{n-1}], \quad (4.15)$$

which plugged into (4.14), with  $x \equiv x \in \text{dom } F$ , gives (iii). Moreover, taking the expectation in (4.15), we obtain

$$\frac{\text{P}_{\min}}{2}(2 + \sigma_{\Gamma-1} - \delta)\mathbb{E}[\|x^{n+1} - x^n\|_{\mathbb{W}}^2] \leq \mathbb{E}[F(x^n)] - \mathbb{E}[F(x^{n+1})], \quad (4.16)$$

which gives (i). Finally, set for all  $n \in \mathbb{N}$ ,  $\xi_n = \mathbb{E}[F(x^n) - F(x^{n+1}) | \mathfrak{E}_{n-1}] \geq 0$ . Then

$$\mathbb{E}\left[\sum_{n=0}^{+\infty} \xi_n\right] = \sum_{n=0}^{+\infty} \mathbb{E}[\xi_n] = \sum_{n=0}^{+\infty} \mathbb{E}[F(x^n)] - \mathbb{E}[F(x^{n+1})] \leq \mathbb{E}[F(x^0)] - \inf_{n \in \mathbb{N}} \mathbb{E}[F(x^n)].$$

This shows that if  $\inf_{n \in \mathbb{N}} \mathbb{E}[F(x^n)] > 0$ , then  $\sum_{n=0}^{+\infty} \xi_n$  is P-integrable and hence it is P-a.s. finite. Then (ii) follows from (4.15) and Proposition 4.5.  $\square$

**Proposition 4.7.** *Under the same assumptions of Proposition 4.6, suppose that condition S1 is replaced by condition S2. Then*

$$(\forall n \in \mathbb{N}) \quad \frac{2 + \sigma_{\Gamma-1} - \delta}{2} \|x^{n+1} - x^n\|_{\Gamma-1}^2 \leq F(x^n) - F(x^{n+1}) \quad \text{P a.s.}$$

*Proof.* We derive from S2 (since  $\varepsilon^n$  has the same distribution of  $\varepsilon$ ) and (4.3) that

$$\langle \nabla f(x^n), x^{n+1} - x^n \rangle \leq f(x^n) - f(x^{n+1}) + \sum_{i=1}^m \frac{1}{2} \varepsilon_i^n \nu_i \|\Delta_i^n\|^2 \quad \text{P a.s.} \quad (4.17)$$

Therefore, summing (4.10), from  $i = 1$  to  $m$ , we have

$$\begin{aligned} \frac{2 + \sigma_{\Gamma-1}}{2} \sum_{i=1}^m \varepsilon_i^n \frac{1}{\gamma_i} \|\Delta_i^n\|^2 &\leq h(x^n) - h(x^{n+1}) + \langle \nabla f(x^n), x^n - x^{n+1} \rangle \\ &\leq h(x^n) - h(x^{n+1}) + f(x^n) - f(x^{n+1}) + \frac{1}{2} \sum_{i=1}^m \varepsilon_i^n \nu_i \|\Delta_i^n\|^2 \quad \text{P a.s.} \end{aligned}$$

Hence  $(1/2) \sum_{i=1}^m (2 + \sigma_{\Gamma-1} - \gamma_i \nu_i) \gamma_i^{-1} \varepsilon_i^n \|\Delta_i^n\|^2 \leq F(x^n) - F(x^{n+1})$  P-a.s.  $\square$

**Proposition 4.8.** *Under the assumptions of Proposition 4.6, suppose in addition that F is bounded from below. Then, there exist  $(y^n)_{n \in \mathbb{N}}$  and  $(v^n)_{n \in \mathbb{N}}$ , sequences of H-valued random variables, such that the following hold.*

(i)  $v^n \in \partial F(y^n)$  P-a.s.

(ii)  $y^n - x^n \rightarrow 0$  and  $v^n \rightarrow 0$  P-a.s.

*Proof.* It follows from (4.2) that,  $(x_i^n(\omega) - \bar{x}_i^{n+1}(\omega))/\gamma_i - \nabla_i f(x^n(\omega)) \in \partial h_i(\bar{x}_i^{n+1}(\omega))$ , for all  $i \in [m]$  and  $\omega \in \Omega$ . Hence

$$\left( \frac{x_i^n(\omega) - \bar{x}_i^{n+1}(\omega)}{\gamma_i} \right)_{1 \leq i \leq m} - \nabla f(x^n(\omega)) \in \partial h(\bar{x}^{n+1}).$$

Set  $y^n = \bar{x}^{n+1}$  and let  $v^n : \Omega \rightarrow H$  be such that, for every  $\omega \in \Omega$ ,

$$\begin{aligned} v^n(\omega) &= \left( \frac{x_i^n(\omega) - y_i^n(\omega)}{\gamma_i} \right)_{1 \leq i \leq m} + \nabla f(y^n(\omega)) - \nabla f(x^n(\omega)) \\ &\in \partial h(y^n(\omega)) + \nabla f(y^n(\omega)) = \partial F(y^n(\omega)). \end{aligned}$$

Clearly  $v^n$  is measurable and hence it is a random variable. Moreover, for every  $\omega \in \Omega$ ,

$$\|v^n(\omega)\| \leq \frac{1}{\gamma_{\min}} \|x^n(\omega) - y^n(\omega)\| + \|\nabla f(y^n(\omega)) - \nabla f(x^n(\omega))\|.$$

Now, since  $F$  is bounded from below, Proposition 4.6(ii) yields that  $(\|y^n - x^n\|_{\Gamma^{-1}}^2)_{n \in \mathbb{N}}$  is summable P-a.s. and hence  $y^n - x^n \rightarrow 0$  P-a.s. The statement follows from the fact that  $\nabla f$  is Lipschitz continuous (see Theorem 3.1(iv)).  $\square$

Now we are ready to state one of the main convergence results of this paper. From one hand, it extends to the stochastic setting a well-known convergence rate of the (deterministic) forward-backward algorithm [7, 15, 36]. On the other hand, it proves the almost sure weak convergence of the iterates of Algorithm 1.1 in the convex case. We stress that none of the works [21, 23, 31, 33, 34, 35, 38] addresses this latter aspect. To the best of our knowledge, [4] is the only work that proves almost sure weak convergence of the iterates. However, in [4, Corollary 5.11] the stepsize is set according to the (global) Lipschitz constant of  $\nabla f$  which, in general, leads to smaller stepsizes and worse upper bounds on convergence rates. See the subsequent discussion.

**Theorem 4.9.** *Let H1–H3 be satisfied. Let  $(\nu_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  and suppose that S1 holds. Let  $(x^n)_{n \in \mathbb{N}}$  be generated by Algorithm 1.1 with, for every  $i \in [m]$ ,  $\gamma_i < 2/\nu_i$ . Set  $\delta = \max_{1 \leq i \leq m} \gamma_i \nu_i$  and  $\rho_{\min} = \min_{1 \leq i \leq m} \rho_i$ . Let  $W$  be as in (4.1) and set  $F = f + h$ ,  $F_* = \inf F$ , and  $S_* = \operatorname{argmin} F \subset H$ . Then, the following hold.*

(i)  $E[F(x^n)] \rightarrow F_*$ .

(ii) Suppose that  $S_* \neq \emptyset$ . Then  $E[F(x^n)] - F_* = o(1/n)$  and, for every integer  $n \geq 1$ ,

$$E[F(x^n)] - F_* \leq \left[ \frac{\operatorname{dist}_W^2(x^0, S_*)}{2} + \left( \frac{\max\{1, (2 - \delta)^{-1}\}}{\rho_{\min}} - 1 \right) (F(x^0) - F_*) \right] \frac{1}{n}. \quad (4.18)$$

Moreover, there exists a random variable  $x_*$  taking values in  $S_*$  such that  $x^n \rightarrow x_*$  P-a.s.

*Proof.* Proposition 4.6(iii) with  $\mu_{\Gamma^{-1}} = \sigma_{\Gamma^{-1}} = 0$  gives, for all  $x \in \operatorname{dom} F$  and  $n \in \mathbb{N}$ ,

$$E[\|x^{n+1} - x\|_W^2 | \mathfrak{E}_{n-1}] \leq \|x^n - x\|_W^2 + 2E[F(x) - F(x^{n+1}) | \mathfrak{E}_{n-1}] + \xi_n,$$

where

$$\xi_n = b_1 E[F(x^n) - F(x^{n+1}) | \mathfrak{E}_{n-1}], \quad b_1 = 2 \left( \frac{\max\{1, 1/(2 - \delta)\}}{\rho_{\min}} - 1 \right).$$

Note that the random variables  $x^n$ 's are discrete with finite range and  $(E[F(x^n)])_{n \in \mathbb{N}}$  is decreasing. Moreover,  $\sum_{n \in \mathbb{N}} E[\xi_n] \leq b_1 (F(x^0) - \inf_{n \in \mathbb{N}} E[F(x^n)])$ . Therefore, the statement follows from Theorem 4.1 and Proposition 4.8.  $\square$

**Discussion.** In the following we examine some crucial aspects related to Algorithm 1.1. We suppose that H4 and L1 hold and that for every  $i \in [m]$ ,  $\gamma_i = \delta/\nu_i$  with  $\delta \in ]0, 2[$ .

*The benefit of a parallel block update.* Here we discuss the advantage of updating multiple blocks in parallel instead of just a single block. We consider the setting of a  $\tau$ -nice uniform block sampling, which was described in Section 3.2. In this case, for every  $i, j \in [m]$ , with  $i \neq j$ ,  $\rho_i = \rho_j$  and, since

$\tau = \sum_{i=1}^m \mathbb{E}[\varepsilon_i] = \sum_{i=1}^m p_i$ , we have, for every  $i \in [m]$ ,  $p_i = p := \tau/m$ . Moreover, we can set for every  $i \in [m]$ ,  $\nu_i = \beta_1 L_i$  with  $\beta_1$  defined as in (3.5). In order to compare different choices of  $\tau$ , we normalize the iterations so to match the same computational cost per iteration of the standard (full parallel) forward-backward algorithm (FB). It follows from (4.18) that after  $n_\tau = \lceil m\bar{n}/\tau \rceil$  iterations of Algorithm 1.1, which have the same total computational cost of  $\bar{n}$  iterations of FB, we have

$$\mathbb{E}[F(x^{\lceil m\bar{n}/\tau \rceil})] - F_* \leq \left[ \beta_1 \frac{\text{dist}_\Lambda^2(x^0, S_*)}{2\delta} + \left( \max \left\{ 1, \frac{1}{2-\delta} \right\} - p \right) (F(x^0) - F_*) \right] \frac{1}{\bar{n}}, \quad (4.19)$$

where  $\Lambda = \bigoplus_{i=1}^m L_i \text{Id}_i$ . Now, since,  $\beta_1 = 1 + (\tau - 1)(\eta - 1)/(m - 1)$ , we see that if  $\eta \ll m$  and  $\tau \ll m$ , then  $\beta_1$  is close to 1 (and  $p$  is close to zero), so that  $\beta_1$  nearly does not depend on  $\tau$ , as long as  $\tau$  remains sufficiently small. For instance, in Section 6 we consider the setting where  $m = 10^5$  and  $\eta = 148$ . In such case, if we let  $\tau = 1, 5, 10, 50$ , the corresponding  $\beta_1$ 's and  $p$ 's are essentially the same so that the right hand side of (4.19) does not change much. Therefore, the above options for  $\tau$  require the same total amount of computations (i.e.,  $m\bar{n}$  block-coordinate updates) and lead essentially to the same improvement in the objective function. However, a parallel implementation, say with  $\tau = 50$  on a CPU with 50 cores, will be 50 times faster than a serial implementation ( $\tau = 1$ ) which uses only one core per iteration. In summary, in the large scale ( $m$  large) and sparse ( $\eta \ll m$ ) setting, the parallel strategy ( $\tau > 1$  and  $\tau$  equal to the number of CPU cores) is definitely advantageous provided that  $\tau$  is sufficiently small compared to  $m$ .

*Comparison with [4].* The almost sure weak convergence of the iterates of Algorithm 1.1 is also obtained in [4], but with stepsizes set according to the global Lipschitz constant of the gradient of  $f$ . Let  $L$  be the Lipschitz constant of  $\nabla f$  and note that  $f$  is also Lipschitz smooth in the norm  $\|\cdot\|_\Lambda$ , defined by the operator  $\Lambda = \bigoplus_{i=1}^m L_i \text{Id}_i$ , with constant  $\eta$  (see Corollary A.2(iv)). Therefore, the results in [4] can be applied in the original norm  $\|\cdot\|$  or in the norm  $\|\cdot\|_\Lambda$ . In this respect we note that since

$$\text{prox}_{\alpha h}^\Lambda(x^n - \alpha \nabla^\Lambda f(x^n)) = (\text{prox}_{(\alpha/L_i)h_i}(x_i^n - (\alpha/L_i)\nabla_i f(x^n)))_{1 \leq i \leq m} \quad (\alpha = \delta/\eta, 0 < \delta < 2),$$

the implementation in the norm  $\|\cdot\|_\Lambda$  is nothing but Algorithm 1.1 with stepsizes  $\gamma_i = \delta/(\eta L_i)$ . In both cases Corollary 5.11 in [4], applied in the corresponding norms, proves weak convergence of the iterates for Algorithm 1.1 with stepsizes  $\gamma_i \equiv \delta/L$  and  $\gamma_i = \delta/(\eta L_i)$  respectively. However, Theorem 4.9, together with Theorem 3.1, allows to set the stepsizes as  $\gamma_i = \delta/(\beta_{1,i} L_i)$ . Since  $\beta_{1,i} L_i$  may be much smaller than  $L$  and, in view of Remark 3.3, it is always smaller than  $\eta L_i$ , Theorem 4.9 provides a significant improvement over [4] in terms of flexibility in the stepsizes.

*The advantage over the standard FB.* We consider the forward-backward algorithm (FB) in the original norm of  $H$  and in the norm  $\|\cdot\|_\Lambda$ . All the remarks about the stepsizes discussed in the previous paragraph apply also here. Moreover, in the case  $\delta = 1$ , standard convergence rate for FB (see e.g., [7]) yields that after  $\bar{n}$  iterations, we have

$$F(x^{\bar{n}}) - F_* \leq \frac{L \text{dist}^2(x^0, S_*)}{2\bar{n}} \quad \text{or} \quad F(x^{\bar{n}}) - F_* \leq \frac{\eta \text{dist}_\Lambda^2(x^0, S_*)}{2\bar{n}}, \quad (4.20)$$

depending on which of the two above implementations of FB we consider. In order to appropriately compare the rates (4.20) with that of Algorithm 1.1 given in Theorem 4.9 in the following we set  $\delta = 1$  and analyze two choices of the block sampling.

- (i) Assume that we perform a  $\tau$ -nice block sampling. Then we saw that the (normalized) convergence rate of Algorithm 1.1 is (4.19). We first note that (4.19) reduces to the second inequality



in (4.20) when  $\delta = 1$  and  $\tau = m$ . Comparing the bounds in (4.20) with (4.19) (with  $\delta = 1$ ) we see that, if we assume that the terms  $\text{dist}_\Lambda^2(x^0, S_*)$  and  $F(x^0) - F_*$  are about of the same magnitude, then Algorithm 1.1 features always a better rate than FB if implemented in the norm  $\|\cdot\|_\Lambda$  (since  $\beta_1 \leq \eta$ ), whereas if FB is implemented in the original norm of  $H$ , Algorithm 1.1 is still a better choice provided that  $\beta_1 \max_{1 \leq i \leq m} L_i \leq L$ .

- (ii) Suppose that the block sampling performs on average  $\tau$  updates per iteration and that, for every  $i \in [m]$ ,  $\rho_i$  is proportional to the Lipschitz constant  $L_i$ , that is,  $\rho_i = \tau L_i / (\sum_{j=1}^m L_j)$  (provided that  $\tau \leq (\sum_{j=1}^m L_j) / \max_{1 \leq j \leq m} L_j$ ). In this case, as stated at the beginning of Section 3.2, we can let  $\nu_i = \beta_2 L_i$  and  $\gamma_i = 1/\nu_i$ . Then, (4.18) becomes for  $\delta = 1$  and  $n = \lceil m\bar{n}/\tau \rceil$ ,

$$\mathbb{E}[F(x^{\lceil m\bar{n}/\tau \rceil})] - F_* \leq \left[ \beta_2 \bar{L} \frac{\text{dist}^2(x^0, S_*)}{2} + \left( \frac{\bar{L}}{L_{\min}} - \frac{\tau}{m} \right) (F(x^0) - F_*) \right] \frac{1}{\bar{n}},$$

where  $\bar{L} = \sum_{i=1}^m L_i/m$  and  $L_{\min} = \min_{1 \leq i \leq m} L_i$ . Here we see that Algorithm 1.1 can be superior to FB if  $(\beta_2 + 2\bar{L}/L_{\min})\bar{L} \leq L$ , under the assumption that  $\bar{L}\text{dist}^2(x^0, S_*) \simeq F(x^0) - F_*$ .

We now provide an additional convergence theorem, analyzing the strongly convex case, which extends [21, Theorem 1] and [38, Theorem 3] to an arbitrary (not necessarily uniform) sampling and to the more general stepsize rule (1.5). The proof is still based on Proposition 4.6(iii) and is postponed to Appendix B.

**Theorem 4.10.** *Under the same assumptions of Theorem 4.9, let  $\mu_{\Gamma-1}$  and  $\sigma_{\Gamma-1}$  be the moduli of strong convexity of  $f$  and  $h$  respectively, in the norm  $\|\cdot\|_{\Gamma-1}$ , and suppose that  $\mu_{\Gamma-1} + \sigma_{\Gamma-1} > 0$  and that  $S_* = \{x_*\}$ . Then, for every  $n \in \mathbb{N}$ ,*

$$\mathbb{E}[F(x^n)] - F_* \leq (1 - \rho_{\min} \bar{\lambda})^n \left( \rho_{\min} (1 + \sigma_{\Gamma-1} - (\delta - 1)_+) \frac{\|x^0 - x_*\|_{\mathbb{W}}^2}{2} + F(x^0) - F_* \right),$$

where

$$\bar{\lambda} = \begin{cases} \frac{2 - \delta + \sigma_{\Gamma-1}}{1 + \sigma_{\Gamma-1}} & \text{if } \delta > 1 \text{ and } \mu_{\Gamma-1} \geq 2 - \delta \\ \frac{2(\mu_{\Gamma-1} + \sigma_{\Gamma-1})}{1 + \sigma_{\Gamma-1} + (\mu_{\Gamma-1} + \sigma_{\Gamma-1})(1 + \sigma_{\Gamma-1})/(1 + \sigma_{\Gamma-1} - (\delta - 1)_+)} & \text{otherwise.} \end{cases} \quad (4.21)$$

**Remark 4.11.** Let  $\gamma_i = \delta/\nu_i$  and  $V = \bigoplus_{i=1}^m \nu_i \text{Id}_i$ . Let  $\mu_V$  and  $\sigma_V$  be the moduli of strong convexity of  $f$  and  $h$  respectively, in the norm  $\|\cdot\|_V$ . Then it is easy to see that  $\mu_{\Gamma-1} = \delta\mu_V$  and  $\sigma_{\Gamma-1} = \delta\sigma_V$ . Moreover, as in Remark 2.1, one can also see that  $\mu_V \leq 1$ . Then, (4.21) becomes

$$\bar{\lambda} = \begin{cases} \frac{2/\delta - 1 + \sigma_V}{1/\delta + \sigma_V} & \text{if } \delta \geq \frac{2}{1 + \mu_V} \\ \frac{2(\mu_V + \sigma_V)}{1/\delta + \sigma_V + (\mu_V + \sigma_V)(1/\delta + \sigma_V)/(2/\delta - 1 + \sigma_V)} & \text{if } 1 < \delta \leq \frac{2}{1 + \mu_V} \\ \frac{2(\mu_V + \sigma_V)}{1/\delta + \sigma_V + (\mu_V + \sigma_V)} & \delta \leq 1. \end{cases} \quad (4.22)$$

One can check that the maximum of  $\bar{\lambda}$  with respect to  $\delta \in ]0, 2[$  is

$$\bar{\lambda}_{\text{opt}} = \frac{4(\mu_V + \sigma_V)}{(\sqrt{1 + \sigma_V} + \sqrt{\mu_V + \sigma_V})^2} \in ]0, 1], \quad (4.23)$$

which is achieved at

$$\delta = \delta_{\text{opt}} := \frac{2}{1 - \sigma_V + \sqrt{(\mu_V + \sigma_V)(1 + \sigma_V)}} \in [1, 2[. \quad (4.24)$$

Note that if  $\mu_V < 1$  (as is normally the case), then  $\bar{\lambda}_{\text{opt}} \in ]0, 1[$  and  $\delta_{\text{opt}} > 1$ .

**Remark 4.12.** If  $\mu$  and  $\sigma$  are the moduli of strong convexity of  $f$  and  $h$  respectively in the original norm. Let, for every  $i \in [m]$ ,  $\gamma_i = \delta/\nu_i$  and set  $\nu_{\max} = \max_{1 \leq i \leq m} \nu_i$ . Then  $\mu_V = \mu/\nu_{\max}$  and  $\sigma_V = \sigma/\nu_{\max}$ . Therefore, the optimal stepsizes are achieved for

$$\delta = \frac{2\nu_{\max}}{\nu_{\max} - \sigma + \sqrt{(\mu + \sigma)(\nu_{\max} + \sigma)}} \quad (4.25)$$

and the corresponding rate in Theorem 4.10 becomes

$$\mathbb{E}[F(x^n)] - F_* \leq \left(1 - \mathfrak{p}_{\min} \frac{4(\mu + \sigma)}{(\sqrt{\nu_{\max} + \sigma} + \sqrt{\mu + \sigma})^2}\right)^n \text{const.}$$

**Remark 4.13.** Suppose that the block sampling is uniform, that is,  $\mathfrak{p}_i = \mathfrak{p}$  for all  $i \in [m]$  and let, for every  $i \in [m]$ ,  $\gamma_i = 1/\nu_i$ . Then  $\delta = 1$  and Theorem 4.10 reduce to

$$\mathbb{E}[F(x^n)] - F_* \leq \left(1 - \mathfrak{p} \frac{2(\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}})}{1 + \mu_{\Gamma^{-1}} + 2\sigma_{\Gamma^{-1}}}\right)^n \left( (1 + \sigma_{\Gamma^{-1}}) \frac{\|x^0 - x_*\|_{\Gamma^{-1}}^2}{2} + F(x^0) - F_* \right). \quad (4.26)$$

This result was obtained in [38, Theorem 3], which is in turn a generalization of [21, Theorem 1], treating the serial case ( $\mathbb{P}(\sum_{i=1}^m \varepsilon_i = 1) = 1$ ). Thus, Theorem 4.10 and the subsequent Remark 4.11 show that the rate in (4.26) can indeed be improved by choosing  $\delta > 1$ .

### 4.3 Linear convergence under error bound conditions

In this section we analyze the convergence of Algorithm 1.1 under error bound conditions. We improve and simplify the results given in [23]. In the rest of the section we assume H1 and H2. Moreover, we let  $X \subset H$ ,  $F = f + h$ ,  $F_* = \inf F$ , and suppose  $S_* := \text{argmin} F \neq \emptyset$ .

We consider the following condition, which was studied in [8] in connection with the proximal gradient method and is known as *Luo-Tseng error bound condition* [22].

EB For some  $c_{X, \Gamma^{-1}} > 0$ , we have

$$(\forall x \in X) \quad \text{dist}_{\Gamma^{-1}}(x, S_*) \leq c_{X, \Gamma^{-1}} \|x - \text{prox}_h^{\Gamma^{-1}}(x - \nabla^{\Gamma^{-1}} f(x))\|_{\Gamma^{-1}}. \quad (4.27)$$

**Remark 4.14.**

- (i) Another popular error bound condition is that of the metric subregularity of the subdifferential. More precisely,  $\partial^{\Gamma^{-1}} F$  is *2-metrically subregular* on  $X$  with respect to the metric  $\|\cdot\|_{\Gamma^{-1}}$  [8, 15] if for some  $\zeta_{X, \Gamma^{-1}} > 0$  the following holds

$$(\forall x \in X) \quad \text{dist}_{\Gamma^{-1}}(x, S_*) \leq \frac{1}{\zeta_{X, \Gamma^{-1}}} \text{dist}_{\Gamma^{-1}}(0, \partial^{\Gamma^{-1}} F(x)). \quad (4.28)$$

- (ii) EB and (4.28) are equivalent if  $h = 0$ , since in that case  $\text{prox}_h^{\Gamma^{-1}} = \text{Id}$  and  $c_{X, \Gamma^{-1}} = \zeta_{X, \Gamma^{-1}}^{-1}$ .

- (iii) Since  $\partial^{\Gamma^{-1}}F(x) = \Gamma\partial F(x)$  and  $\|\cdot\| \geq \gamma_{\min}^{1/2}\|\cdot\|_{\Gamma^{-1}}$ , it follows that if for every  $x \in X$ ,  $\text{dist}(x, S_*) \leq \zeta_{X,\text{Id}}^{-1}\text{dist}(0, \partial F(x))$ , then (4.28) holds with constant  $\zeta_{X,\Gamma^{-1}} = \gamma_{\min}\zeta_{X,\text{Id}}$ .
- (iv) [8, Theorem 3.5] yields that for any Hilbert norm  $\|\cdot\|_W$ ,  $\|x - \text{prox}_h^W(x - \nabla^W f(x))\|_W \leq \text{dist}_W(0, \partial^W F(x))$ . So, if EB holds on  $X$ , then (4.28) holds on  $X$  with  $\zeta_{X,\Gamma^{-1}} = c_{X,\Gamma^{-1}}^{-1}$ . In [8, Theorem 3.4-3.5] also the reverse implication was shown when  $f$  is Lipschitz smooth and  $X$  is a sublevel set of  $F$ .

**Remark 4.15.** In [8, Corollary 3.6] condition EB was shown to be equivalent to the following *quadratic growth condition* (also called 2-conditioning in [15])

$$(\exists \alpha_{X,\Gamma^{-1}} > 0)(\forall x \in X) \quad F(x) - \inf F \geq \frac{\alpha_{X,\Gamma^{-1}}}{2} \text{dist}_{\Gamma^{-1}}^2(x, S_*), \quad (4.29)$$

on every sublevel set  $X = \{x \in H \mid F(x) - F_* \leq r\}$ . Moreover, the relationships between the constants are  $c_{X,\Gamma^{-1}} = (1 + 2/\alpha_{X,\Gamma^{-1}})(1 + L_{\|\cdot\|_{\Gamma^{-1}}})$  and  $\alpha_{X,\Gamma^{-1}} < 1/c_{X,\Gamma^{-1}}$ . Finally, if the quadratic growth condition (4.29) holds, then (4.28) holds on  $X$ , with  $\zeta_{X,\Gamma^{-1}} = \alpha_{X,\Gamma^{-1}}/2$ .

We now analyze the convergence of Algorithm 1.1 under condition EB.

**Theorem 4.16.** *Under the assumptions of Theorem 4.9, suppose that  $S_* \neq \emptyset$  and that EB holds on a set  $X$  such that  $X \supset \{x^n \mid n \in \mathbb{N}\}$  P-a.s. with  $c_{X,\Gamma^{-1}} > 0$ . Then,*

$$(\forall n \in \mathbb{N}) \quad \mathbb{E}[F(x^n)] - F_* \leq \left(1 - \rho_{\min} \min\left\{1, \frac{2 - \delta}{2c_{X,\Gamma^{-1}}}\right\}\right)^n (\mathbb{E}[F(x^0)] - F_*). \quad (4.30)$$

Moreover, there exists a random variable  $x_*$  which takes values in  $S_*$  such that  $x^n \rightarrow x_*$  P-a.s. and  $\mathbb{E}[\|x^n - x_*\|_W] = O((1 - \rho_{\min} \min\{1, (2 - \delta)/(2c_{X,\Gamma^{-1}})\})^{n/2})$ .

*Proof.* Let  $n \in \mathbb{N}$  and  $x \in S_*$ . Then, (4.8) with  $\mu_{\Gamma^{-1}} = \sigma_{\Gamma^{-1}} = 0$ , yields

$$\begin{aligned} & \frac{1}{\rho_{\min}} \mathbb{E}[F(x^{n+1}) - F(x^n) \mid \mathfrak{E}_{n-1}] \\ & \leq \frac{\delta - 2}{2} \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}}^2 + \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}} \|x^n - x\|_{\Gamma^{-1}} - (F(x^n) - F_*). \end{aligned} \quad (4.31)$$

Since (4.31) holds for all  $x \in S_*$ , using EB, (4.15), and Proposition 4.5, we have

$$\begin{aligned} & \frac{2}{\rho_{\min}} \mathbb{E}[F(x^{n+1}) - F(x^n) \mid \mathfrak{E}_{n-1}] \\ & \leq (\delta - 2) \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}}^2 + 2 \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}} \text{dist}_{\Gamma^{-1}}(x^n, S_*) - 2(F(x^n) - F_*) \\ & \leq \frac{2(2c_{X,\Gamma^{-1}} + \delta - 2)_+}{\rho_{\min}(2 - \delta)} \mathbb{E}[F(x^n) - F(x^{n+1}) \mid \mathfrak{E}_{n-1}] - 2(F(x^n) - F_*) \quad \text{P a.s.}, \end{aligned} \quad (4.32)$$

which can be equivalently written as

$$\begin{aligned} & \mathbb{E}[F(x^{n+1}) - F_* \mid \mathfrak{E}_{n-1}] - (F(x^n) - F_*) \\ & \leq \left(\frac{2c_{X,\Gamma^{-1}}}{2 - \delta} - 1\right)_+ \left(F(x^n) - F_* - \mathbb{E}[F(x^{n+1}) - F_* \mid \mathfrak{E}_{n-1}]\right) - \rho_{\min}(F(x^n) - F_*). \end{aligned}$$

Therefore,

$$\max \left\{ 1, \frac{2c_{\mathcal{X}, \Gamma^{-1}}}{2 - \delta} \right\} \mathbb{E}[F(x^{n+1}) - F_* \mid \mathfrak{E}_{n-1}] \leq \left( \max \left\{ 1, \frac{2c_{\mathcal{X}, \Gamma^{-1}}}{2 - \delta} \right\} - \mathfrak{p}_{\min} \right) (F(x^n) - F_*),$$

which gives (4.30). Now we set  $\rho = 1 - \mathfrak{p}_{\min} \min \{1, (2 - \delta)/(2c_{\mathcal{X}, \Gamma^{-1}})\}$  and  $\theta = \mathfrak{p}_{\min}(2 - \delta)/2$ . Then, Jensen inequality, (4.16), and (4.30) yield

$$\mathbb{E}[\|x^n - x^{n+1}\|_{\mathcal{W}}] \leq \theta^{-1/2} \sqrt{\mathbb{E}[F(x^n)] - F_*} \leq \theta^{-1/2} \rho^{n/2} \sqrt{\mathbb{E}[F(x^0)] - F_*}. \quad (4.33)$$

Therefore, since  $\rho^{1/2} < 1$ , we have  $\mathbb{E}[\sum_{n \in \mathbb{N}} \|x^n - x^{n+1}\|_{\mathcal{W}}] = \sum_{n \in \mathbb{N}} \mathbb{E}[\|x^n - x^{n+1}\|_{\mathcal{W}}] < +\infty$ . Hence  $\sum_{n \in \mathbb{N}} \|x^n - x^{n+1}\|_{\mathcal{W}} < +\infty$  P-a.s., which means that  $(x^n)_{n \in \mathbb{N}}$  is a Cauchy sequence P-a.s. Now, Theorem 4.9(ii) yields that there exists a random variable  $x_*$  with values in  $S_*$  such that  $x^n \rightarrow x_*$  P-a.s. Therefore,  $x^n \rightarrow x_*$  P-a.s. Finally, let  $n \in \mathbb{N}$ . Then, for every  $p \in \mathbb{N}$ ,

$$\|x^n - x^{n+p}\|_{\mathcal{W}} \leq \sum_{i=0}^{p-1} \|x^{n+i} - x^{n+i+1}\|_{\mathcal{W}} \leq \sum_{i=0}^{+\infty} \|x^{n+i} - x^{n+i+1}\|_{\mathcal{W}}.$$

Hence, letting  $p \rightarrow +\infty$ , we have  $\mathbb{E}[\|x^n - x_*\|_{\mathcal{W}}] \leq \sum_{i=0}^{+\infty} \mathbb{E}[\|x^{n+i} - x^{n+i+1}\|_{\mathcal{W}}]$ . Therefore, it follows from (4.33) that

$$\mathbb{E}[\|x^n - x_*\|_{\mathcal{W}}] \leq \theta^{-1/2} \sqrt{\mathbb{E}[F(x^0)] - F_*} \sum_{i=0}^{+\infty} \rho^{(n+i)/2} = \theta^{-1/2} \sqrt{\mathbb{E}[F(x^0)] - F_*} \frac{\rho^{n/2}}{1 - \rho^{1/2}}. \quad \square$$

**Remark 4.17.**

- (i) The rate given in Theorem 4.16 matches the one given in [8, Theorem 3.2] for the deterministic case ( $\mathfrak{p}_{\min} = 1$ ).
- (ii) In Theorem 4.16, the constant  $c_{\mathcal{X}, \Gamma^{-1}}$  depends on the stepsizes  $\gamma_i$ 's which in turn depend on  $\delta$  (usually  $\gamma_i = \delta/\nu_i$  with  $0 < \delta < 2$ ). Therefore, the optimal value of  $\delta$  in the rate (4.30) can be determined after specifying the expression of  $c_{\mathcal{X}, \Gamma^{-1}}$ . We did so in the special application of Section 5.2.
- (iii) In [23, Definition 5.2], in relation to Algorithm 1.1 but with uniform block sampling and assuming S3 and  $R_{\Gamma^{-1}}(x^0) := \sup_{x \in \{F \leq F(x^0)\}} \text{dist}_{\Gamma^{-1}}(x, S_*) < +\infty$ , the following error bound condition is considered

$$\text{dist}_{\Gamma^{-1}}(x, S_*) \leq (\kappa_{1, \mathcal{X}, \Gamma^{-1}} + \kappa_{2, \mathcal{X}, \Gamma^{-1}} \text{dist}_{\Gamma^{-1}}^2(x, S_*)) \|x - \text{prox}_{\mathfrak{h}}^{\Gamma^{-1}}(x - \nabla^{\Gamma^{-1}} f(x))\|_{\Gamma^{-1}}, \quad (4.34)$$

for some constants  $\kappa_{1, \mathcal{X}, \Gamma^{-1}} > 0$  and  $\kappa_{2, \mathcal{X}, \Gamma^{-1}} \geq 0$ . The authors show several examples in which such condition is satisfied with  $\mathcal{X} = \text{dom } \mathfrak{h}$  and possibly  $\kappa_{2, \mathcal{X}, \Gamma^{-1}} > 0$ . The above error bound looks more general than EB. However, for the purpose of analyzing Algorithm 1.1 and under the assumptions considered in [23] this is not the case. Indeed, in [23, equation (3.11)] it was shown that the algorithm is descending almost surely<sup>2</sup>, so  $\{x^n \mid n \in \mathbb{N}\} \subset \{F \leq F(x^0)\}$  P-a.s. Therefore, since  $\sup_{x \in \{F \leq F(x^0)\}} \text{dist}_{\Gamma^{-1}}(x, S_*) = R_{\Gamma^{-1}}(x^0) < +\infty$ , if (4.34) holds on a set  $\mathcal{X}$  containing P-a.s. the set  $\{x^n \mid n \in \mathbb{N}\}$ , then EB holds on  $\mathcal{X}' := \mathcal{X} \cap \{F \leq F(x^0)\} \supset \{x^n \mid n \in \mathbb{N}\}$

<sup>2</sup>Alternatively, note that S3 implies S2 which in turn, in view of Proposition 4.7, ensures the descending property.

P-a.s. with  $c_{\mathcal{X}', \Gamma-1} := \kappa_{1, \mathcal{X}, \Gamma-1} + \kappa_{2, \mathcal{X}, \Gamma-1} R_{\Gamma-1}(x^0)^2$ . Thus, Theorem 4.16 applies accordingly. Moreover, [23, Theorem 5.5] gives the linear rate

$$\begin{aligned} \mathbb{E}[F(x^n)] - F_* &\leq \left(1 - \frac{1}{1 + \bar{c}}\right)^n (\mathbb{E}[F(x^0)] - F_*), \quad \text{where} \\ \bar{c} &= \frac{1}{\rho} \left(2 + \frac{2c_{\mathcal{X}', \Gamma-1}}{\sqrt{\rho}} + (1 - \rho) \frac{c_{\mathcal{X}', \Gamma-1}^2}{\rho} + 2c_{\mathcal{X}', \Gamma-1} + 1 - \rho\right). \end{aligned}$$

Then, we have  $1 + \bar{c} \geq (3 + 4c_{\mathcal{X}', \Gamma-1})/\rho$  and hence

$$\frac{1}{1 + \bar{c}} \leq \frac{\rho}{3 + 4c_{\mathcal{X}', \Gamma-1}} < \frac{\rho}{\max\{1, 2c_{\mathcal{X}', \Gamma-1}\}}.$$

This shows that Theorem 4.16 improves the rate in [23, Theorem 5.5]. Moreover, the analysis given here, relying on Proposition 4.6, relaxes the assumptions and is significantly simpler.

- (iv) It follows from [23, Theorem 6.8] that if  $f$  is a quadratic function and  $h$  is an indicator function of a polyhedral set, then (4.34) is satisfied on  $\text{dom } h$ . Therefore, if  $\text{dom } h$  is bounded, then EB holds on  $\mathcal{X} = \text{dom } h$  with  $c_{\mathcal{X}, \Gamma-1} := \kappa_{1, \mathcal{X}, \Gamma-1} + \kappa_{2, \mathcal{X}, \Gamma-1} \text{diam}_{\Gamma-1}^2(\text{dom } h)$  and Theorem 4.16 can be applied, since  $\{x^n \mid n \in \mathbb{N}\} \subset \text{dom } h$ .
- (v) Several works address the convergence of random coordinate descent methods under error bound conditions. We mention [16] which considers a serial sampling and stepsizes set according to the global Lipschitz constant of  $\nabla f$  and [14], which analyzes restarting procedures for accelerated and parallel coordinate descent methods using assumptions S1 and (4.29).

**Remark 4.18.** Often error bound conditions or quadratic growth conditions are satisfied when  $\mathcal{X}$  is a sublevel set (see Remark 4.15). So, in such scenarios, apart when  $\text{dom } h$  is a sublevel set of  $F$ , in order to fulfill the assumption  $\mathcal{X} \supset \{x^n \mid n \in \mathbb{N}\}$  P-a.s. in Theorem 4.16, it is desirable for Algorithm 1.1 to be (a.s.) descending. This occurs if condition S2 holds (Proposition 4.7), whereas, in general, S1 does not guarantee any such descending property. However, especially when  $\eta \ll m$ , condition S2 may be much more restrictive than S1, thus leading to a significant reduction of the stepsizes, which ultimately slows down the convergence. The next result shows that Algorithm 1.1 can be slightly modified so to ensure the descending property while keeping the validity of Theorem 4.16.

**Theorem 4.19.** *Let H1–H3 be satisfied. Let  $(\nu_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  and suppose that S1 holds. Suppose in addition that  $S_* \neq \emptyset$ , and that EB holds on the set  $\mathcal{X} = \{F \leq F(x^0)\}$  with  $c_{\mathcal{X}, \Gamma-1} > 0$ . Let  $(x^n)_{n \in \mathbb{N}}$  be generated by the following variation of Algorithm 1.1*

$$\begin{array}{l} \text{for } n = 0, 1, \dots \\ \quad \left[ \begin{array}{l} \text{for } i = 1, \dots, m \\ \quad \left[ \tilde{x}_i^{n+1} = x_i^n + \varepsilon_i^n (\text{prox}_{\gamma_i h_i}(x_i^n - \gamma_i \nabla_i f(x^n)) - x_i^n) \right. \\ \quad \text{if } F(\tilde{x}^{n+1}) \leq F(x^n) \\ \quad \quad x^{n+1} = \tilde{x}^{n+1} \\ \quad \text{else} \\ \quad \quad x^{n+1} = x^n. \end{array} \right. \end{array} \quad (4.35)$$

Then the conclusions of Theorem 4.16 still hold.

*Proof.* It follows from the definition of  $x^{n+1}$  that  $F(x^{n+1}) \leq F(x^n)$ . Therefore  $\mathsf{X} \supset \{x^n \mid n \in \mathbb{N}\}$ . Recalling (4.2) and (4.3), algorithm (4.35) can be alternatively written as

$$\begin{cases} \text{for } n = 0, 1, \dots \\ \quad \text{for } i = 1, \dots, m \\ \quad \quad \bar{x}_i^{n+1} = \text{prox}_{\gamma_i h_i}(x_i^n - \gamma_i \nabla_i f(x^n)), \\ \quad \tilde{x}^{n+1} = x^n + \varepsilon^n \odot (\bar{x}^{n+1} - x^n) \\ \quad x^{n+1} = \tilde{x}^{n+1} \mathbf{1}_{\{F(\tilde{x}^{n+1}) \leq F(x^n)\}} + x^n \mathbf{1}_{\{F(\tilde{x}^{n+1}) > F(x^n)\}} \end{cases} \quad (4.36)$$

and we have  $F(x^{n+1}) \leq F(\tilde{x}^{n+1})$ . Then we can essentially repeat the argument in the proof of Theorem 4.16. First we note that (4.8) and hence (4.31) holds with  $x^{n+1}$  replaced by  $\tilde{x}^{n+1}$ . This follows from the definition of  $\tilde{x}^{n+1}$ . Moreover, also (4.13) holds with  $x^{n+1}$  replaced by  $\tilde{x}^{n+1}$  and hence we derive (with  $x = x^n$  and  $\sigma_{\Gamma^{-1}} = 0$ ) that

$$(2 - \delta) \|\bar{x}^{n+1} - x^n\|_{\Gamma^{-1}}^2 \leq \frac{2}{\rho_{\min}} \mathbb{E}[F(x^n) - F(\tilde{x}^{n+1}) \mid \mathfrak{E}_{n-1}]. \quad (4.37)$$

Then, we have

$$\begin{aligned} & \frac{2}{\rho_{\min}} \mathbb{E}[F(\tilde{x}^{n+1}) - F(x^n) \mid \mathfrak{E}_{n-1}] \\ & \leq (\delta - 2) \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}}^2 + 2 \|x^n - \bar{x}^{n+1}\|_{\Gamma^{-1}} \text{dist}_{\Gamma^{-1}}(x^n, \mathsf{S}_*) - 2(F(x^n) - F_*) \\ & \leq \frac{2(2c_{\mathsf{X}, \Gamma^{-1}} + \delta - 2)_+}{\rho_{\min}(2 - \delta)} \mathbb{E}[F(x^n) - F(\tilde{x}^{n+1}) \mid \mathfrak{E}_{n-1}] - 2(F(x^n) - F_*) \quad \text{P a.s.} \end{aligned} \quad (4.38)$$

and hence, since  $F(x^{n+1}) \leq F(\tilde{x}^{n+1})$ , (4.32) still holds (for the new definition of  $x^{n+1}$ ). Thus, (4.30) follows. As for the second part of the statement, we note that, since  $F(x^{n+1}) \leq F(\tilde{x}^{n+1})$ , by (4.37), we have  $(2 - \delta) \|\bar{x}^{n+1} - x^n\|_{\Gamma^{-1}}^2 \leq (2/\rho_{\min}) \mathbb{E}[F(x^n) - F(x^{n+1}) \mid \mathfrak{E}_{n-1}]$ . Moreover, it follows from the definitions of  $x^{n+1}$  and  $\tilde{x}^{n+1}$  in algorithm (4.36) that

$$x^{n+1} - x^n = (\tilde{x}^{n+1} - x^n) \mathbf{1}_{\{F(\tilde{x}^{n+1}) \leq F(x^n)\}}$$

and hence, by Proposition 4.5, we have

$$\mathbb{E}[\|x^{n+1} - x^n\|_{\mathbb{W}}^2 \mid \mathfrak{E}_{n-1}] \leq \mathbb{E}[\|\tilde{x}^{n+1} - x^n\|_{\mathbb{W}}^2 \mid \mathfrak{E}_{n-1}] = \|\bar{x}^{n+1} - x^n\|_{\Gamma^{-1}}^2.$$

In the end (4.16) with  $\sigma_{\Gamma^{-1}} = 0$  still holds and the proof can continue as in that of Theorem 4.16.  $\square$

## 5 Applications

In this section we show some relevant optimization problems for which the theoretical analysis of Algorithm 1.1 can be particularly useful.

### 5.1 The Lasso problem

Many papers study the convergence of coordinate descent methods for the Lasso problem and recent works prove linear convergence (see e.g., [16, 23, 25]). In [16] a random serial update of

blocks is considered while in [25] the general framework of feasible descend methods is analyzed which include (nonrandom) cyclic coordinate methods. In the following we discuss our contribution comparing with [23]. Let  $A \in \mathbb{R}^{p \times m}$  and  $\mathbf{b} \in \mathbb{R}^p$ . We consider the problem

$$\min_{\mathbf{x} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (\lambda > 0). \quad (5.1)$$

We denote by  $\mathbf{a}^i$  and  $\mathbf{a}_k$  the  $i$ -th column and  $k$ -th row of  $A$  respectively. Since

$$\frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 = \frac{1}{2} \sum_{k=1}^p (\langle \mathbf{a}_k, \mathbf{x} \rangle - \mathbf{b}_k)^2 = \frac{1}{2} \sum_{k=1}^p \left( \sum_{i=1}^m \mathbf{a}_k^i x_i - \mathbf{b}_k \right)^2,$$

H4 holds and  $\eta = \max_{1 \leq k \leq p} \text{card}(\text{spt}(\mathbf{a}_k))$ . Moreover, since  $\nabla_i f(\mathbf{x}) = \langle \mathbf{a}^i, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle$  and recalling Remark 3.2(i), conditions L1 and L2 are satisfied with  $L_i = \|\mathbf{a}^i\|^2$  and  $L^{(k)} = \|\mathbf{a}_k\|^2$  respectively. Then, Algorithm 1.1 (assuming that each block is made of one coordinate only) writes as

$$x^{n+1} = x^n + \sum_{i=1}^m \varepsilon_i^n [\text{soft}_{\gamma_i \lambda}(x_i^n - \gamma_i \mathbf{a}^{i \top} (\mathbf{A}x^n - \mathbf{b})) - x_i^n] \mathbf{e}_i, \quad (5.2)$$

where the *soft thresholding operator*  $\text{soft}_{\gamma_i \lambda}$  is defined as  $\text{soft}_{\gamma_i \lambda}(t) = \text{sign}(t) \max\{0, |t| - \gamma_i \lambda\}$  and  $(\mathbf{e}_i)_{1 \leq i \leq m}$  is the canonical basis of  $\mathbb{R}^m$ . Now, define  $u^n = \mathbf{A}x^n - \mathbf{b}$ . Then multiplying the equation in (5.2) by  $A$  and subtracting  $\mathbf{b}$  by both terms, the algorithm is equivalently written as

$$\begin{cases} \text{for } i \in \text{spt}(\varepsilon^n) \\ \left[ \begin{array}{l} \xi_i = \text{soft}_{\gamma_i \lambda}(x_i^n - \gamma_i \mathbf{a}^{i \top} u^n) - x_i^n, \\ x^{n+1} = x^n + \sum_{i \in \text{spt}(\varepsilon^n)} \xi_i \mathbf{e}_i \\ u^{n+1} = u^n + \sum_{i \in \text{spt}(\varepsilon^n)} \xi_i \mathbf{a}^i, \end{array} \right. \end{cases} \quad (5.3)$$

showing that each iteration costs  $O(p\tau_{\max})$  multiplications, where  $\tau_{\max}$  is the maximum number of block updates per iteration. We now address the determination of the smoothness parameters  $(\nu_i)_{1 \leq i \leq m}$ . We first give a general rule which holds for any arbitrary sampling. Recalling Theorem 3.1(iii) and Remark 3.2(i) and noting that  $\{k | i \in I_k\} = \{k | i \in \text{spt}(\mathbf{a}_k)\} = \text{spt}(\mathbf{a}^i)$ , if we set  $\nu_i = \sum_{k \in \text{spt}(\mathbf{a}^i)} \|\mathbf{a}_k\|^2$ , then S3 and hence S2 holds. This choice was considered in [23]. Moreover, according to the discussion at the beginning of Section 3.2 other options for satisfying S2 are  $\nu_i = \min\{\eta, \tau_{\max}\} \|\mathbf{a}^i\|^2$  or  $\nu_i = \sum_{k \in \text{spt}(\mathbf{a}^i)} \min\{\text{card}(\text{spt}(\mathbf{a}_k)), \tau_{\max}\} (\mathbf{a}_k^i)^2$ . This latter choice is better than the second one and, if we assume that the nonzero entries of  $A$  are about of the same magnitude, it is also better than the first one. Next, we face the special case of the  $\tau$ -nice sampling which allows to reduce the  $\nu_i$ 's while satisfying S1. Recalling the corresponding discussion in Section 3.2, we have the following alternatives: (1) set, for every  $i \in [m]$ ,  $\nu_i = (1 + (\tau - 1)(\eta - 1)/(m - 1)) \|\mathbf{a}^i\|^2$ ; (2) set for every  $i \in [m]$ ,  $\nu_i = \sum_{k \in \text{spt}(\mathbf{a}^i)} (1 + (\tau - 1)(\text{card}(\text{spt}(\mathbf{a}_k)) - 1)/(m - 1)) (\mathbf{a}_k^i)^2$ . Finally, we make few remarks on the convergence properties of algorithm (5.3). Since the objective function in (5.1) satisfies a quadratic growth condition on its sublevel sets [15, Example 3.8], then Remark 4.15, Remark 4.18, and Theorem 4.16, yield linear convergence of algorithm (5.3) provided that S2 holds. Whereas Theorem 4.19 ensures that if we modify algorithm (5.3) so that we accept the next iterate  $x^{n+1}$  only if  $\|u^{n+1}\|^2 + 2\lambda \|x^{n+1}\|_1 \leq \|u^n\|^2 + 2\lambda \|x^n\|_1$ , then the resulting algorithm converges linearly under condition S1. If the violation of the monotonicity condition above occurs few times along all the iterations, this modification does not increase much the computational cost of the algorithm (see also Section 6). We stress that both Theorem 4.16 and Theorem 4.19 ensure also almost sure and linear convergence in mean of the iterates of (5.2). This latter result is new and is especially relevant in this context, since the iterates carry sparsity information.

## 5.2 Computing the minimal norm solution of a linear system

Let  $A \in \mathbb{R}^{m \times p}$  and  $\mathbf{b} \in R(A)$  (the range of  $A$ ). Let us consider the problem

$$\underset{\substack{\mathbf{x} \in \mathbb{R}^p \\ A\mathbf{x} = \mathbf{b}}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x}\|^2. \quad (5.4)$$

Here, we denote by  $\mathbf{a}_i \in \mathbb{R}^p$  and  $\mathbf{a}^k \in \mathbb{R}^m$  the  $i$ -th row and the  $k$ -th column of  $A$ . The dual problem is

$$\underset{\mathbf{u} \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} \|A^\top \mathbf{u}\|^2 - \langle \mathbf{u}, \mathbf{b} \rangle := \mathcal{D}(\mathbf{u}), \quad (5.5)$$

which is a smooth convex optimization problem. Moreover, if  $\mathbf{x}_*$  is the solution of (5.4) and  $\mathbf{x} = A^\top \mathbf{u}$  (the primal-dual relationship), then we have

$$\frac{1}{2} \|\mathbf{x} - \mathbf{x}_*\|^2 \leq \mathcal{D}(\mathbf{u}) - \inf \mathcal{D}. \quad (5.6)$$

Then, the dual problem is clearly of the form (1.1), with  $\mathbf{h} = 0$ , and **H4** and **L1** are satisfied, assuming that each block is made of one coordinate only, with  $L_i = \|\mathbf{a}_i\|^2$  and  $\eta = \max_{1 \leq k \leq p} \text{card}(\text{spt}(\mathbf{a}^k))$ . So, Algorithm 1.1 applied to (5.5), turns into

$$\mathbf{u}^{n+1} = \mathbf{u}^n - \sum_{i=1}^m \varepsilon_i^n \gamma_i (\langle \mathbf{a}_i, A^\top \mathbf{u}^n \rangle - \mathbf{b}_i) \mathbf{e}_i.$$

Now, setting  $\mathbf{x}^n = A^\top \mathbf{u}^n$  and multiplying the above equality by  $A^\top$ , we have

$$\mathbf{x}^{n+1} = \mathbf{x}^n - \sum_{i=1}^m \varepsilon_i^n \gamma_i (\langle \mathbf{a}_i, \mathbf{x}^n \rangle - \mathbf{b}_i) \mathbf{a}_i. \quad (5.7)$$

Since,  $\mathbf{b} \in R(A)$ , it is easy to see, through a singular value decomposition of  $A$ , that, for every  $\mathbf{u} \in \mathbb{R}^m$ ,  $\sigma_{\min}^2(A) \text{dist}(\mathbf{u}, \text{argmin } \mathcal{D}) \leq \|\nabla \mathcal{D}(\mathbf{u})\|$  (where  $\sigma_{\min}(A)$  is the minimum singular value of  $A$ ) [15, Example 3.6]. So, in view of Remark 4.14(ii)-(iii), **EB** is satisfied on the entire space with constant  $c_{\mathbb{R}^m, \Gamma^{-1}} = (\gamma_{\min} \sigma_{\min}^2(A))^{-1}$ . Therefore, if, for every  $i \in [m]$ ,  $\gamma_i = \delta / (\beta_{1,i} \|\mathbf{a}_i\|^2)$  with  $0 < \delta < 2$ , Theorem 4.16 and (5.6) ensure the linear convergence of the iterates  $\mathbf{x}^n$ 's towards the solution of (5.4) with rate  $(1 - \rho_{\min} \min \{1, \gamma_{\min} \sigma_{\min}^2(A) (2 - \delta) / 2\})^{1/2}$ . We remark that (5.7) is nothing but a stochastic gradient descent algorithm on the problem

$$\underset{\mathbf{x} \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{2} \sum_{i=1}^m (\langle \mathbf{a}_i, \mathbf{x} \rangle - \mathbf{b}_i)^2.$$

Since  $\|A\mathbf{x} - \mathbf{b}\|^2 \leq \|A\|^2 \|\mathbf{x} - \mathbf{x}_*\|^2$ , we have then showed the linear convergence rate

$$\frac{1}{2} \|A\mathbf{x}^n - \mathbf{b}\|^2 - \frac{1}{2} \|A\mathbf{x}_* - \mathbf{b}\|^2 = O\left(\sigma_{\max}^2(A) \left(1 - \rho_{\min} \min \left\{1, \frac{\sigma_{\min}^2(A) \delta (2 - \delta)}{2 \max_i \beta_{1,i} \|\mathbf{a}_i\|^2}\right\}\right)^n\right),$$

of the stochastic gradient descent with arbitrary and possibly variable batch size for least squares problems. This also shows that the best rate is achieved for  $\delta = 1$ . We finally note that in the serial case, that is, if for every  $n \in \mathbb{N}$   $\text{spt}(\varepsilon^n) = \{i_n\}$ , multiplying equation (5.7) by  $\mathbf{a}_{i_n}^\top$ , we have

$$\langle \mathbf{a}_{i_n}, \mathbf{x}^{n+1} \rangle = \langle \mathbf{a}_{i_n}, \mathbf{x}^n \rangle - \gamma_{i_n} (\langle \mathbf{a}_{i_n}, \mathbf{x}^n \rangle - \mathbf{b}_{i_n}) \|\mathbf{a}_{i_n}\|^2.$$



Therefore, since in this case  $\beta_{1,i} = \beta_2 = 1$ , we can chose the stepsizes such that  $\gamma_i \|a_i\|^2 = 1$  (so that  $\delta = 1$ ) and hence  $x^{n+1}$  is a solution of the  $i_n$ -th equation of the linear system  $Ax = b$ . Moreover,  $x^{n+1}$  is the projection of  $x^n$  onto the affine space defined by the equation  $a_{i_n}x = b_{i_n}$  [41]. Thus, this method is nothing but the *randomized Kaczmarz method* [37] and we proved linear convergence for general probabilities  $p_i$ 's, although the constants we derive are not optimal (see [19, 37, 41]).

### 5.3 Regularized empirical risk minimization

Let  $H$  be a separable real Hilbert space. Regularized empirical risk estimation solves the following optimization problem

$$\underset{w \in H}{\text{minimize}} \quad \frac{1}{\lambda m} \sum_{i=1}^m \ell(y_i, \langle w, x_i \rangle) + \frac{1}{2} \|w\|^2 := \mathcal{P}(w), \quad (5.8)$$

where  $(x_i, y_i)_{1 \leq i \leq m}$  is the training set (input-output pairs),  $\ell: \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}_+$ ,  $\mathcal{Y} \subset \mathbb{R}$ , is the *loss* function, which is convex in the second variable, and  $\lambda > 0$  is a regularization parameter. The dual problem of (5.8) is

$$\underset{u \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} u^\top K u + \frac{1}{\lambda m} \sum_{i=1}^m \ell^*(y_i, -u_i \lambda m) := \mathcal{D}(u), \quad (5.9)$$

where  $\ell^*(y_i, \cdot)$  is the Fenchel conjugate of  $\ell(y_i, \cdot)$  and  $K = XX^\top \in \mathbb{R}^{m \times m}$  is the Gram matrix of  $(x_i)_{1 \leq i \leq m}$ . Moreover, the solutions  $(\bar{w}, \bar{u})$  of the primal and dual problems are characterized by the following KKT conditions

$$\begin{cases} \bar{w} = X^\top \bar{u} = \sum_{i=1}^m \bar{u}_i x_i, \\ \forall i \in \{1, \dots, m\} \quad -\bar{u}_i m \lambda \in \partial \ell(y_i, \langle x_i, \bar{w} \rangle), \end{cases} \quad (5.10)$$

where  $\partial \ell(y_i, \cdot)$  is the subdifferential of  $\ell(y_i, \cdot)$ . Note also that the first of (5.10) gives the link between the dual and the primal variable and, if  $w = X^\top u$ , then it holds  $(1/2) \|w - \bar{w}\|^2 \leq \mathcal{D}(u) - \inf \mathcal{D}$ . Now, the dual problem (5.9) is of the form (1.1) and hence Algorithm 1.1 can be applied. The following examples give implementation details for two specific losses.

**Example 5.1** (Ridge regression). The least squares loss is  $\ell(s, t) = (1/2)|s - t|^2$ . Then  $\ell^*(s, r) = (1/2)r^2 + rs$  and, in this case, (5.9) reduces to

$$\underset{u \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} u^\top (K + \lambda m \text{Id}) u - y^\top u := \mathcal{D}(u)$$

which is strongly convex with modulus  $\lambda m$  and has solution  $\bar{u} = (K + \lambda m \text{Id})^{-1} y$ . Since  $\mathcal{D}$  is smooth and  $\nabla_i \mathcal{D}(u) = e_i^\top (K + \lambda m \text{Id}) u - y_i$ , conditions H4 and L1 hold with  $L_i = \|x_i\|^2 + \lambda m = K_{i,i} + \lambda m$  and Algorithm 1.1 (with  $h = 0$ ) becomes

$$u^{n+1} = u^n - \sum_{i=1}^m \varepsilon_i^n \gamma_i (e_i^\top K u^n + \lambda m u_i^n - y_i) e_i. \quad (5.11)$$

Moreover, multiplying (5.11) by  $X^\top$ , defining  $w^n = X^\top u^n$ , and recalling that  $K = XX^\top$ , we have

$$w^{n+1} = w^n - \sum_{i=1}^m \varepsilon_i^n \gamma_i (\langle w^n, x_i \rangle - y_i) x_i - \lambda m \sum_{i=1}^m \varepsilon_i^n \gamma_i u_i^n x_i. \quad (5.12)$$

Note that, since the dual problem is strongly convex with modulus  $\lambda m$ , then it follows from Theorem 4.10, Remark 4.12, and Theorem C.1(i) that, setting, for every  $i \in [m]$ ,  $\nu_i \geq \beta_{1,i}(K_{i,i} + \lambda m)$  and  $\gamma_i = 1/\nu_i$ , we have

$$E[\mathcal{P}(w^n)] - \inf \mathcal{P} \leq \left(1 + \frac{\|K\|}{\lambda m}\right) \left(1 - \rho_{\min} \frac{2\lambda m}{\nu_{\max} + \lambda m}\right)^n \text{const.}$$

Now, we compare algorithm (5.12) with the stochastic gradient descent on problem (5.8). Assume that  $P(\sum_{i=1}^m \varepsilon_i^n = \tau) = 1$  for some  $\tau \in [m]$ . Then,  $\rho_{\min} = \tau/m$  and we can take  $\zeta_i \leq (K_{i,i} + \lambda m)^{-1}$ , and set  $\nu_i = \tau/\zeta_i$  and  $\gamma_i = 1/\nu_i$ , so that algorithm (5.12) turns into

$$w^{n+1} = w^n - \sum_{i=1}^m \varepsilon_i^n \frac{\zeta_i}{\tau} (\langle w^n, x_i \rangle - y_i) x_i - \lambda m \sum_{i=1}^m \varepsilon_i^n \frac{\zeta_i}{\tau} u_i^n x_i. \quad (5.13)$$

If we apply stochastic gradient descent with batch size  $\tau \in [m]$  and stepsize  $\zeta > 0$  directly on the primal problem (5.8) (multiplied by  $\lambda m$ ), and recalling that  $w^n = \sum_{i=1}^m u_i^n x_i$ , we have

$$w^{n+1} = w^n - \frac{\zeta}{\tau} \sum_{i=1}^m \varepsilon_i^n (\langle w^n, x_i \rangle - y_i) x_i - \lambda m \frac{\zeta}{m} \sum_{i=1}^m u_i^n x_i. \quad (5.14)$$

Then, comparing (5.13) and (5.14) we see that, provided that  $\zeta_i = \zeta$  for every  $i \in [m]$ , they only differ for the replacement  $(1/m) \sum_{i=1}^m u_i^n x_i \leftrightarrow (1/\tau) \sum_{i=1}^m \varepsilon_i^n u_i^n x_i$ . We stress that the stepsize  $\zeta$  in the stochastic gradient descent algorithm (5.14) is normally set according to the spectral norm of  $K + \lambda m \text{Id}$ , which may be difficult to compute. On the contrary in algorithm (5.13) the stepsizes  $\zeta_i$ 's are simply set as  $\zeta_i \leq 1/(K_{i,i} + \lambda m)$ , so they allow possibly much longer steps and also do not require any SVD computation.

**Example 5.2** (Support vector machines). The hinge loss is  $\ell(s, t) = (1 - st)_+$ . Then we have  $\ell^*(s, r) = r + \iota_{[0,1]}(sr)$  and the dual problem (5.9) is

$$\underset{u \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} u^\top K u - y^\top u + \iota_{(\lambda m)^{-1}[0,1]^m}(y \odot u). \quad (5.15)$$

Then Algorithm 1.1 on the dual turns into a parallel random block-coordinate projected gradient descent method. Moreover, it follows from Remark 4.17(iv) that the objective in (5.15) satisfies EB on its domain. Therefore, it follows from Theorem 4.16, Theorem 3.1(i), and Theorem C.1(ii) that  $E[\mathcal{P}(w^n)] - \inf \mathcal{P}$  converges linearly to zero, provided that, for all  $i \in [m]$ ,  $\nu_i \geq \beta_{1,i} K_{ii}$  and  $\gamma_i < 2/\nu_i$ .

## 6 Numerical Experiments

In this section we consider a Lasso problem, that is,

$$\min_{x \in \mathbb{R}^m} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \quad (\lambda > 0), \quad (6.1)$$

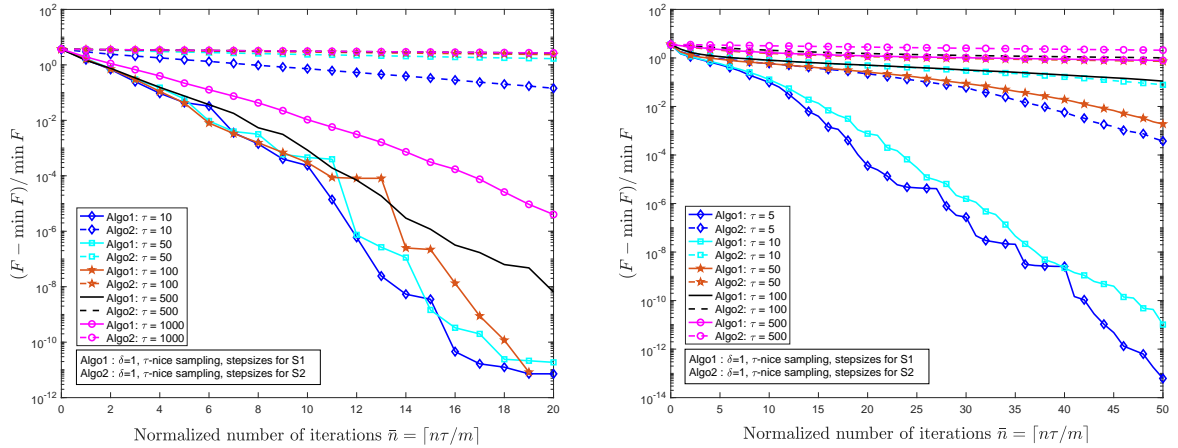


Figure 1: Comparison between **S1** and **S2** for the the stepsizes selection in a Lasso problem. Left:  $5 \cdot 10^4$  equations in  $10^5$  unknowns; degree of partial separability  $\eta = 148$ . Right:  $10^3$  equations in  $5 \cdot 10^3$  unknowns; degree of partial separability  $\eta = 563$ .

where  $A \in \mathbb{R}^{p \times m}$  is generated with random entries uniformly distributed in  $[-1, 1]$  so that each row is sparse and  $b = A\bar{x} + 0.06 \cdot \alpha$  with  $\alpha \sim N(0, \text{Id}_p)$  and  $\bar{x}$  a sparse vector in  $\mathbb{R}^m$ . We implement Algorithm 1.1 with  $\gamma_i = \delta/\nu_i$  ( $0 < \delta < 2$ ) and a  $\tau$ -nice uniform sampling, as described in Section 5.1. We present two experiments. The first compares conditions **S1** and **S2** for the determination of the stepsizes  $\gamma_i$ . The second one investigates the role played by  $\delta$ . In all the experiments we empirically checked that the algorithm is essentially descending in the sense that during the iterations there are very few violations of the descent property and with low magnitude. So, since the objective function in (6.1) satisfies **EB** on the sublevel sets, in virtue of Theorem 4.19, linear convergence holds.

### Condition **S1** vs **S2** and the effectiveness of the parallel strategy.

We compare the conditions **S1** and **S2** for the stepsizes selection and we checked the critical role played by **S1** for the effectiveness of the parallel strategy on problems with sparse structure. Here we set  $\delta = 1$ . In Figure 1, Algo1 uses smoothness parameters specifically designed for the  $\tau$ -nice sampling, that is,  $\nu_i = \beta_1 \|a^i\|^2$  with  $\beta_1 = 1 + (\tau - 1)(\eta - 1)/(m - 1)$  (making **S1** satisfied), while Algo2 uses a more conservative choice for the smoothness parameters which is valid for any sampling updating a maximum of  $\tau$  blocks per iteration, that is  $\nu_i = \beta_2 \|a^i\|^2$  with  $\beta_2 = \min\{\tau, \eta\}$  (making also **S2** satisfied). In the left diagram, we considered a large scale setting with  $\eta \ll m$ . In that case  $\beta_1$  may be much smaller than  $\beta_2$  leading to significantly larger stepsizes. Moreover, and more importantly, we note that as long as  $\tau$  is small enough the behavior of Algo1 does not depend on  $\tau$  (indeed  $\tau = 10, 50, 100$  perform equally well), whereas this is not true for Algo2. This feature of **S1**, first noted in [34], has been already discussed after Theorem 4.9 and is at the basis of the effectiveness of the parallel strategy. Indeed in the small- $\tau$  regime described above, the various versions of Algo1 depicted in Figure 1(left) have the same total computation cost ( $\bar{n}m$  block-coordinate updates), but the parallel implementation on  $\tau$  cores is  $\tau$  times faster than the serial one ( $\tau = 1$ ). Finally, in the right diagram of Figure 1 we show a scenario in which  $\eta/m$  is larger (the problem is less sparse). In such situation we see that the difference between the two stepsize selection criteria is less evident for  $\tau \geq 50$ . Moreover, Algo1 is more sensitive to  $\tau$  (compare  $\tau = 10$  and  $\tau = 50$ ), so that the benefit of the parallel scheme is reduced.

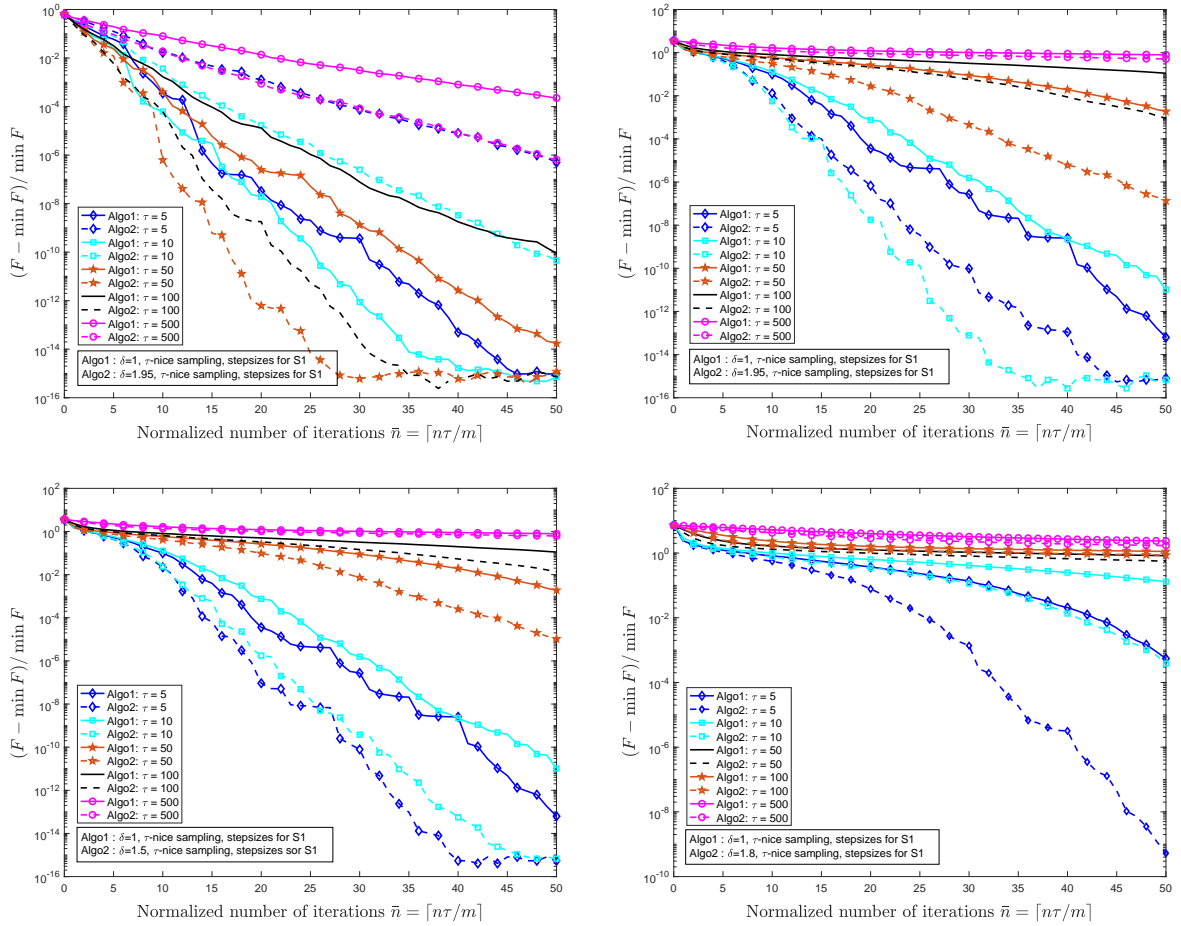


Figure 2: Comparison between  $\delta = 1$  and  $\delta > 1$ . Lasso problem with  $n = 10^3$  equations in  $m = 5 \cdot 10^3$  unknowns. Degree of partial separability  $\eta = 71$  (top left),  $\eta = 563$  (top right and bottom left), and  $\eta = 2594$  (bottom right). The choice  $\delta = 1$  is better than  $\delta > 1$  only for  $\eta = 71$  and  $\tau = 5, 10$ .

### The effect of $\delta > 1$

Here we study the effect of over-relaxing the stepsizes, meaning choosing  $\delta > 1$ . We compare Algorithm 1.1 with  $\delta = 1$  and several choices of  $\delta > 1$ . Figure 2 considers different scenarios for the degree of separability  $\eta$  of  $f$ . In those cases we see that choosing  $\delta > 1$  usually speeds up the convergence, depending on the parameter  $\eta$  of partial separability of  $f$  and the number  $\tau$  of parallel block updates. This fact seems not to occur when both  $\eta/m$  and  $\tau/m$  are very small.

### References

- [1] H.H. Bauschke, P.L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces. 2nd Ed.*, Springer, New York, 2017.
- [2] A. Beck, L. Tetruashvili, On the convergence of the block coordinate descent type methods, *SIAM J. Optim.*, vol. 23, n.4, pp. 2037–2060, 2013.

- [3] D.P. Bertsekas, Incremental proximal methods for large scale convex optimization, *Math. Program. Ser. B*, pp. 129–163, 2011.
- [4] P.L. Combettes, J-C. Pesquet, Stochastic quasi-Fejér block-coordinate fixed point iterations with random sweeping, *SIAM J. Optim.*, vol. 25, n.2, pp. 1121–1248, 2015.
- [5] P.L. Combettes, J-C. Pesquet, Stochastic quasi-Fejér block-coordinate fixed point iterations with random sweeping II: mean-square and linear convergence, *Math. Program. Ser. B*, pp. 1–19, 2018.
- [6] P.L. Combettes, V.R. Wajs, Signal recovery by proximal forward-backward splitting, *Multiscale Model. Simul.*, vol. 4, pp. 1168–1200, 2005.
- [7] D. Davis, Y. Yin, Convergence rate analysis of several splitting schemes. In *Splitting Methods in Communication, Imaging, Science, and Engineering* (R. Glowinski, S.J. Osher, and W. Yin, Eds.), pp. 115–163, Springer, Cham, 2016.
- [8] D. Drusvyatskiy, A.S. Lewis, Error bounds, quadratic growth, and linear convergence of proximal methods, *Math. Oper. Res.*, Vol. 43, pp. 919–948, 2018.
- [9] C. Dünnér, S. Forte, M. Takáč, M. Jaggi, Primal-dual rates and certificates, *International Conference on Machine Learning, PMLR*, 48, pp. 783–792, 2016.
- [10] R. Durrett, *Probability. Theory and Examples. 4th Ed.*, Cambridge University Press, New York, 2010.
- [11] Yu M. Ermol'ev, On the method of generalized stochastic gradients and quasi-Fejér sequences, *Cybernetics*, Vol. 5, pp. 208–220, 1969.
- [12] O. Fercoq, P. Bianchi, A Coordinate-descent primal-dual algorithm with large step size and possibly nonseparable functions. *SIAM J. Optim.*, vol. 29, pp. 100–134, 2019.
- [13] O. Fercoq, P. Richtàrik, Accelerated, parallel, and proximal coordinate descent. *SIAM J. Optim.*, vol. 25, pp. 1997–2023, 2015.
- [14] O. Fercoq, Z. Qu, Restarting the accelerated coordinate descent method with a rough strong convexity estimate. *Computational Optimization and Applications volume*, vol. 75, pp. 63–91, 2020.
- [15] G. Garrigos, L. Rosasco, S. Villa, Convergence of the forward-backward algorithm: beyond the worst-case with the help of geometry. [arXiv:1703.09477](https://arxiv.org/abs/1703.09477), 2017.
- [16] H. Karimi, J. Nutini, M. Schmidt, Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Machine Learning and Knowledge Discovery in Databases* (P. Frasconi, N. Landwehr, G. Manco and J. Vreeken Eds.), pp. 795–811, Springer International Publishing, Cham, 2016.
- [17] K. Knopp, *Infinite Sequences and Series.*, Dover Publications, Inc., New York, 1956.
- [18] K.K. Kiwiel, Convergence of approximate and incremental subgradient methods for convex optimization, *SIAM J. Optim.*, vol. 14, pp. 807–840, 2006.

- [19] D. Leventhal, A.S. Lewis, Randomized method for linear constraints: Convergence rates and conditioning. *Math. Oper. Res.*, vol. 35, pp. 641–654, 2010.
- [20] J. Lin, L. Rosasco, S Villa, D-X. Zhou, Modified Fejér sequences and applications, *Comput. Optim. Appl.* vol. 71, pp. 95–113, 2018.
- [21] Z. Lu, L. Xiao, On the complexity analysis of randomized block-coordinate descent methods. *Math. Program. Ser. A*, vol. 152, pp. 615–642, 2015.
- [22] Z-Q. Luo, P. Tseng, Error bounds and convergence analysis of feasible descent methods: a general approach. *Ann. Oper. Res.*, vol. 46, pp. 157–178, 1993.
- [23] I. Necoara, D. Clipici, Parallel random coordinate descent method for composite minimization: convergence analysis and error bounds. *SIAM J. Optim.*, vol. 26, pp. 197–226, 2016.
- [24] I. Necoara, Y. Nesterov, F. Glineur, Random block coordinate descent methods for linearly constrained optimization over networks. *J. Optim. Theory Appl.*, vol. 173, pp. 227–254, 2017.
- [25] I. Necoara, Y. Nesterov, F. Glineur, Linear convergence of first order methods for non-strongly convex optimization. *Math. Program.*, vol. 175, pp. 69–107, 2019.
- [26] A. Nemirovski, A. Juditsky, G. Lan, A. Shapiro, Robust stochastic approximation approach to stochastic programming, *SIAM J. Optim.*, vol. 19, pp. 1574–1609, 2009.
- [27] Y. Nesterov, *Introductory Lectures on Convex Optimization. A Basic Course*, Kluwer Academic Publishers, Boston, MA, 2004.
- [28] Y. Nesterov, Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM J. Optim.*, vol. 22, pp. 341–362, 2012.
- [29] Z. Qu, P. Richtàrik, T. Zhang, Quartz: randomized dual coordinate ascent with arbitrary sampling, *Advances in Neural Information Processing Systems*, 28, pp. 865–873, 2015.
- [30] Z. Qu, P. Richtàrik, Coordinate descent with arbitrary sampling I: algorithms and complexity, *Optim. Method Softw.* vol. 31, pp. 829–857, 2016.
- [31] Z. Qu, P. Richtàrik, Coordinate descent with arbitrary sampling II: expected separable overapproximation *Optim. Methods Softw.* vol. 31, pp. 858–884, 2016.
- [32] P. Richtàrik, M Takàč, Distributed coordinate descent method for learning with big data, *J Mach. Learn. Res.* vol. 17, pp. 1–25, 2016.
- [33] P. Richtàrik, M Takàč, Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function, *Math. Program. Ser. A*, vol. 144, pp. 1–38, 2014.
- [34] P. Richtàrik, M Takàč, Parallel coordinate descent methods for big data optimization, *Math. Program. Ser. A*, vol. 156, pp. 156–484, 2016.
- [35] P. Richtàrik, M Takàč, On optimal probabilities in stochastic coordinate descent methods, *Optim. Lett.* vol. 10, pp. 1233–1243, 2016.
- [36] S. Salzo, The variable metric forward-backward splitting algorithm under mild differentiability assumptions. *SIAM J. Optim.*, vol. 27, pp. 2153–2181, 2017.

- [37] T. Strohmer, R. Vershynin, A randomized Kaczmarz algorithm with exponential convergence, *J. Fourier Anal. Appl.*, vol. 15, n.2, pp. 262–278, 2009.
- [38] R. Tappenden, M Takáč, P. Richtárik, On the complexity of parallel coordinate descent, *Optim. Methods Softw.* vol. 33, pp. 373–395, 2018.
- [39] A. B. Taylor, J. M. Hendrickx, F. Glineur, Smooth strongly convex interpolation and exact worst-case performance of first-order methods. *Math. Program.*, vol. 161, pp. 307–345, 2017.
- [40] P-W. Wang, C-J. Lin, Iteration complexity of feasible descent methods for convex optimization, *J. Mach. Learn. Res.* vol. 15, pp. 1523–1548, 2014.
- [41] S. Wright, Coordinate descent algorithms, *Math. Program. Ser. B*, vol. 151, pp. 3–34, 2015.

## A Structured Lipschitz smoothness

In this section we discuss the Lipschitz smoothness properties of  $f$  under the hypotheses **H4** and **L1** and we prove Theorem 3.1. Most of the results presented in this section are basically given in [34]. However, here they are rephrased in our notation and extended to our more general assumptions.

**Proposition A.1.** *Let  $f: \mathbb{H} \rightarrow \mathbb{R}$  be a convex function satisfying assumptions **H4** and **L1**. Let  $I$  be a nonempty subset of  $[m]$  and let  $(q_i)_{i \in I} \in \mathbb{R}_{++}$  be such that  $\sum_{i \in I \cap I_k} q_i \leq 1$ , for every  $k \in [p]$ . Then for every  $x$  and  $y \in \mathbb{H}$  such that  $\text{spt}(x - y) \subset I$ , we have*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \sum_{i \in I} \frac{L_i}{q_i} \|y_i - x_i\|^2. \quad (\text{A.1})$$

*Proof.* Let  $v = y - x$  and, for every  $k \in [p]$ , set  $z^{(k)} = \sum_{i=1}^m U_{k,i} x_i$ . Then

$$f(y) = \sum_{k=1}^p g_k \left( \sum_{i=1}^m U_{k,i} (x_i + v_i) \right) = \sum_{k=1}^p g_k \left( z^{(k)} + \sum_{i \in I} U_{k,i} v_i \right) = \sum_{k=1}^p g_k \left( z^{(k)} + \sum_{i \in I \cap I_k} U_{k,i} v_i \right).$$

Now, for every  $k \in [p]$ , we have

$$z^{(k)} + \sum_{i \in I \cap I_k} U_{k,i} v_i = \left( 1 - \sum_{i \in I \cap I_k} q_i \right) z^{(k)} + \sum_{i \in I \cap I_k} q_i (z^{(k)} + q_i^{-1} U_{k,i} v_i).$$

Therefore, using the convexity of each  $g_k$  we have

$$f(y) = \sum_{k=1}^p g_k \left( z^{(k)} + \sum_{i \in I \cap I_k} U_{k,i} v_i \right) \leq \sum_{k=1}^p \left[ \left( 1 - \sum_{i \in I \cap I_k} q_i \right) g_k(z^{(k)}) + \sum_{i \in I \cap I_k} q_i g_k(z^{(k)} + q_i^{-1} U_{k,i} v_i) \right].$$

It follows from the definition of  $I_k$  that  $\sum_{i \in I \cap I_k} q_i g_k(z^{(k)} + q_i^{-1} U_{k,i} v_i) = \sum_{i \in I \setminus I_k} q_i g_k(z^{(k)})$ . Hence, switching the order of summation, and using the fact that  $f(x_1, \dots, x_{i-1}, \cdot, x_{i+1}, \dots, x_m)$  is Lipschitz

smooth with constant  $L_i$ , we have

$$\begin{aligned}
f(y) &\leq \sum_{k=1}^p \left[ \left(1 - \sum_{i \in I} q_i\right) \mathbf{g}_k(z^{(k)}) + \sum_{i \in I} q_i \mathbf{g}_k(z^{(k)} + q_i^{-1} \mathbf{U}_{k,i} \mathbf{v}_i) \right] \\
&= \left(1 - \sum_{i \in I} q_i\right) f(x) + \sum_{i \in I} q_i \sum_{k=1}^p \mathbf{g}_k \left( \sum_{j=1}^m \mathbf{U}_{k,j} (x_j + q_i^{-1} (\mathbf{J}_i \mathbf{v}_i)_j) \right) \\
&= \left(1 - \sum_{i \in I} q_i\right) f(x) + \sum_{i \in I} q_i f(x + q_i^{-1} \mathbf{J}_i \mathbf{v}_i) \\
&\leq \left(1 - \sum_{i \in I} q_i\right) f(x) + \sum_{i \in I} q_i \left[ f(x) + \langle \nabla_i f(x), q_i^{-1} \mathbf{v}_i \rangle + \frac{L_i}{2} \|q_i^{-1} \mathbf{v}_i\|^2 \right] \\
&= f(x) + \langle \nabla f(x), \mathbf{v} \rangle + \frac{1}{2} \sum_{i \in I} \frac{L_i}{q_i} \|\mathbf{v}_i\|^2. \quad \square
\end{aligned}$$

**Corollary A.2.** Let  $f: \mathbb{H} \rightarrow \mathbb{R}$  be a convex function satisfying **H4** and **L1**. Let  $\eta = \max_{1 \leq k \leq m} \text{card}(I_k)$ ,  $(\gamma_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$ , and  $(q_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$  be such that, for every  $k \in [p]$ ,  $\sum_{i \in I_k} q_i \leq 1$ . Let  $\Lambda = \bigoplus_{i=1}^m L_i \text{Id}_i$ ,  $\Gamma = \bigoplus_{i=1}^m \gamma_i \text{Id}_i$ , and  $\mathbf{Q} = \bigoplus_{i=1}^m q_i \text{Id}_i$ . Then, the function  $f$  is Lipschitz smooth

- (i) in the metric  $\|\cdot\|_W$  defined by  $W = \Lambda \mathbf{Q}^{-1}$  with constant 1;
- (ii) in the metric  $\|\cdot\|_{\Gamma^{-1}}$ , with constant<sup>3</sup>  $\max_{1 \leq k \leq p} \sum_{i \in I_k} \gamma_i L_i$ .
- (iii) in the (original) metric of  $\mathbb{H}$  with constant  $\max_{1 \leq k \leq p} \sum_{i \in I_k} L_i$ ;
- (iv) in the metric  $\|\cdot\|_\Lambda$ , with constant  $\eta$ .

*Proof.* (i): It follows from Proposition A.1 with  $I = [m]$  and noting that  $\langle \nabla f(x), y - x \rangle = \langle \nabla^W f(x), y - x \rangle_W$  and then invoking the characterization of the Lipschitz continuity of the gradient throughout the descent lemma (see [1, Theorem 18.15(iii)]).

(ii): It follows from (A.1) by choosing  $I = [m]$ ,  $q_i = \gamma_i L_i / (\max_{1 \leq k \leq p} \sum_{j \in I_k} \gamma_j L_j)$ , and noting that  $\langle \nabla f(x), y - x \rangle = \langle \nabla^{\Gamma^{-1}} f(x), y - x \rangle_{\Gamma^{-1}}$  and then invoking [1, Theorem 18.15(iii)].

(iii): It follows from (ii) with  $\gamma_i = 1$ .

(iv): It follows from (ii) with  $\gamma_i = 1/L_i$ . □

**Remark A.3.** If  $\eta = m$  ( $f$  is not partially separable), Corollary A.2-(iii)-(iv) establishes that

$$L_i \leq L \leq \sum_{i=1}^m L_i \quad \text{and} \quad L_{\|\cdot\|_\Lambda} \leq m.$$

We show that the above bounds are tight. Indeed, if we consider  $f(x) = (1/2) \|Ax - b\|_2^2$ , where  $A \in \mathbb{R}^{n \times m}$  and  $b \in \mathbb{R}^n$ , then we have  $L_i = \|a_i\|^2$  (where  $a_i$  is the  $i$ -th column of  $A$ ), so that  $\sum_{i=1}^m L_i = \|A\|_F^2$ . Instead, since  $\nabla f(x) = A^*(Ax - b)$ , the Lipschitz constant of  $\nabla f$  is  $\|A\|^2$ . It is well-known that if  $A$  is rank one, then  $\|A\|^2 = \|A\|_F^2$ , so in this case the Lipschitz constant of  $\nabla f$  is exactly  $\sum_{i=1}^m L_i$ . Moreover, if in addition the columns of  $A$  have the same norm, then  $L = \sum_{j=1}^m L_j = mL_i$  and hence  $L_{\|\cdot\|_\Lambda} = m$ . We finally note that if  $A$  is an orthonormal matrix, then  $\|A\|^2 = 1$  and hence  $1 = L_i = L = L_{\|\cdot\|_\Lambda} < \sum_{i=1}^m L_i = m$ .

<sup>3</sup>This constant is  $\leq \delta\eta / \min_{1 \leq i \leq m} \beta_{1,i}$  if the  $\gamma_i$ 's are set according to (1.5) with the  $\nu_i$ 's as in Theorem 3.1(i).



**Corollary A.4.** Let  $f: \mathbb{H} \rightarrow \mathbb{R}$  be a function satisfying **H4** and **L1**. Let  $\epsilon \in \{0, 1\}^m$  and  $\mathbf{x}, \mathbf{v} \in \mathbb{H}$ . Then,

$$f(\mathbf{x} + \epsilon \odot \mathbf{v}) \leq f(\mathbf{x}) + \sum_{i=1}^m \epsilon_i \langle \nabla_i f(\mathbf{x}), \mathbf{v}_i \rangle + \max_{1 \leq k \leq p} \left( \sum_{i \in I_k} \epsilon_i \right) \sum_{i=1}^m \epsilon_i \frac{L_i}{2} \|\mathbf{v}_i\|^2. \quad (\text{A.2})$$

*Proof.* It follows from Proposition **A.1** with  $\mathbf{y} = \mathbf{x} + \epsilon \odot \mathbf{v}$ ,  $I = \text{spt}(\epsilon)$ , and  $q_i = 1 / (\max_{1 \leq k \leq p} \text{card}(I \cap I_k)) = 1 / (\max_{1 \leq k \leq p} \sum_{i \in I_k} \epsilon_i)$ .  $\square$

**Remark A.5.** Most of the above results, appears in [34] for the special case that  $\mathbf{U}_{k,i} = \mathbf{J}_i$  for  $i \in I_k$  and  $\mathbf{U}_{k,i} = 0$  for  $i \notin I_k$ . In particular, see [34, Theorem 8].

**Proposition A.6.** Let  $f: \mathbb{H} \rightarrow \mathbb{R}$  be a function satisfying **H4** and suppose that, for every  $k \in [p]$ ,  $\mathbf{g}_k$  is  $L^{(k)}$ -Lipschitz smooth. Set for every  $i \in [m]$ ,  $\tilde{L}_i = \|\sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i}\|$ . Then the following holds.

- (i)  $f$  is Lipschitz smooth with constant  $\|\sum_{k=1}^p L^{(k)} \mathbf{U}_k^\top \mathbf{U}_k\|$  in the original metric of  $\mathbb{H}$ ;<sup>4</sup>
- (ii)  $f$  satisfies assumption **L1** with  $L_i = \tilde{L}_i$ ;
- (iii) Suppose that, for every  $k \in [p]$  and for every  $i, j \in [m]$ ,  $i \neq j$ , the range of the operators  $\mathbf{U}_{k,i}$  and  $\mathbf{U}_{k,j}$  are orthogonal. Then, for every  $\mathbf{x}, \mathbf{v} \in \mathbb{H}$ ,

$$f(\mathbf{x} + \mathbf{v}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \sum_{i=1}^m \tilde{L}_i \|\mathbf{v}_i\|^2. \quad (\text{A.3})$$

*Proof.* (i): For every  $k \in [p]$ , let  $\mathbf{U}_k: \mathbb{H} \rightarrow \mathbb{G}_k$ ,  $\mathbf{U}_k \mathbf{x} = \sum_{i=1}^m \mathbf{U}_{k,i} \mathbf{x}_i$ . Let  $\mathbf{x}, \mathbf{v} \in \mathbb{H}$ . We have

$$f(\mathbf{x} + \mathbf{v}) = \sum_{k=1}^p \mathbf{g}_k(\mathbf{U}_k \mathbf{x} + \mathbf{U}_k \mathbf{v}) \leq \sum_{k=1}^p \left( f_k(\mathbf{U}_k \mathbf{x}) + \langle \nabla f_k(\mathbf{U}_k \mathbf{x}), \mathbf{U}_k \mathbf{v} \rangle + \frac{L^{(k)}}{2} \|\mathbf{U}_k \mathbf{v}\|^2 \right). \quad (\text{A.4})$$

Therefore, we have

$$\begin{aligned} f(\mathbf{x} + \mathbf{v}) &\leq f(\mathbf{x}) + \sum_{k=1}^p \langle \mathbf{U}_k^\top \nabla f_k(\mathbf{U}_k \mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \sum_{k=1}^p L^{(k)} \langle \mathbf{U}_k^\top \mathbf{U}_k \mathbf{v}, \mathbf{v} \rangle \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \left\langle \sum_{k=1}^p L^{(k)} \mathbf{U}_k^\top \mathbf{U}_k \mathbf{v}, \mathbf{v} \right\rangle \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \left\| \sum_{k=1}^p L^{(k)} \mathbf{U}_k^\top \mathbf{U}_k \right\| \|\mathbf{v}\|^2. \end{aligned} \quad (\text{A.5})$$

(ii): It follows from (A.5) with  $\mathbf{v} = \mathbf{J}_i \mathbf{v}_i$  that

$$\begin{aligned} f(\mathbf{x} + \mathbf{J}_i \mathbf{v}_i) &\leq f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), \mathbf{v}_i \rangle + \frac{1}{2} \left\langle \sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \mathbf{v}_i, \mathbf{v}_i \right\rangle \\ &\leq f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), \mathbf{v}_i \rangle + \frac{1}{2} \left\| \sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \right\| \|\mathbf{v}_i\|^2 \end{aligned}$$

<sup>4</sup>In [4, Corollary 5.11] the worse constant  $\sum_{k=1}^p L^{(k)} \|\mathbf{U}_k \mathbf{U}_k^\top\|$  was considered.

hence  $\tilde{L}_i = \|\sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i}\|$  is a Lipschitz constant of  $\nabla_i f(\mathbf{x}_i, \dots, \mathbf{x}_{i-1}, \cdot, \mathbf{x}_{i+1}, \dots, \mathbf{x}_m)$ .

(iii): Since  $\langle \mathbf{U}_{k,i} \mathbf{v}_i, \mathbf{U}_{k,j} \mathbf{v}_j \rangle = 0$  if  $i \neq j$ , it follows from (A.4) that

$$\begin{aligned} f(\mathbf{x} + \mathbf{v}) &\leq f(\mathbf{x}) + \sum_{k=1}^p \langle \mathbf{U}_k^\top \nabla f_k(\mathbf{U}_k \mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \sum_{i=1}^m \sum_{k=1}^p L^{(k)} \langle \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \mathbf{v}_i, \mathbf{v}_i \rangle \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \sum_{i=1}^m \left\langle \sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \mathbf{v}_i, \mathbf{v}_i \right\rangle \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2} \sum_{i=1}^m \left\| \sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \right\| \|\mathbf{v}_i\|^2. \end{aligned} \quad \square$$

**Remark A.7.** If  $R(\mathbf{U}_{k,i})$  and  $R(\mathbf{U}_{k,j})$  are orthogonal to each other, then

$$\sum_{k=1}^p L^{(k)} \mathbf{U}_k^\top \mathbf{U}_k = \sum_{i=1}^m \mathbf{J}_i \underbrace{\left( \sum_{k=1}^p L^{(k)} \mathbf{U}_{k,i}^\top \mathbf{U}_{k,i} \right)}_{\mathbf{H}_i \rightarrow \mathbf{H}_i} \mathbf{J}_i^\top,$$

and hence  $\|\sum_{k=1}^p L^{(k)} \mathbf{U}_k^\top \mathbf{U}_k\| = \max_{1 \leq i \leq m} \tilde{L}_i$ .

## B Additional proofs

**Proof. of Theorem 3.1**

(ii): It follows from (A.2) that, point-wise it holds

$$f(\mathbf{x} + \varepsilon \odot \mathbf{v}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \varepsilon \odot \mathbf{v} \rangle + \max_{1 \leq k \leq p} \left( \sum_{i \in I_k} \varepsilon_i \right) \sum_{i=1}^m \varepsilon_i \frac{L_i}{2} \|\mathbf{v}_i\|^2. \quad (\text{B.1})$$

Moreover, since  $\beta_2 = \text{ess sup} \left( \max_{1 \leq k \leq p} \sum_{i \in I_k} \varepsilon_i \right)$ , we have that  $L_i \max_{1 \leq k \leq p} \sum_{i \in I_k} \varepsilon_i \leq L_i \beta_2 \leq \nu_i$  P-a.s. The statement follows.

(i) It follows by taking the expectation in (B.1) and noting that

$$\begin{aligned} \mathbb{E} \left[ \max_{1 \leq k \leq p} \left( \sum_{i \in I_k} \varepsilon_i \right) \sum_{i=1}^m \varepsilon_i \frac{L_i}{2} \|\mathbf{v}_i\|^2 \right] &= \sum_{i=1}^m \mathbb{E} \left[ \varepsilon_i \max_{1 \leq k \leq p} \left( \sum_{i \in I_k} \varepsilon_i \right) \right] \frac{L_i}{2} \|\mathbf{v}_i\|^2 \\ &= \sum_{i=1}^m \mathbb{E} \left[ \max_{1 \leq k \leq p} \left( \sum_{i \in I_k} \varepsilon_i \right) \mid \varepsilon_i = 1 \right] \mathbf{p}_i \frac{L_i}{2} \|\mathbf{v}_i\|^2, \end{aligned}$$

where we used the fact that for every discrete random variable  $\zeta$ ,  $\mathbb{E}[\varepsilon_i \zeta] = \mathbb{E}[\zeta \mid \varepsilon_i = 1] \mathbf{p}_i$ .

(iii) It follows from Proposition A.6(iii) that S3 holds with  $\nu_i = \tilde{L}_i$ .

(iv): Let  $i \in [m]$  and  $\mathbf{v}_i \in \mathbf{H}_i$  and set  $\mathbf{v} = \mathbf{J}_i \mathbf{v}_i = (0, \dots, 0, \mathbf{v}_i, 0, \dots, 0)$ . Then

$$\begin{aligned} \mathbb{E}[f(\mathbf{x} + \varepsilon \odot \mathbf{v})] &= \mathbb{E}[f(\mathbf{x} + \mathbf{J}_i(\varepsilon_i \mathbf{v}_i))] = \mathbf{p}_i f(\mathbf{x} + \mathbf{J}_i \mathbf{v}_i) + (1 - \mathbf{p}_i) f(\mathbf{x}) \\ \mathbb{E}[\langle \nabla f(\mathbf{x}), \varepsilon \odot \mathbf{v} \rangle] &= \mathbb{E}[\langle \nabla_i f(\mathbf{x}), \varepsilon_i \mathbf{v}_i \rangle] = \mathbf{p}_i \langle \nabla_i f(\mathbf{x}), \mathbf{v}_i \rangle \end{aligned}$$

Hence, it follows from S1 that  $f(\mathbf{x} + \mathbf{J}_i \mathbf{v}_i) \leq f(\mathbf{x}) + \langle \nabla_i f(\mathbf{x}), \mathbf{v}_i \rangle + (1/2) \nu_i \|\mathbf{v}_i\|^2$ . This shows that  $f$  is Lipschitz smooth w.r.t. the  $i$ -th block coordinate with Lipschitz constant  $\nu_i$ . The global Lipschitz smoothness of  $f$  follows from Corollary A.2.  $\square$

**Proof. of Remark 3.2(iv).**

We have  $f(\mathbf{x}) = \sum_{k=1}^p \mathbf{g}_k(\mathbf{U}_k \mathbf{x})$  and  $f(\mathbf{x} + \mathbf{v}) = \sum_{k=1}^p \mathbf{g}_k(\mathbf{U}_k \mathbf{x} + \sum_{i=1}^m \mathbf{U}_{k,i} \mathbf{v}_i)$ . Moreover,  $\nabla f(\mathbf{x}) = \sum_{k=1}^p \mathbf{U}_k^\top \nabla \mathbf{g}_k(\mathbf{U}_k \mathbf{x})$ . We let  $\mathbf{x} \in \mathbf{H}$  and define

$$\begin{aligned}\varphi(\mathbf{v}) &:= f(\mathbf{x} + \mathbf{v}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle \\ \psi_k(\mathbf{u}) &:= \mathbf{g}_k(\mathbf{U}_k \mathbf{x} + \mathbf{u}) - \mathbf{g}_k(\mathbf{U}_k \mathbf{x}) - \langle \nabla \mathbf{g}_k(\mathbf{U}_k \mathbf{x}), \mathbf{u} \rangle.\end{aligned}$$

We clearly have  $\varphi(\mathbf{v}) \geq 0$ ,  $\varphi(0) = 0$  and  $\psi_k(\mathbf{u}) \geq 0$ ,  $\psi_k(0) = 0$ . Moreover,

$$\varphi(\mathbf{v}) = \sum_{k=1}^p \mathbf{g}_k(\mathbf{U}_k \mathbf{x} + \mathbf{U}_k \mathbf{v}) - \sum_{k=1}^p \mathbf{g}_k(\mathbf{U}_k \mathbf{x}) - \sum_{k=1}^p \langle \nabla \mathbf{g}_k(\mathbf{U}_k \mathbf{x}), \mathbf{U}_k \mathbf{v} \rangle = \sum_{k=1}^p \psi_k(\mathbf{U}_k \mathbf{v}). \quad (\text{B.2})$$

Therefore,

$$\begin{aligned}\mathbb{E}[\varphi(\varepsilon \odot \mathbf{v})] &= \sum_{k=1}^p \mathbb{E} \left[ \psi_k \left( \sum_{i \in I_k} \mathbf{U}_{k,i} \varepsilon_i \mathbf{v}_i \right) \right] \\ &= \sum_{k=1}^p \sum_{t=1}^{\eta} \mathbb{E} \left[ \psi_k \left( \sum_{i \in I_k} \mathbf{U}_{k,i} \varepsilon_i \mathbf{v}_i \right) \mid \sum_{i \in I_k} \varepsilon_i = t \right] \mathbb{P} \left( \sum_{i \in I_k} \varepsilon_i = t \right).\end{aligned} \quad (\text{B.3})$$

Now, if  $\omega \in \Omega$  is such that  $\sum_{i \in I_k} \varepsilon_i(\omega) = t$  and we set  $I = \{i \in [m] \mid \varepsilon_i(\omega) = 1\}$ , then we have  $\text{card}(I \cap I_k) = t$  and

$$\psi_k \left( \sum_{i \in I_k} \mathbf{U}_{k,i} \varepsilon_i(\omega) \mathbf{v}_i \right) = \psi_k \left( \sum_{i \in I \cap I_k} \mathbf{U}_{k,i} \mathbf{v}_i \right) \leq \frac{1}{t} \sum_{i \in I \cap I_k} \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i) = \frac{1}{t} \sum_{i=1}^m \varepsilon_i(\omega) \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i).$$

Hence  $\psi_k \left( \sum_{i \in I_k} \mathbf{U}_{k,i} \varepsilon_i \mathbf{v}_i \right) \leq (1/t) \sum_{i=1}^m \varepsilon_i \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i)$  on the event  $\sum_{i \in I_k} \varepsilon_i = t$ . Then,

$$\begin{aligned}\mathbb{E} \left[ \psi_k \left( \sum_{i \in I_k} \mathbf{U}_{k,i} \varepsilon_i \mathbf{v}_i \right) \mid \sum_{i \in I_k} \varepsilon_i = t \right] &\leq \frac{1}{t} \sum_{i=1}^m \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i) \mathbb{E} \left[ \varepsilon_i \mid \sum_{i \in I_k} \varepsilon_i = t \right] \\ &= \frac{1}{t} \sum_{i=1}^m \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i) \mathbb{P} \left( \varepsilon_i = 1 \mid \sum_{i \in I_k} \varepsilon_i = t \right).\end{aligned} \quad (\text{B.4})$$

Plugging the above inequality in (B.3) we get

$$\begin{aligned}\mathbb{E}[\varphi(\varepsilon \odot \mathbf{v})] &\leq \sum_{k=1}^p \sum_{t=1}^{\eta} \frac{1}{t} \sum_{i=1}^m \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i) \mathbb{P} \left( \varepsilon_i = 1 \mid \sum_{i \in I_k} \varepsilon_i = t \right) \mathbb{P} \left( \sum_{i \in I_k} \varepsilon_i = t \right) \\ &= \sum_{i=1}^m \sum_{t=1}^{\eta} \frac{1}{t} \sum_{k=1}^p \psi_k(t \mathbf{U}_{k,i} \mathbf{v}_i) \mathbb{P} \left( \sum_{i \in I_k} \varepsilon_i = t \mid \varepsilon_i = 1 \right) \mathbf{p}_i \\ &\leq \sum_{i=1}^m \sum_{t=1}^{\eta} \frac{1}{t} \max_{\substack{1 \leq k \leq p \\ i \in I_k}} \mathbb{P} \left( \sum_{i \in I_k} \varepsilon_i = t \mid \varepsilon_i = 1 \right) \mathbf{p}_i \sum_{k=1}^p \psi_k(\mathbf{U}_{k,i} t \mathbf{v}_i) \\ &= \sum_{i=1}^m \sum_{t=1}^{\eta} \frac{1}{t} \max_{\substack{1 \leq k \leq p \\ i \in I_k}} \mathbb{P} \left( \sum_{i \in I_k} \varepsilon_i = t \mid \varepsilon_i = 1 \right) \mathbf{p}_i \varphi(\mathbf{J}_i t \mathbf{v}_i) \\ &\leq \sum_{i=1}^m \sum_{t=1}^{\eta} \frac{1}{t} \max_{\substack{1 \leq k \leq p \\ i \in I_k}} \mathbb{P} \left( \sum_{i \in I_k} \varepsilon_i = t \mid \varepsilon_i = 1 \right) \mathbf{p}_i \frac{L_i}{2} t^2 \|\mathbf{v}_i\|^2,\end{aligned} \quad (\text{B.5})$$

where in the last inequality we used **L1**. So, setting  $\beta_{1,i}$  as in (3.3), if  $\nu_i \geq \beta_{1,i}L_i$ , then **S1** holds. Note that in deriving (B.5), if for some  $i \in [m]$  there are no  $k \in [p]$  such that  $i \in I_k$ , the corresponding term  $\max_{k \in \emptyset} \mathbb{P}(\sum_{i \in I_k} \varepsilon_i = t \mid \varepsilon_i = 1)$  can be set to zero.  $\square$

**Proof. of formula (3.4).**

Since in the proof of (3.3) in Remark 3.2(iv), we only use the fact that  $i \notin I_k \Rightarrow U_{k,i} = 0$ , we can assume, without loss of generality, that, for every  $k \in [p]$ ,  $\text{card}(I_k) = \eta$ . Let  $i \in \bigcup_{k=1}^p I_k$ . Then, since the block sampling is doubly uniform, we have that  $\mathbb{P}(\sum_{i \in I_k} \varepsilon_i = t \mid \varepsilon_i = 1)$  does not depend on  $k$  such that  $i \in I_k$ . Therefore, (3.3) becomes  $\beta_{1,i} = \sum_{t=1}^{\eta} t \mathbb{P}(\sum_{j \in I_k} \varepsilon_j = t \mid \varepsilon_i = 1)$ , for some  $k \in [p]$  such that  $i \in I_k$ . Hence

$$\begin{aligned} \beta_{1,i} &= \mathbb{E} \left[ \sum_{j \in I_k} \varepsilon_j \mid \varepsilon_i = 1 \right] = \sum_{j \in I_k} \mathbb{E}[\varepsilon_j \mid \varepsilon_i = 1] = \sum_{j \in I_k} \mathbb{P}(\varepsilon_j = 1 \mid \varepsilon_i = 1) \\ &= \frac{1}{\mathfrak{p}} \sum_{i \in I_k} \mathbb{P}(\varepsilon_j = 1, \varepsilon_i = 1) = 1 + (\eta - 1) \frac{\tilde{\mathfrak{p}}}{\mathfrak{p}}, \end{aligned} \quad (\text{B.6})$$

where  $\mathfrak{p} = \mathbb{P}(\varepsilon_i = 1)$  and  $\tilde{\mathfrak{p}} = \mathbb{P}(\varepsilon_j = 1, \varepsilon_i = 1)$  (with  $i \neq j$ ). Now, since  $\mathbb{E}[\sum_{i=1}^m \varepsilon_i] = m\mathfrak{p}$  and  $\mathbb{E}[(\sum_{i=1}^m \varepsilon_i)^2] = \sum_{i=1}^m \mathbb{E}[\varepsilon_i^2] + \sum_{i \neq j} \mathbb{E}[\varepsilon_i \varepsilon_j] = m\mathfrak{p} + m(m-1)\tilde{\mathfrak{p}}$ , we have that  $\tilde{\mathfrak{p}} = (\mathbb{E}[(\sum_{i=1}^m \varepsilon_i)^2] - m\mathfrak{p}) / (m(m-1))$ , which plugged into (B.6) gives (3.4).

**Proof. of Lemma 4.3**

Let  $z \in \mathbb{H}$ . It follows from the definition of  $x^+$  that  $x - x^+ - \nabla\varphi(x) \in \partial\psi(x^+)$ . Therefore,  $\psi(z) \geq \psi(x^+) + \langle x - x^+ - \nabla\varphi(x), z - x^+ \rangle + (\mu_\psi/2)\|z - x^+\|^2$ , hence

$$\langle x - x^+, z - x^+ \rangle \leq \psi(z) - \psi(x^+) + \langle \nabla\varphi(x), z - x^+ \rangle - \frac{\mu_\psi}{2}\|z - x^+\|^2.$$

Now, we note that  $\|x^+ - z\|^2 = \|x^+ - x\|^2 + \|x - z\|^2 + 2\langle x^+ - x, x - z \rangle$ . Then,

$$\begin{aligned} \langle x - x^+, z - x \rangle + \langle x - x^+, x - x^+ \rangle &\leq \psi(z) - \psi(x^+) + \langle \nabla\varphi(x), z - x \rangle + \langle \nabla\varphi(x), x - x^+ \rangle \\ &\quad - \frac{\mu_\psi}{2}\|z - x\|^2 - \frac{\mu_\psi}{2}\|x - x^+\|^2 - \mu_\psi \langle x - x^+, z - x \rangle \end{aligned}$$

and hence

$$\begin{aligned} (1 + \mu_\psi)\langle x - x^+, z - x \rangle &\leq \psi(z) - \psi(x) + \langle \nabla\varphi(x), z - x \rangle - \frac{\mu_\psi}{2}\|z - x\|^2 + \psi(x) - \psi(x^+) \\ &\quad + \langle \nabla\varphi(x), x - x^+ \rangle - \left(1 + \frac{\mu_\psi}{2}\right)\|x - x^+\|^2. \end{aligned}$$

Since  $\langle \nabla\varphi(x), z - x \rangle \leq \varphi(z) - \varphi(x) - (\mu_\varphi/2)\|z - x\|^2$ , the statement follows.  $\square$

**Lemma B.1.** Let  $a, b, c \in \mathbb{R}_{++}$ . Then the largest constant  $\bar{\lambda} > 0$  satisfying the following inequality

$$\forall (s, t) \in \mathbb{R}_+^2, \text{ with } t \geq cs, \quad cs + t \geq \bar{\lambda}(as + bt) \quad (\text{B.7})$$

is

$$\bar{\lambda} = \min \left\{ \frac{1}{b}, \frac{2c}{a+bc} \right\} = \begin{cases} \frac{1}{b} & \text{if } b \geq \frac{a}{c} \\ \frac{2c}{a+bc} & \text{if } b \leq \frac{a}{c}. \end{cases} \quad (\text{B.8})$$

*Proof.* Property (B.7) is equivalent to

$$\forall (s, t) \in \mathbb{R}_+^2, \text{ with } t \geq cs \text{ and } cs + t > 0, \quad \frac{as + bt}{cs + t} \leq \frac{1}{\lambda}.$$

Therefore,

$$\begin{aligned} \frac{1}{\lambda} &= \sup \left\{ \frac{as + bt}{cs + t} \mid s, t \in \mathbb{R}_+, cs + t > 0, t \geq cs \right\} \\ &= \sup \{ as + bt \mid s, t \in \mathbb{R}_+, cs + t = 1, t \geq cs \}. \end{aligned} \quad (\text{B.9})$$

Now, since

$$(cs + t = 1 \text{ and } t \geq cs) \Leftrightarrow (cs = 1 - t \text{ and } t \geq 1 - t) \Leftrightarrow (cs = 1 - t \text{ and } t \geq 1/2),$$

it follows from (B.9) that

$$\frac{1}{\lambda} = \sup_{t \in [1/2, 1]} \frac{a}{c}(1 - t) + bt = \max \left\{ b, \frac{1}{2} \left( \frac{a}{c} + b \right) \right\} = \begin{cases} b & \text{if } b \geq \frac{a}{c} \\ \frac{1}{2} \left( \frac{a}{c} + b \right) & \text{if } b \leq \frac{a}{c}. \end{cases} \quad (\text{B.10})$$

Therefore, the statement follows.  $\square$

### **Proof. of Theorem 4.10**

We first note that, since,  $\|\cdot\|_{\Gamma^{-1}} \geq \mathfrak{p}_{\min} \|\cdot\|_{\mathbb{W}}$ , the conclusion of Proposition 4.6(iii) can be stated as follows:

$$\begin{aligned} &\mathbb{E} \left[ \mathfrak{p}_{\min} \frac{1 + \sigma_{\Gamma^{-1}}}{2} \|x^{n+1} - x\|_{\mathbb{W}}^2 + F(x^{n+1}) - F(x) \mid \mathfrak{E}_{n-1} \right] \\ &\leq \mathfrak{p}_{\min} \frac{1 + \sigma_{\Gamma^{-1}}}{2} \|x^n - x\|_{\mathbb{W}}^2 + F(x^n) - F(x) \\ &\quad - \mathfrak{p}_{\min} \left( \frac{\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}}}{2} \mathfrak{p}_{\min} \|x^n - x\|_{\mathbb{W}}^2 + F(x^n) - F(x) \right) \\ &\quad + \frac{(\delta - 1)_+}{2 + \sigma_{\Gamma^{-1}} - \delta} \mathbb{E}[F(x^n) - F(x^{n+1}) \mid \mathfrak{E}_{n-1}]. \end{aligned} \quad (\text{B.11})$$

Let  $x = x_*$  and set for brevity  $r_n^2 = (\mathfrak{p}_{\min}/2) \|x^n - x_*\|_{\mathbb{W}}^2$ , and  $F_n = F(x^n)$ . Then, (B.11) yields

$$\begin{aligned} \mathbb{E}[(1 + \sigma_{\Gamma^{-1}})r_{n+1}^2 + F_{n+1} - F_* \mid \mathfrak{E}_{n-1}] &\leq (1 + \sigma_{\Gamma^{-1}})r_n^2 + F_n - F_* \\ &\quad - \mathfrak{p}_{\min}((\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}})r_n^2 + F_n - F_*) \\ &\quad + \frac{(\delta - 1)_+}{2 + \sigma_{\Gamma^{-1}} - \delta} \mathbb{E}[F_n - F_{n+1} \mid \mathfrak{E}_{n-1}]. \end{aligned}$$

Let  $b = 1 + (\delta - 1)_+ / (2 + \sigma_{\Gamma^{-1}} - \delta)$ . Then the above inequality can be rewritten as

$$\begin{aligned} \mathbb{E}[(1 + \sigma_{\Gamma^{-1}})r_{n+1}^2 + b(F_{n+1} - F_*) \mid \mathfrak{E}_{n-1}] &\leq (1 + \sigma_{\Gamma^{-1}})r_n^2 + b(F_n - F_*) \\ &\quad - \mathfrak{p}_{\min}((\mu_{\Gamma^{-1}} + \sigma_{\Gamma^{-1}})r_n^2 + F_n - F_*). \end{aligned} \quad (\text{B.12})$$

Now, we derive from (2.1)-(2.2) that

$$F_n - F_* \geq \frac{\mu_{\Gamma-1} + \sigma_{\Gamma-1}}{2} \sum_{i=1}^m \frac{1}{\gamma_i} \|x_i^n - x_i\|^2 \geq \frac{\mu_{\Gamma-1} + \sigma_{\Gamma-1}}{2} \mathbf{p}_{\min} \|x^n - \mathbf{x}\|_{\mathbb{W}}^2 = (\mu_{\Gamma-1} + \sigma_{\Gamma-1}) r_n^2.$$

Therefore, it follows from Lemma B.1 (with  $c = \mu_{\Gamma-1} + \sigma_{\Gamma-1}$  and  $a = 1 + \sigma_{\Gamma-1}$ ) that

$$(\mu_{\Gamma-1} + \sigma_{\Gamma-1}) r_n^2 + F_n - F_* \geq \bar{\lambda} ((1 + \sigma_{\Gamma-1}) r_n^2 + b(F_n - F_*)), \quad (\text{B.13})$$

where

$$\bar{\lambda} = \begin{cases} \frac{1}{b} & \text{if } b \geq \frac{1 + \sigma_{\Gamma-1}}{\mu_{\Gamma-1} + \sigma_{\Gamma-1}} \\ \frac{2(\mu_{\Gamma-1} + \sigma_{\Gamma-1})}{1 + \sigma_{\Gamma-1} + b(\mu_{\Gamma-1} + \sigma_{\Gamma-1})} & \text{if } b \leq \frac{1 + \sigma_{\Gamma-1}}{\mu_{\Gamma-1} + \sigma_{\Gamma-1}}. \end{cases} \quad (\text{B.14})$$

Then, by (B.12) and (B.13), we have that

$$\mathbb{E}[(1 + \sigma_{\Gamma-1}) r_{n+1}^2 + b(F_{n+1} - F_*) \mid \mathfrak{E}_{n-1}] \leq (1 - \mathbf{p}_{\min} \bar{\lambda}) ((1 + \sigma_{\Gamma-1}) r_n^2 + b(F_n - F_*))$$

and hence, taking the expectation, and applying the resulting inequality recursively, we have,

$$b(\mathbb{E}[F_n] - F_*) \leq \mathbb{E}[(1 + \sigma_{\Gamma-1}) r_n^2 + b(F_n - F_*)] \leq (1 - \mathbf{p}_{\min} \bar{\lambda})^n ((1 + \sigma_{\Gamma-1}) r_0^2 + b(F_0 - F_*)). \quad (\text{B.15})$$

To conclude it is sufficient to note that, since

$$b = \max \left\{ 1, \frac{1 + \sigma_{\Gamma-1}}{2 - \delta + \sigma_{\Gamma-1}} \right\} \quad \text{and (in virtue of (2.4))} \quad \mu_{\Gamma-1} \leq \delta,$$

we have

$$\bar{\lambda} = \begin{cases} \frac{2 - \delta + \sigma_{\Gamma-1}}{1 + \sigma_{\Gamma-1}} & \text{if } \delta > 1 \text{ and } \mu_{\Gamma-1} \geq 2 - \delta \\ \frac{2(\mu_{\Gamma-1} + \sigma_{\Gamma-1})}{1 + \sigma_{\Gamma-1} + (\mu_{\Gamma-1} + \sigma_{\Gamma-1})(1 + \sigma_{\Gamma-1})/(2 - \delta + \sigma_{\Gamma-1})} & \text{if } \delta > 1 \text{ and } \mu_{\Gamma-1} \leq 2 - \delta \\ \frac{2(\mu_{\Gamma-1} + \sigma_{\Gamma-1})}{1 + \sigma_{\Gamma-1} + (\mu_{\Gamma-1} + \sigma_{\Gamma-1})} & \text{if } \delta \leq 1. \end{cases} \quad \square$$

## C Some results on duality theory

In this section, for the reader's convenience, we recap the results obtained in [9]. Let  $\varphi: \mathbb{H} \rightarrow \mathbb{R}$  and  $\psi: \mathbb{G} \rightarrow ]-\infty, +\infty]$  be two lower semicontinuous and convex functions defined on Hilbert spaces, and let  $A: \mathbb{H} \rightarrow \mathbb{G}$  be a bounded linear operator. In this section we suppose that  $\varphi$  is  $\mu$ -strongly convex. We consider the following optimization problems in duality (in the sense of Fenchel-Rockafellar)

$$\min_{x \in \mathbb{H}} \varphi(x) + \psi(Ax) := \mathcal{P}(x) \quad \text{and} \quad \min_{u \in \mathbb{G}} \psi^*(u) + \varphi^*(-A^\top u) := \mathcal{D}(u) \quad (\text{C.1})$$

We define the *duality gap function*  $G: \mathbb{H} \times \mathbb{G} \rightarrow ]-\infty, +\infty]$ ,  $G(x, u) = \mathcal{P}(x) + \mathcal{D}(u)$  and recall that

$$(\mathcal{P}(x) - \inf \mathcal{P}) + (\mathcal{D}(u) - \inf \mathcal{D}) \leq G(x, u).$$

So, the duality gap function bounds the primal and dual objectives. We have the following theorem

**Theorem C.1.** *Suppose that  $R(\mathbf{A}) \subset \text{dom } \partial\psi$ . Then the following holds:*

(i) *Suppose that  $\psi^*$  is  $\alpha$ -strongly convex. Let  $\mathbf{u} \in \text{dom } \psi^*$  and set  $\mathbf{x} = \nabla\varphi^*(-\mathbf{A}^\top \mathbf{u})$ . Then,*

$$G(\mathbf{x}, \mathbf{u}) \leq \left(1 + \frac{\|\mathbf{A}\|^2}{\alpha\mu}\right) (\mathcal{D}(\mathbf{u}) - \inf \mathcal{D}). \quad (\text{C.2})$$

(ii) *Suppose that  $\psi$  is  $\theta$ -Lipschitz continuous. Let  $\mathbf{u} \in \text{dom } \psi^*$  be such that  $\mathcal{D}(\mathbf{u}) - \inf \mathcal{D} < \|\mathbf{A}\|^2 L^2 / \mu$  and set  $\mathbf{x} = \nabla\varphi^*(-\mathbf{A}^\top \mathbf{u})$ . Then, we have*

$$G(\mathbf{x}, \mathbf{u}) \leq 2 \frac{\|\mathbf{A}\|\theta}{\mu^{1/2}} (\mathcal{D}(\mathbf{u}) - \inf \mathcal{D})^{1/2}. \quad (\text{C.3})$$

*Moreover, if  $u$  is a random variable taking values in  $\text{dom } \psi^*$  and such that  $\mathbb{E}[\mathcal{D}(u)] - \inf \mathcal{D} < \|\mathbf{A}\|^2 L^2 / \mu$  and we set  $x = \nabla\varphi^*(-\mathbf{A}^\top u)$ , then  $\mathbb{E}[G(x, u)] \leq 2\|\mathbf{A}\|\theta/\mu^{1/2} (\mathbb{E}[\mathcal{D}(u)] - \inf \mathcal{D})^{1/2}$ .*