

Block-proximal methods with spatially adapted acceleration

Tuomo Valkonen*

September 30, 2016

We propose a class of primal–dual block-coordinate descent methods based on blockwise proximal steps. The methods can be executed either stochastically, randomly executing updates of each block of variables, or the methods can be executed deterministically, featuring performance-improving blockwise adapted step lengths. Our work improves upon several previous stochastic primal–dual methods in that we do not require a full dual update in an accelerated method; both the primal and dual can be updated stochastically in blocks. We moreover provide convergence rates without the strong assumption of all functions being strongly convex or smooth: a mixed $O(1/N^2) + O(1/N)$ for an ergodic gap and individual strongly convex primal blocks. Within the context of deterministic methods, our blockwise adapted methods provide improvements on our earlier work on subspace accelerated methods. We test the proposed methods on various image processing problems.

1. Introduction

We want to efficiently solve optimisation problems of the form

$$\min_x G(x) + F(Kx), \quad (1.1)$$

arising from the variational regularisation of image processing and inverse problems. We assume $G : X \rightarrow \overline{\mathbb{R}}$ and $F : Y \rightarrow \overline{\mathbb{R}}$ to be convex, proper, and lower semicontinuous functionals on Hilbert spaces X and Y , respectively, and $K \in \mathcal{L}(X; Y)$ to be a bounded linear operator.

Several first-order optimisation methods have been developed for (1.1), typically with both G and F convex, and K linear, but recently also accepting a level of non-convexity and non-linearity [1–4]. Since at least one of G or F is typically non-smooth—popular regularisers for image processing are non-smooth, such as [5–7]—effective primal algorithms operating directly on the form (1.1), are typically a form of classical forward–backward splitting, occasionally going by the name of iterative soft-thresholding [8–10].

If G is separable, $G(x_1, \dots, x_m) = \sum_{j=1}^m G_j(x_j)$, problem (1.1) becomes

$$\min_x \sum_{j=1}^m G_j(x_j) + F(Kx). \quad (1.2)$$

*Department of Mathematical Sciences, University of Liverpool, United Kingdom.
tuomo.valkonen@liverpool.ac.uk

In big data optimisation, various forward–backward block-coordinate descent methods have been developed for this type of problems, or an alternative dual formulation. At each step of the optimisation method, they only update a subset of the blocks x_j , randomly in parallel, see e.g. the review [11] and the original articles [12–22]. Typically F is assumed smooth. Often, in these methods, each of the functions G_j is assumed strongly convex. Besides parallelism, one advantage of these methods is them being able to exploit the *local* blockwise factor of smoothness (Lipschitz gradient) of F and K . This can be better than the global factor, and helps to obtain improved convergence compared to standard methods.

Unfortunately, primal-only and dual-only stochastic approaches are rarely applicable to imaging problems that do not satisfy the separability and smoothness requirements simultaneously, at least not without additional Moreau–Yosida (aka. Huber, aka. Nesterov) regularisation. Generally, even without the splitting of the problem into blocks, primal-only or dual-only approaches, as discussed above, can be inefficient on more complicated problems, as the steps of the algorithms become very expensive optimisation problems themselves. This difficulty can often be circumvented through primal-dual approaches. If F is convex, and F^* denotes the conjugate of F , the problem (1.1) can be written in the min-max form

$$\min_x \max_y G(x) + \langle Kx, y \rangle - F^*(x), \quad (1.3)$$

If G is also convex, a popular algorithm for (1.3) is the Chambolle–Pock method [23, 24], also classified as the Primal-Dual Hybrid Gradient Method (Modified) or PDHGM in [25]. The method consists of alternating proximal steps on x and y , combined with an over-relaxation step that ensures convergence. The PDHGM is closely-related to the classical Alternating Direction Method of Multipliers (ADMM, [26]). Both the PDHGM and ADMM bear close relationship to Douglas–Rachford splitting [27, 28] and the split Bregman method [29, 30]. These relationships are discussed in detail in [25, 31].

While early work on block-coordinate descent methods concentrated on primal-only or dual-only algorithms, recently primal-dual algorithms based on the ADMM and the PDHGM have been proposed [32–38]. Besides [32, 33, 38] that have restrictive smoothness and strong convexity requirements, little is known about the convergence rates of these algorithms. Interestingly, a direct blockwise version of ADMM for (1.2) is however infeasible [39].

The convergence of the PDHGM can be accelerated from the basic rate $O(1/N)$ to $O(1/N^2)$ if G (equivalently F^*) is strongly convex [23]. However, saddle-point formulations (1.3) of important problems often lack strong convexity on the entire domain of G . These include any problem with higher-order TGV [6] or ICTV [7] regularisation, as well as such simple problems as deblurring with any regularisation term. Motivated by this, we recently showed in [3] that acceleration schemes can still produce convergence with a mixed rate, $O(1/N^2)$ with respect to initialisation, and $O(1/N)$ with respect to “residual variables”, if G is only *partially strongly convex*. This means with $x = (x_1, \dots, x_m)$ for some block k and a factor $\gamma_k > 0$ the condition

$$G(x') - G(x) \geq \langle z, x' - x \rangle + \frac{\gamma_k}{2} \|x'_k - x_k\|^2, \quad (1.4)$$

over all x' and subgradients $z \in \partial G(x)$. This property can be compared to partial smoothness on manifolds [40, 41], used to study the fast convergence of standard methods to submanifolds [42].

Under the condition (1.4), the iterates $\{x_k^i\}_{i=0}^\infty$ of the “partially accelerated” methods proposed in [3] will converge fast to an optimal solution \hat{x}_k , while nothing is known about the convergence of non-strongly-convex blocks. In Section 2 of this work, we

improve this analysis to be able to deal with *partial steps*, technically non-invertible step-length operators. We also allow the steps to be chosen stochastically. Some of the abstract proofs that are relatively minor generalisations of those in [3], we have left to Appendix A. From this abstract analysis, we derive in Sections 3 and 4 both stochastic and deterministic block-proximal primal-dual methods with the novelty of having *local or blockwise step lengths*. These can either be adapted dynamically to the actual sequence of block updates taken by the method, or taken deterministically with the goal of reducing communication in parallel implementations. As the most extreme case, which we consider in several of our image processing experiments in the final Section 5, our methods can have pixelwise-adapted step lengths.

The stochastic variants of our methods do not require the entire dual variable to be updated, it can also be randomised under light compatibility conditions on the primal and dual blocks. In the words of [38], our methods are “doubly-stochastic”. Our additional advances here are the convergence analysis—proving mixed $O(1/N^2) + O(1/N)$ rates of both the (squared) iterates and an ergodic gap—not being restricted to single-block updates, and not demanding strong convexity or smoothness from the entire problem, only individual blocks of interest.

2. A general method with non-invertible step operators

2.1. Background

To make the notation definite, we write $\mathcal{L}(X; Y)$ for the space of bounded linear operators between Hilbert spaces X and Y . The identity operator we denote by I . For $T, S \in \mathcal{L}(X; X)$, we use $T \geq S$ to mean that $T - S$ is positive semidefinite; in particular $T \geq 0$ means that T is positive semidefinite. Also for possibly non-self-adjoint T , we introduce the inner product and norm-like notations

$$\langle x, z \rangle_T := \langle Tx, z \rangle, \quad \text{and} \quad \|x\|_T := \sqrt{\langle x, x \rangle_T}, \quad (2.1)$$

the latter only defined for positive semi-definite T . We write $T \simeq T'$ if $\langle x, x \rangle_{T'-T} = 0$ for all x .

Denoting $\overline{\mathbb{R}} := [-\infty, \infty]$, we now let $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ be given convex, proper, lower semicontinuous functionals $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ on Hilbert spaces X and Y . We also let $K \in \mathcal{L}(X; Y)$ be a bounded linear operator. We then wish to solve the minimax problem

$$\min_{x \in X} \max_{y \in Y} G(x) + \langle Kx, y \rangle - F^*(y), \quad (\text{P})$$

assuming the existence of a solution $\hat{u} = (\hat{x}, \hat{y})$ satisfying the optimality conditions

$$-K^*\hat{y} \in \partial G(\hat{x}), \quad \text{and} \quad K\hat{x} \in \partial F^*(\hat{y}). \quad (\text{OC})$$

The primal-dual method of Chambolle and Pock [23] for (P) consists of iterating the system

$$x^{i+1} := (I + \tau_i \partial G)^{-1}(x^i - \tau_i K^* y^i), \quad (2.2a)$$

$$\bar{x}^{i+1} := \omega_i(x^{i+1} - x^i) + x^{i+1}, \quad (2.2b)$$

$$y^{i+1} := (I + \sigma_{i+1} \partial F^*)^{-1}(y^i + \sigma_{i+1} K \bar{x}^{i+1}). \quad (2.2c)$$

In the basic version of the algorithm, $\omega_i = 1$, $\tau_i \equiv \tau_0$, and $\sigma_i \equiv \sigma_0$, assuming that the step length parameters satisfy $\tau_0 \sigma_0 \|K\|^2 < 1$. The method has $O(1/N)$ rate for the ergodic duality gap. If G is strongly convex with factor γ , we may accelerate

$$\omega_i := 1/\sqrt{1 + 2\gamma\tau_i}, \quad \tau_{i+1} := \tau_i \omega_i, \quad \text{and} \quad \sigma_{i+1} := \sigma_i/\omega_i, \quad (2.3)$$

to achieve $O(1/N^2)$ convergence rates. To motivate our later choices, we observe that σ_0 is never needed if we equivalently parametrise the algorithm by $\delta = 1 - \|K\|^2 \tau_0 \sigma_0 > 0$.

In [3], we extended the algorithm (2.2) & (2.3) to partially strongly convex problems. For suitable step length operators $T_i \in \mathcal{L}(X; X)$ and $\Sigma_i \in \mathcal{L}(Y; Y)$, as well as an over-relaxation parameter $\tilde{\omega}_i > 0$, it consists of the iterations

$$x^{i+1} := (I + T_i \partial G)^{-1}(x^i - T_i K^* y^i), \quad (2.4a)$$

$$\tilde{x}^{i+1} := \tilde{\omega}_i(x^{i+1} - x^i) + x^{i+1}, \quad (2.4b)$$

$$y^{i+1} := (I + \Sigma_{i+1} \partial F^*)^{-1}(y^i + \Sigma_{i+1} K \tilde{x}^{i+1}), \quad (2.4c)$$

The diagonally preconditioned algorithm of [43] also fits into this form. For specific choices of T_i , Σ_i , and $\tilde{\omega}_i$, and under additional conditions on G and F^* , we were able to obtain mixed $O(1/N^2) + O(1/N)$ convergence rates of an ergodic duality gap, as well as the squared distance of the primal iterates within a subspace on which G is strongly convex.

Let us momentarily assume that the operators T_i and Σ_{i+1} are invertible. Following [3, 44], we may with the general variable splitting notation

$$u = (x, y),$$

write the system (2.4) as

$$0 \in H(u^{i+1}) + M_i(u^{i+1} - u^i), \quad (\text{PP}_0)$$

for the monotone operator

$$H(u) := \begin{pmatrix} \partial G(x) + K^* y \\ \partial F^*(y) - Kx \end{pmatrix}, \quad (2.5)$$

and the *preconditioning* or *step-length operator*

$$M_i := \begin{pmatrix} T_i^{-1} & -K^* \\ -\tilde{\omega}_i K & \Sigma_{i+1}^{-1} \end{pmatrix}. \quad (2.6)$$

With these operators, the optimality conditions (OC) can also be encoded as $0 \in H(\hat{u})$.

Remark 2.1 (A word about the indexing). The reader may have noticed that the steps of (2.2) and (2.4) involve T_i and Σ_{i+1} , with distinct indices. This is mainly to maintain in (2.3) the identity $\tau_i \sigma_i = \tau_{i+1} \sigma_{i+1}$ from [23]. On the other hand the proximal point formulation necessitates x -first ordering of the steps (2.2) in contrast to the y -first order in [23]. Thus our y^{i+1} is their y^i .

2.2. Non-invertible step length operators

What if T_i and Σ_{i+1} are non-invertible? The algorithm (2.4) of course works, but how about the proximal point version (PP₀)? We want to use the form (PP₀), because it greatly eases the analysis of the method.

Defining

$$W_{i+1} := \begin{pmatrix} T_i & 0 \\ 0 & \Sigma_{i+1} \end{pmatrix}, \quad \text{and (for now)} \quad L_{i+1} = \begin{pmatrix} I & -T_i K^* \\ -\tilde{\omega}_i \Sigma_{i+1} K & I \end{pmatrix}, \quad (2.7)$$

the method (2.4) can also be written

$$W_{i+1} H(u^{i+1}) + L_{i+1}(u^{i+1} - u^i) \ni 0. \quad (\text{PP})$$

To study the convergence of (PP), we apply the concept of *testing* that we introduced in [3]. The idea is to analyse the inclusion obtained by multiplying (PP) by the testing operator

$$Z_{i+1} := \begin{pmatrix} \Phi_i & 0 \\ 0 & \Psi_{i+1} \end{pmatrix} \quad (2.8)$$

for some primal test $\Phi_i \in \mathcal{L}(X; X)$ and dual test $\Psi_{i+1} \in \mathcal{L}(Y; Y)$. To employ the general estimates of Appendix A, we need $Z_{i+1}L_{i+1}$ to be self-adjoint and positive semi-definite. Therefore, we allow for general $L_{i+1} \in \mathcal{L}(X \times Y; X \times Y)$ instead of fixing the one in (2.7), and assume that

$$Z_{i+1}L_{i+1} = \begin{pmatrix} \Phi_i & -\Lambda_i^* \\ -\Lambda_i & \Psi_{i+1} \end{pmatrix} \geq 0 \text{ and is self-adjoint,} \quad (C0)$$

for some $\Lambda_i \in \mathcal{L}(X; Y)$, ($i \in \mathbb{N}$). If Φ_i and Ψ_{i+1} are invertible, we can solve L_{i+1} from (C0).

For $\Gamma \in \mathcal{L}(X; X)$, we define

$$\Xi_{i+1}(\Gamma) := \begin{pmatrix} 2T_i\Gamma & 2T_iK^* \\ -2\Sigma_{i+1}K & 0 \end{pmatrix}, \quad (2.9)$$

and

$$\Delta_{i+1}(\Gamma) := Z_{i+1}L_{i+1} - Z_i(L_i + \Xi_i(\Gamma)). \quad (2.10)$$

In practise, Γ will in be given by the partial strong monotonicity or convexity of G that we will soon discuss in Sections 2.4 and 2.5. The remainder of this manuscript by and large consists of estimating Δ_{i+1} . It measures with respect to the tests Z_i and Z_{i+1} , the difference between the local metrics induced by the operators L_{i+1} and $L_i + \Xi_i(\Gamma)$. The first is chosen by practical considerations of algorithm realisability, while the second one would be the theoretically desired one from the analysis in Appendix A. We begin this work by expanding Δ_{i+1} .

Lemma 2.1. *Suppose (C0) holds. Then*

$$\Delta_{i+2}(\Gamma) \simeq \begin{pmatrix} \Phi_{i+1} - \Phi_i(I + 2T_i\Gamma) & A_{i+2}^* \\ A_{i+2} & \Psi_{i+2} - \Psi_{i+1} \end{pmatrix} \quad (2.11)$$

for

$$A_{i+2} := (\Psi_{i+1}\Sigma_{i+1}K - \Lambda_{i+1}) + (\Lambda_i - KT_i^*\Phi_i^*). \quad (2.12)$$

Proof. Using (C0), we have

$$Z_{i+1}(L_{i+1} + \Xi_{i+1}(\Gamma)) = \begin{pmatrix} \Phi_i(I + 2T_i\Gamma) & 2\Phi_iT_iK^* - \Lambda_i^* \\ -2\Psi_{i+1}\Sigma_{i+1}K - \Lambda_i & \Psi_{i+1} \end{pmatrix}.$$

In particular

$$\begin{aligned} \Delta_{i+2}(\Gamma) &= \begin{pmatrix} \Phi_{i+1} - \Phi_i(I + 2T_i\Gamma) & \Lambda_i^* - \Lambda_{i+1}^* - 2\Phi_iT_iK^* \\ 2\Psi_{i+1}\Sigma_{i+1}K + \Lambda_i - \Lambda_{i+1} & \Psi_{i+2} - \Psi_{i+1} \end{pmatrix} \\ &\simeq \begin{pmatrix} \Phi_{i+1} - \Phi_i(I + 2T_i\Gamma) & A_{i+2}^* \\ A_{i+2} & \Psi_{i+2} - \Psi_{i+1} \end{pmatrix}. \quad \square \end{aligned}$$

Example 2.1. In the standard algorithm (2.2) & (2.3), we have $T_i = \tau_i I$, $\Phi_i = \tau_i^{-2} I$, and $\Gamma = \gamma I$, as well as $\Sigma_i = \sigma_i I$ and $\Psi_{i+1} = 1/(1 - \delta) = 1/(\tau_0 \sigma_0 \|K\|^2)$. It follows that $\Delta_{i+1} = 0$.

Remark 2.2. Later on, we will need to force various expectations of A_{i+2} to equal zero. It may already be instructive for the reader to consider how Λ_i will be constrained when one sets $A_{i+2} = 0$, as will happen in deterministic variants of our algorithm.

2.3. Stochastic variants

Just before commencing with the i :th iteration of (PP), let us choose T_i and Σ_{i+1} randomly, only based on the information we have gathered beforehand. We denote this information by \mathcal{O}_{i-1} , including the specific random realisations of T_{i-1} and Σ_i . Technically \mathcal{O}_{i-1} is a σ -algebra, and satisfies $\mathcal{O}_{i-1} \subset \mathcal{O}_i$. For details on the measure-theoretic approach to probability, we refer the reader to [45], here we merely recall some basic concepts.

Definition 2.1. We denote by $(\Omega, \mathcal{O}, \mathbb{P})$ the *probability space* consisting of the set Ω of possible realisation of a random experiment, by \mathcal{O} a σ -algebra on Ω , and by \mathbb{P} a probability measure on (Ω, \mathcal{O}) . We denote the expectation corresponding to \mathbb{P} by \mathbb{E} , the conditional probability with respect to a sub- σ -algebra $\mathcal{O}' \subset \mathcal{O}$ by $\mathbb{P}[\cdot | \mathcal{O}']$, and the conditional expectation by $\mathbb{E}[\cdot | \mathcal{O}']$.

We also use the next non-standard notation.

Definition 2.2. If \mathcal{O} is a σ -algebra on the space Ω , we denote by $\mathcal{R}(\mathcal{O}; V)$ the space of V -valued random variables A , such that $A : \Omega \rightarrow V$ is \mathcal{O} -measurable.

We frequently abuse notation, and use A as if it were a variable in V . For example, $T_i \in \mathcal{R}(\mathcal{O}_i; \mathcal{L}(X; X))$, but we generally think of T_i directly as an operator in $\mathcal{L}(X; X)$. Indeed, in our work, random variables and probability spaces only become apparent in the expectations \mathbb{E} , and probabilities \mathbb{P} , otherwise everything is happening in the underlying space V . The spaces $\mathcal{R}(\mathcal{O}; V)$ mainly serve to clarify the measurability with respect to different σ -algebras \mathcal{O}_i , that is, the algorithm iteration on which the realisation of a random variable is known.

\mathcal{O}_i thus includes all our knowledge prior to iteration i , with the interpretation that the random realisations of T_i and Σ_{i+1} are also known just before iteration i begins. Formally

$$T_i \in \mathcal{R}(\mathcal{O}_i; \mathcal{L}(X; X)), \quad \text{and} \quad \Sigma_{i+1} \in \mathcal{R}(\mathcal{O}_i; \mathcal{L}(Y; Y)).$$

Assuming that also $\Phi_i \in \mathcal{R}(\mathcal{O}_i; \mathcal{L}(X; X))$ and $\Psi_{i+1} \in \mathcal{R}(\mathcal{O}_i; \mathcal{L}(Y; Y))$, we deduce from (PP) that $x^{i+1} \in \mathcal{R}(\mathcal{O}_i; X)$ and $y^{i+1} \in \mathcal{R}(\mathcal{O}_i; Y)$. We say that any variable, e.g., $x^{i+1} \in \mathcal{R}(\mathcal{O}_i; X)$, that can be computed from the information in \mathcal{O}_i to be *deterministic with respect to \mathcal{O}_i* . Then in particular $\mathbb{E}[x^{i+1} | \mathcal{O}_k] = x^{i+1}$ for all $k \geq i$.

2.4. Basic estimates on the abstract proximal point iteration

To derive convergence rate estimates, we start by formulating abstract forms of partial strong monotonicity. As a first step, we take subsets of operators

$$\mathcal{T} \subset \mathcal{L}(X; X), \quad \text{and} \quad \mathcal{S} \subset \mathcal{L}(Y; Y).$$

We then suppose that ∂G is *partially strongly \mathcal{T} -monotone*, which we take to mean

$$\langle \partial G(x') - \partial G(x), \tilde{T}^*(x' - x) \rangle \geq \|x' - x\|_{\tilde{T}\Gamma}^2, \quad (x, x' \in X; \tilde{T} \in \mathcal{T}) \quad (\text{G-PM})$$

for some linear operator $0 \leq \Gamma \in \mathcal{L}(X; X)$. The operator $\tilde{T} \in \mathcal{T}$ acts as a testing operator. Similarly, we assume that ∂F^* is *\mathcal{S} -monotone* in the sense

$$\langle \partial F^*(y') - \partial F^*(y), \tilde{\Sigma}^*(y' - y) \rangle \geq 0, \quad (y, y' \in Y; \tilde{\Sigma} \in \mathcal{S}). \quad (\text{F}^*\text{-PM})$$

Example 2.2. $G(x) = \frac{1}{2}\|f - Ax\|^2$ satisfies (G-PM) with $\Gamma = A^*A$ and $\mathcal{T} = \mathcal{L}(X; X)$. Indeed, we calculate $\langle \nabla G(x') - \nabla G(x), \tilde{T}^*(x' - x) \rangle = \langle A^*A(x' - x), \tilde{T}^*(x' - x) \rangle = \|x' - x\|_{\tilde{T}\Gamma}^2$.

Example 2.3. If $F^*(y) = \sum_{\ell=1}^n F_\ell^*(y_\ell)$ for $y = (y_1, \dots, y_n)$ with each F_ℓ^* convex, then F^* satisfies **(F*-PM)** with $\mathcal{S} = \{\sum_{\ell=1}^n \beta_\ell Q_\ell \mid \beta_\ell \geq 0\}$ for $Q_\ell y := y_\ell$. Indeed, the condition **(F*-PM)** simply splits into separate monotonicity conditions for each ∂F_ℓ^* , ($\ell = 1, \dots, n$).

The fundamental estimate that forms the basis of all our convergence results, is the following. We note that the theorem does not yet prove convergence, but to do so we have to estimate the “penalty sum” involving the operators Δ_{i+2} , as well as the operator $Z_{N+1}L_{N+1}$. This is the content of the much of the rest of this paper. Moreover, to derive more meaningful “on average” results in the stochastic setting, we will take the expectation of (2.13).

Theorem 2.1. *Let us be given $K \in \mathcal{L}(X; Y)$, and convex, proper, lower semicontinuous functionals $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ on Hilbert spaces X and Y , satisfying **(G-PM)** and **(F*-PM)** for some $0 \leq \Gamma \in \mathcal{L}(X; X)$. Suppose the (random) operators $T_i, \Phi_i \in \mathcal{R}(O_i; \mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(O_i; \mathcal{L}(Y; Y))$ satisfy $\Phi_i T_i \in \mathcal{R}(O_i; \mathcal{T})$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{R}(O_i; \mathcal{S})$ for each $i \in \mathbb{N}$. If, moreover, **(C0)** holds, then the iterates $u^i = (x^i, y^i)$ of the proximal point iteration **(PP)** satisfy for all $N \geq 1$ the estimate*

$$\|u^N - \hat{u}\|_{Z_{N+1}L_{N+1}}^2 + \sum_{i=0}^{N-1} \|u^{i+1} - u^i\|_{Z_{i+1}L_{i+1}}^2 \leq \|u^0 - \hat{u}\|_{Z_0L_0}^2 + \sum_{i=0}^{N-1} \|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\Gamma)}^2. \quad (2.13)$$

As the proof is a relatively straightforward improvement of [3, Theorem 2.1] to non-invertible T_i and Σ_i , we relegate it to Appendix A.

2.5. Estimates on an ergodic duality gap

We may also prove the convergence of an ergodic duality gap. For this, the abstract monotonicity assumptions **(G-PM)** and **(F*-PM)** are not enough, and we need analogous convexity formulations. We find it most straightforward to formulate these conditions directly in the stochastic setting. Namely we assume for all $N \geq 1$ that whenever $\tilde{T}_i \in \mathcal{R}(O_i; \mathcal{T})$ and $x^{i+1} \in \mathcal{R}(O_i; X)$ for each $i = 0, \dots, N-1$, and that $\sum_{i=0}^{N-1} \mathbb{E}[\tilde{T}_i] = I$, then for some $0 \leq \Gamma \in \mathcal{L}(X; X)$ holds

$$G\left(\sum_{i=0}^{N-1} \mathbb{E}[\tilde{T}_i^* x^{i+1}]\right) - G(\hat{x}) \geq \sum_{i=0}^{N-1} \mathbb{E}\left[\langle \partial G(x^{i+1}), \tilde{T}_i^*(x^{i+1} - \hat{x}) \rangle + \frac{1}{2} \|x^{i+1} - \hat{x}\|_{\tilde{T}_i \Gamma}^2\right]. \quad (\text{G-EC})$$

Analogously, we assume whenever $\tilde{\Sigma}_{i+1} \in \mathcal{R}(O_i; \mathcal{S})$ and $y^{i+1} \in \mathcal{R}(O_i; Y)$ for each $i = 0, \dots, N-1$ with $\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\Sigma}_{i+1}] = I$ that

$$F^*\left(\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\Sigma}_{i+1}^* y^{i+1}]\right) - F^*(\hat{y}) \geq \sum_{i=0}^{N-1} \mathbb{E}\left[\langle \partial F^*(y^{i+1}), \tilde{\Sigma}_{i+1}^*(y^{i+1} - \hat{y}) \rangle\right]. \quad (\text{F*-EC})$$

If everything is deterministic, **(G-EC)** and **(F*-EC)** with $N = 1$ imply **(G-PM)** and **(F*-PM)**.

Example 2.4. If $\tilde{\Sigma}_{i+1} = \tilde{\sigma}_{i+1} I$ for some (random) scalar $\tilde{\sigma}_{i+1}$, then it is easy to verify **(F*-EC)** using Jensen’s inequality. We will generalise this example in Section 3.2.

Further, we assume for some $\bar{\eta}_i > 0$ that

$$\mathbb{E}[T_i^* \Phi_i^*] = \bar{\eta}_i I, \quad \text{and} \quad \mathbb{E}[\Psi_{i+1} \Sigma_{i+1}] = \bar{\eta}_i I, \quad (i \geq 1), \quad (\text{CG})$$

and for

$$\zeta_N := \sum_{i=0}^{N-1} \bar{\eta}_i \quad (2.14)$$

define the ergodic sequences

$$\tilde{x}_N := \zeta_N^{-1} \mathbb{E} \left[\sum_{i=0}^{N-1} T_i^* \Phi_i^* x^{i+1} \right], \quad \text{and} \quad \tilde{y}_N := \zeta_N^{-1} \mathbb{E} \left[\sum_{i=0}^{N-1} \Sigma_{i+1}^* \Psi_{i+1}^* y^{i+1} \right]. \quad (2.15)$$

These sequences will eventually be generated through the application of (G-EC) and (F*-EC) with $\tilde{T}_i := \Phi_i T_i$ and $\tilde{\Sigma}_i := \Psi_i \Sigma_i$. Finally, we introduce the duality gap

$$\mathcal{G}(x, y) := (G(x) + \langle \hat{y}, Kx \rangle - F(\hat{y})) - (G(\hat{x}) + \langle y, K\hat{x} \rangle - F^*(y)). \quad (2.16)$$

Then we have:

Theorem 2.2. *Let us be given $K \in \mathcal{L}(X; Y)$, and convex, proper, lower semicontinuous functionals $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ on Hilbert spaces X and Y , satisfying (G-PM), (F*-PM), (G-EC) and (F*-EC) for some $0 \leq \Gamma \in \mathcal{L}(X; X)$. Suppose the (random) operators $T_i, \Phi_i \in \mathcal{R}(O_i; \mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(O_i; \mathcal{L}(Y; Y))$ satisfy $\Phi_i T_i \in \mathcal{R}(O_i; \mathcal{T})$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{R}(O_i; \mathcal{S})$ for each $i \in \mathbb{N}$. If, moreover, (C0) and (CG) hold, then the iterates $u^i = (x^i, y^i)$ of the proximal point iteration (PP) satisfy for all $N \geq 1$ the estimate*

$$\begin{aligned} \mathbb{E} \left[\|x^N - \hat{x}\|_{Z_{N+1}L_{N+1}}^2 + \sum_{i=0}^{N-1} \|u^{i+1} - u^i\|_{Z_{i+1}L_{i+1}}^2 \right] + \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N) \\ \leq \|u^0 - \hat{u}\|_{Z_0L_0}^2 + \sum_{i=0}^{N-1} \mathbb{E}[\|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\Gamma/2)}^2]. \end{aligned} \quad (2.17)$$

Again, the proof is in Appendix A.

Remark 2.3. The difference between $\Delta_{i+2}(\Gamma/2)$ and $\Delta_{i+2}(\Gamma)$ in (2.17) and (2.13) corresponds to the fact [23, 46] that in (2.2) we need slower $\omega_i = 1/\sqrt{1 + \gamma\tau_i}$ compared to (2.3), to get $O(1/N^2)$ convergence of the gap. The scheme (2.3) is only known to yield convergence of the iterates.

2.6. Estimates on another ergodic duality gap

The condition (CG) does not hold for the basic non-stochastic accelerated method (2.2) & (2.3), which would satisfy $\tau_i \phi_i = \psi \sigma_i$ for a suitable constant ψ and $\phi_i = \tau_i^{-2}$; see Example 2.1. We therefore assume for some $\bar{\eta}_i \in \mathbb{R}$ that

$$\mathbb{E}[T_i^* \Phi_i^*] = \bar{\eta}_i I, \quad \text{and} \quad \mathbb{E}[\Psi_i \Sigma_i] = \bar{\eta}_i I, \quad (i \geq 1), \quad (\text{CG}_*)$$

and for

$$\zeta_{*,N} := \sum_{i=1}^{N-1} \bar{\eta}_i \quad (2.18)$$

define

$$\tilde{x}_{*,N} := \zeta_{*,N}^{-1} \mathbb{E} \left[\sum_{i=1}^{N-1} T_i^* \Phi_i^* x^{i+1} \right], \quad \text{and} \quad \tilde{y}_{*,N} := \zeta_{*,N}^{-1} \mathbb{E} \left[\sum_{i=1}^{N-1} \Sigma_i^* \Psi_i^* y^i \right]. \quad (2.19)$$

These variables give our new duality gap. Our original ergodic variables will however turn out to be the only possibility for doubly-stochastic methods. The convergence result is:

Theorem 2.3. *In Theorem 2.2, let us assume Eq. (CG_{*}) in place of Eq. (CG). Then (2.17) holds with $\zeta_{*,N} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N})$ in place of $\zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N)$.*

Once again, the proof is in Appendix A.

Remark 2.4. Technically, in Theorems 2.1 to 2.3, we do not need T_i , Φ_i , Σ_{i+1} , and Ψ_{i+1} deterministic with respect to O_i . It is sufficient that $T_i, \Phi_i \in \mathcal{R}(O; \mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(O; \mathcal{L}(Y; Y))$ satisfy $\Phi_i T_i \in \mathcal{R}(O; \mathcal{T})$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{R}(O; \mathcal{S})$ for each $i \in \mathbb{N}$. For consistency, we have however made the stronger assumptions, which we will start needing from now on.

2.7. Simplifications and summary so far

We assume the existence of some fixed $\delta > 0$, and bounds $b_{i+2}(\tilde{\Gamma})$ dependent on our choice of $\tilde{\Gamma} \in \mathcal{L}(X; X)$ such that

$$\mathbb{E}[Z_{i+2}L_{i+2}|O_i] \geq \begin{pmatrix} \delta\Phi_{i+1} & 0 \\ 0 & 0 \end{pmatrix}, \quad \text{and} \quad (\text{C1})$$

$$\mathbb{E}[\|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\tilde{\Gamma})}^2] \leq \mathbb{E}[\|u^{i+1} - u^i\|_{Z_{i+1}L_{i+1}}^2] + b_{i+2}(\tilde{\Gamma}), \quad (i \in \mathbb{N}). \quad (\text{C2})$$

Then we obtain the following corollary.

Corollary 2.1. *Let us be given $K \in \mathcal{L}(X; Y)$, and convex, proper, lower semicontinuous functionals $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ on Hilbert spaces X and Y , satisfying (G-PM) and (F*-PM) for some $0 \leq \Gamma \in \mathcal{L}(X; X)$. Suppose (C0), (C1), and (C2) hold for some random operators $T_i, \Phi_i \in \mathcal{R}(O_i; \mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(O_i; \mathcal{L}(Y; Y))$ with $\Phi_i T_i \in \mathcal{R}(O_i; \mathcal{T})$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{R}(O_i; \mathcal{S})$ for each $i \in \mathbb{N}$. Assuming one of the following cases to hold, let*

$$\tilde{g}_N := \begin{cases} 0, & \tilde{\Gamma} = \Gamma, \\ \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N), & \tilde{\Gamma} = \Gamma/2; (\text{G-EC}), (\text{F}^*\text{-EC}) \text{ and } (\text{CG}) \text{ hold} \\ \zeta_{*,N} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}), & \tilde{\Gamma} = \Gamma/2; (\text{G-EC}), (\text{F}^*\text{-EC}) \text{ and } (\text{CG}_*) \text{ hold.} \end{cases}$$

Then the iterates $u^i = (x^i, y^i)$ of the proximal point iteration (PP) satisfy

$$\delta \mathbb{E}[\|x^N - \hat{x}\|_{\Phi_N}^2] + \tilde{g}_N \leq \|u^0 - \hat{u}\|_{Z_0 L_0}^2 + \sum_{i=0}^{N-1} b_{i+2}(\Gamma) \quad (2.20)$$

Proof. We take the expectation of the estimate of Theorem 2.1, and directly use the estimates of Theorems 2.2 and 2.3. Then we use (C1) and (C2). \square

2.8. Interpreting the conditions

What do the conditions (C0), (C1) and (C2) mean? The condition (C1) is actually a stronger version of the positivity in (C0). If Φ_{i+1} is self-adjoint and positive definite, using (C0), for any $\delta \in (0, 1)$ we deduce

$$Z_{i+2}L_{i+2} \simeq \begin{pmatrix} \Phi_{i+1} & -\Lambda_{i+1}^* \\ -\Lambda_{i+1} & \Psi_{i+2} \end{pmatrix} \geq \begin{pmatrix} \delta\Phi_{i+1} & 0 \\ 0 & \Psi_{i+2} - \frac{1}{1-\delta}\Lambda_{i+1}\Phi_{i+1}^{-1}\Lambda_{i+1}^* \end{pmatrix}. \quad (2.21)$$

Thus the positivity in (C0) is verified along with (C1) if for some $\delta \in (0, 1)$ holds

$$(1 - \delta)\Psi_{i+1} \geq \Lambda_i \Phi_i^{-1} \Lambda_i^* \quad \text{and} \quad \Phi_i > 0, \quad (i \in \mathbb{N}). \quad (\text{C1}')$$

Regarding (C2), we expand

$$\|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\tilde{\Gamma})}^2 = \|u^{i+1} - u^i\|_{\Delta_{i+2}(\tilde{\Gamma})}^2 + 2\langle u^{i+1} - u^i, u^i - \hat{u} \rangle_{\Delta_{i+2}(\tilde{\Gamma})} + \|u^i - \hat{u}\|_{\Delta_{i+2}(\tilde{\Gamma})}^2. \quad (2.22)$$

Standard nesting properties of conditional expectations, since $u^{i+1} \in \mathcal{R}(O_i; X \times Y)$ and $u^i \in \mathcal{R}(O_{i-1}; X \times Y)$, show

$$\begin{aligned} \mathbb{E} \left[\|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\tilde{\Gamma})}^2 \right] &= \mathbb{E} \left[\|u^{i+1} - u^i\|_{\mathbb{E}[\Delta_{i+2}(\tilde{\Gamma})|O_i]}^2 \right. \\ &\quad + 2\langle u^{i+1} - u^i, u^i - \hat{u} \rangle_{\mathbb{E}[\Delta_{i+2}(\tilde{\Gamma})|O_i]} \\ &\quad \left. + \|u^i - \hat{u}\|_{\mathbb{E}[\Delta_{i+2}(\tilde{\Gamma})|O_{i-1}]}^2 \right]. \end{aligned} \quad (2.23)$$

Further using (2.11), we see for $O' = O_i, O_{i-1}$ that

$$\mathbb{E}[\Delta_{i+2}(\tilde{\Gamma})|O'] = \begin{pmatrix} \mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O'] & \mathbb{E}[A_{i+2}^*|O'] \\ \mathbb{E}[A_{i+2}|O'] & \mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O'] \end{pmatrix}.$$

If the off-diagonal satisfies

$$\mathbb{E}[A_{i+2}|O_i](x^{i+1} - x^i) = 0, \quad (2.24a)$$

$$\mathbb{E}[A_{i+2}^*|O_i](y^{i+1} - y^i) = 0, \quad \text{and} \quad (2.24b)$$

$$\mathbb{E}[A_{i+2}^*|O_{i-1}] = 0, \quad (2.24c)$$

we may simplify (2.23) into

$$\begin{aligned} \mathbb{E} \left[\|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\tilde{\Gamma})}^2 \right] &= \mathbb{E} \left[\|x^{i+1} - x^i\|_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_i]}^2 \right. \\ &\quad + \|y^{i+1} - y^i\|_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]}^2 \\ &\quad + 2\langle x^{i+1} - x^i, x^i - \hat{x} \rangle_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_i]} \\ &\quad + 2\langle y^{i+1} - y^i, y^i - \hat{y} \rangle_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]} \\ &\quad + \|x^i - \hat{x}\|_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_{i-1}]}^2 \\ &\quad \left. + \|y^i - \hat{y}\|_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_{i-1}]}^2 \right]. \end{aligned}$$

For any operator $M \in \mathcal{L}(X; X)$, defining $0 \leq |M| \in \mathcal{L}(X; X)$ as satisfying

$$\langle x, x' \rangle_M \leq \|x\|_{|M|} \|x'\|_{|M|}, \quad (x, x' \in M),$$

we obtain by Cauchy's inequality for any $\alpha_i, \beta_i > 0$ that

$$\begin{aligned} \mathbb{E} \left[\|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\tilde{\Gamma})}^2 \right] &\leq \mathbb{E} \left[\|x^{i+1} - x^i\|_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_i] + \alpha_i |\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_i]|}^2 \right. \\ &\quad + \|y^{i+1} - y^i\|_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i] + \beta_i |\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]|}^2 \\ &\quad + \|x^i - \hat{x}\|_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_{i-1}] + \alpha_i^{-1} |\mathbb{E}[\Phi_{i+1} - \Phi_i(I + 2T_i\tilde{\Gamma})|O_i]|}^2 \\ &\quad \left. + \|y^i - \hat{y}\|_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_{i-1}] + \beta_i^{-1} |\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]|}^2 \right]. \end{aligned} \quad (2.25)$$

Recalling the definition of A_{i+2} in (2.12), the off-diagonal conditions (2.24) may be written

$$\mathbb{E}[\Psi_{i+1}\Sigma_{i+1}K - \Lambda_{i+1} + \Lambda_i - KT_i^*\Phi_i^*|O_i](x^{i+1} - x^i) = 0, \quad (C2\prime.a)$$

$$\mathbb{E}[K^*\Sigma_{i+1}^*\Psi_{i+1}^* - \Lambda_{i+1}^* + \Lambda_i^* - \Phi_i T_i K|O_i]y^{i+1} - y^i = 0, \quad (C2\prime.b)$$

$$\mathbb{E}[\Psi_{i+1}\Sigma_{i+1}K - \Lambda_{i+1} + \Lambda_i - KT_i^*\Phi_i^*|O_{i-1}] = 0. \quad (C2\prime.c)$$

Using (2.21) on the right hand side, and (2.25) on the left hand side of (C2), we see that (C2) is satisfied if in addition to (C2'.a)–(C2'.c), we have some bounds $b_{i+2}^x(\tilde{\Gamma}), b_{i+2}^y \in \mathbb{R}$ such that

$$\mathbb{E}[\|x^{i+1} - x^i\|_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I+2T_i\tilde{\Gamma})|O_i] + \alpha_i|\mathbb{E}[\Phi_{i+1} - \Phi_i(I+2T_i\tilde{\Gamma})|O_i] - \delta\Phi_i}^2 + \|x^i - \hat{x}\|_{\mathbb{E}[\Phi_{i+1} - \Phi_i(I+2T_i\tilde{\Gamma})|O_{i-1}] + \alpha_i^{-1}|\mathbb{E}[\Phi_{i+1} - \Phi_i(I+2T_i\tilde{\Gamma})|O_i]}}^2] \leq b_{i+2}^x(\tilde{\Gamma}), \quad (\text{C2'.d})$$

and

$$\mathbb{E}[\|y^{i+1} - y^i\|_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i] + \beta_i|\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]}}^2 + \|y^i - \hat{y}\|_{\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_{i-1}] + \beta_i^{-1}|\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]}}^2] \leq b_{i+2}^y. \quad (\text{C2'.e})$$

The bounds $b_{i+2}^x(\tilde{\Gamma})$ and b_{i+2}^y clearly exist if the iterates stay bounded, or if the conditional operator expectations stay negative semi-definite. To achieve the latter, we need to choose the testing operators Φ_i and Ψ_i suitably. We do this in the following sections for practical block-coordinate descent algorithms. First, to conclude the present derivations:

Corollary 2.2. *We may in Corollary 2.1 replace (C1) and (C2) by (C1') and (C2').*

3. Block-proximal methods

To derive practical algorithms, we need to satisfy the conditions of Corollary 2.2. To do this, we employ ideas from the stochastic block-coordinate descent methods discussed in the Introduction (Section 1). We construct T_i and Σ_{i+1} to select some sub-blocks of the variables x and y to update. This way, we also seek to gain performance improvements through the local properties of G , F^* , and K .

3.1. Structure of the step length operators

Let P_1, \dots, P_m be a collection of projection operators in X , with $\sum_{j=1}^m P_j = I$ and $P_j P_i = 0$ if $i \neq j$. Likewise, suppose Q_1, \dots, Q_n are projection operators in Y such that $\sum_{\ell=1}^n Q_\ell = I$ and $Q_\ell Q_k = 0$ for $k \neq \ell$. With this, for $j \in \{1, \dots, n\}$ and a subset $S \subset \{1, \dots, m\}$, we denote

$$\mathcal{V}(j) := \{\ell \in \{1, \dots, n\} \mid Q_\ell K P_j \neq 0\}, \quad \text{and} \quad \mathcal{V}(S) = \bigcup_{j \in S} \mathcal{V}(j).$$

For some $\{\phi_{j,i}\}_{j=1}^m, \{\psi_{\ell,i+1}\}_{\ell=1}^n \subset \mathbb{R}(O_i; (0, \infty))$, we then define

$$\Phi_i := \sum_{j=1}^m \phi_{j,i} P_j, \quad \text{and} \quad \Psi_{i+1} := \sum_{\ell=1}^n \psi_{\ell,i+1} Q_\ell. \quad (\text{S-}\Phi\Psi)$$

We take random subsets $S(i) \subset \{1, \dots, m\}$, and $V(i+1) \subset \{1, \dots, n\}$, deterministic with respect to O_i , and set

$$T_i := \sum_{j \in S(i)} \tau_{j,i} P_j, \quad \text{and} \quad \Sigma_{i+1} := \sum_{\ell \in V(i+1)} \sigma_{\ell,i+1} Q_\ell, \quad (i \geq 0). \quad (\text{S-}T\Sigma)$$

We assume that the blockwise step lengths satisfy $\{\tau_{j,i+1}\}_{j=1}^m, \{\sigma_{\ell,i+2}\}_{\ell=1}^n \subset \mathbb{R}(O_i; (0, \infty))$. Then all Φ_i , Ψ_{i+1} , T_i , and Σ_{i+1} are self-adjoint and positive semi-definite.

Finally, we introduce $\hat{\tau}_{j,i} := \tau_{j,i} \chi_{S(i)}(j)$ and $\hat{\sigma}_{\ell,i} := \sigma_{\ell,i} \chi_{V(i)}(\ell)$, as well as denote by

$$\pi_{j,i} := \mathbb{P}[j \in S(i) \mid O_{i-1}], \quad \text{and} \quad \nu_{\ell,i+1} := \mathbb{P}[\ell \in V(i+1) \mid O_{i-1}],$$

the probability that j will be contained in $S(i)$ and, respectively, that ℓ will be contained in $V(i+1)$, given what is known at iteration $i-1$.

3.2. Structure of G and F^*

We assume that G and F^* are (block-)separable in the sense

$$G(x) = \sum_{j=1}^m G_j(P_j x), \quad (\text{S-G})$$

and

$$F^*(y) = \sum_{\ell=1}^n F_\ell^*(Q_\ell y). \quad (\text{S-F}^*)$$

With $\mathcal{T} := \{\sum_{j=1}^m t_j P_j \mid t_j > 0\}$, and $\mathcal{S} := \{\sum_{\ell=1}^n s_\ell Q_\ell \mid s_\ell > 0\}$, the conditions (G-PM) and (F*-PM) then reduce to the strong monotonicity of each ∂G_j with factor γ_j , and the monotonicity of each ∂F_ℓ . In particular (F*-PM) is automatically satisfied. Thus $\Gamma = \sum_{j=1}^m \gamma_j P_j$. We also write $\tilde{\Gamma} = \sum_{j=1}^m \tilde{\gamma}_j P_j$.

Let us temporarily introduce $\tilde{T}_i := \sum_{j=1}^m \tilde{\gamma}_{j,i} \geq 0$, satisfying $\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\gamma}_{j,i}] = 1$ for each $j = 1, \dots, m$. Splitting (G-EC) into separate inequalities over all $j = 1, \dots, m$, and using the strong convexity of G_j , we see that (G-EC) holds if

$$G_j \left(\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\gamma}_{j,i} P_j x^{i+1}] \right) - G_j(P_j \hat{x}) \geq \sum_{i=0}^{N-1} \mathbb{E}[\tilde{\gamma}_i (G_j(x^{i+1}) - G_j(\hat{x}))], \quad (j = 1, \dots, m). \quad (3.1)$$

The right hand side can also be written as $\int_{\Omega^N} G_j(P_j x^i(\omega)) - G_j(P_j \hat{x}) d\mu^N(i, \omega)$ for the measure $\mu^N := \tilde{\gamma}_j \sum_{i=0}^{N-1} \delta_i \times \mathbb{P}$ on the domain $\Omega^N := \{0, \dots, N-1\} \times \Omega$. Using our assumption $\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\gamma}_{j,i}] = 1$, we deduce $\mu^N(\Omega^N) = 1$. An application of Jensen's inequality now shows (3.1). Therefore (G-EC) is automatically satisfied. By similar arguments we see that (F*-EC) also holds.

We now need to satisfy the conditions of Corollary 2.2 for the above structural setup. Namely, we need to satisfy (C0), (C1'), and (C2') to obtain convergence estimates of the primal iterates, and we need to satisfy either (CG) or (CG*) to obtain gap estimates. We divide these verifications into the following subsections Sections 3.3 to 3.6, after which we summarise the results in Section 3.7.

3.3. Satisfaction of the off-diagonal conditions (C2'.a)–(C2'.c) and either (CG) or (CG*)

Expanded, L_{i+1} solved from (C0), and the proximal maps inverted, (PP) states

$$x^{i+1} = (I + T_i \partial G)^{-1} (x^i + \Phi_i^{-1} \Lambda_i^* (y^{i+1} - y^i) - T_i K^* y^{i+1}), \quad (3.2a)$$

$$y^{i+1} = (I + \Sigma_{i+1} \partial F^*)^{-1} (y^i + \Psi_{i+1}^{-1} \Lambda_i (x^{i+1} - x^i) + \Sigma_{i+1} K x^{i+1}). \quad (3.2b)$$

To derive an efficient algorithm, we have to be able to solve this system easily. In particular, we wish to avoid any cross-dependencies on x^{i+1} and y^{i+1} between the two steps. One way to avoid this, is to make the first step independent of y^{i+1} . To do this, we could enforce $\Phi_i^{-1} \Lambda_i^* = T_i K^*$, but this can be refined a little bit, as not all blocks of y^{i+1} are updated.

Moreover, to simplify the treatment of (C2'.a) and (C2'.b), and for $S(i)$ and $V(i+1)$ to correspond exactly to the coordinates that are updated, as one would expect, we enforce

$$x_j^{i+1} = x_j^i, \quad (j \notin S(i)), \quad \text{and likewise} \quad (\text{C-cons.a})$$

$$y_\ell^{i+1} = y_\ell^i, \quad (\ell \notin V(i+1)). \quad (\text{C-cons.b})$$

Let us take

$$\Lambda_i := \sum_{j=1}^m \sum_{\ell \in \mathcal{V}(j)} \lambda_{\ell,j,i} Q_\ell K P_j \quad \text{for some } \lambda_{\ell,j,i} \in \mathcal{R}(\mathcal{O}_i, [0, \infty)). \quad (\text{S-}\Lambda)$$

If (C-cons) and (S- Λ) hold, then (C2'.a)–(C2'.c) follow if

$$\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] = \psi_{\ell,i+1} \hat{\sigma}_{\ell,i+1} + \lambda_{\ell,j,i} - \phi_{j,i} \hat{\tau}_{j,i}, \quad \begin{cases} j \in S(i), \\ \ell \in \mathcal{V}(j), \end{cases} \quad (3.3a)$$

$$\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] = \psi_{\ell,i+1} \hat{\sigma}_{\ell,i+1} + \lambda_{\ell,j,i} - \phi_{j,i} \hat{\tau}_{j,i}, \quad \begin{cases} \ell \in V(i+1), \\ j \in \mathcal{V}^{-1}(\ell), \end{cases} \quad (3.3b)$$

$$\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_{i-1}] = \mathbb{E}[\psi_{\ell,i+1} \hat{\sigma}_{\ell,i+1} + \lambda_{\ell,j,i} - \phi_{j,i} \hat{\tau}_{j,i} | \mathcal{O}_{i-1}], \quad \begin{cases} j = 1, \dots, m, \\ \ell \in \mathcal{V}(j). \end{cases} \quad (3.3c)$$

We set

$$\tilde{\lambda}_{\ell,j,i+1} := \psi_{\ell,i+1} \hat{\sigma}_{\ell,i+1} + \lambda_{\ell,j,i} - \phi_{j,i} \hat{\tau}_{j,i},$$

and using (3.3a) and (3.3b), compute

$$\begin{aligned} \mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_{i-1}] &= \mathbb{E}[\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] | \mathcal{O}_{i-1}] \\ &= \mathbb{E}[\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] \chi_{V(i+1)}(\ell) | \mathcal{O}_{i-1}] \\ &\quad + \mathbb{E}[\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] (1 - \chi_{V(i+1)}(\ell)) \chi_{S(i)}(j) | \mathcal{O}_{i-1}] \\ &\quad + \mathbb{E}[\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] (1 - \chi_{V(i+1)}(\ell)) (1 - \chi_{S(i)}(j)) | \mathcal{O}_{i-1}] \\ &= \mathbb{E}[\tilde{\lambda}_{\ell,j,i+1} \chi_{V(i+1)}(\ell) | \mathcal{O}_{i-1}] \\ &\quad + \mathbb{E}[\tilde{\lambda}_{\ell,j,i+1} (1 - \chi_{V(i+1)}(\ell)) \chi_{S(i)}(j) | \mathcal{O}_{i-1}] \\ &\quad + \mathbb{E}[\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] (1 - \chi_{V(i+1)}(\ell)) (1 - \chi_{S(i)}(j)) | \mathcal{O}_{i-1}]. \end{aligned} \quad (3.4)$$

If

$$\lambda_{\ell,j,i} = 0, \quad (j \notin S(i) \text{ or } \ell \notin V(i+1)), \quad (3.5a)$$

we obtain

$$\begin{aligned} \mathbb{E}[\tilde{\lambda}_{\ell,j,i+1} | \mathcal{O}_{i-1}] &= \mathbb{E}[\tilde{\lambda}_{\ell,j,i+1} \chi_{V(i+1)}(\ell) | \mathcal{O}_{i-1}] \\ &\quad + \mathbb{E}[\tilde{\lambda}_{\ell,j,i+1} (1 - \chi_{V(i+1)}(\ell)) \chi_{S(i)}(j) | \mathcal{O}_{i-1}]. \end{aligned}$$

Together this and (3.4) show that (3.3c) holds if (3.3a) and (3.3b) do along with

$$\mathbb{E}[\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] (1 - \chi_{V(i+1)}(\ell)) (1 - \chi_{S(i)}(j)) | \mathcal{O}_{i-1}] = 0.$$

From this, it is easy to see that (3.3) holds if (3.5a) does, and

$$\mathbb{E}[\lambda_{\ell,j,i+1} | \mathcal{O}_i] = \psi_{\ell,i+1} \hat{\sigma}_{\ell,i+1} + \lambda_{\ell,j,i} - \phi_{j,i} \hat{\tau}_{j,i}, \quad (j = 1, \dots, m; \ell \in \mathcal{V}(j)). \quad (3.5b)$$

For a specific choice of $\lambda_{\ell,j,i}$, we take $\mathring{S}(i) \subset S(i)$, $\mathring{V}(i+1) \subset V(i+1)$, and set

$$\lambda_{\ell,j,i} := \phi_{j,i} \hat{\tau}_{j,i} \chi_{\mathring{S}(i)}(j) - \psi_{\ell,i+1} \hat{\sigma}_{\ell,i+1} \chi_{\mathring{V}(i+1)}(\ell), \quad (\ell \in \mathcal{V}(j)). \quad (\text{R-}\lambda)$$

For $\lambda_{\ell,j,i}$ to take values that allow the straightforward computation of (3.2), we assume

$$\mathcal{V}^{-1}(\mathring{V}(i+1)) \cap \mathcal{V}^{-1}(\mathcal{V}(\mathring{S}(i))) = 0. \quad (\text{C-nest.a})$$

It is then easy to see that (C-cons) and the easy computability of (3.2) demand

$$S(i) = \mathring{S}(i) \cup \mathcal{V}^{-1}(\mathring{V}(i+1)), \quad \text{and} \quad V(i+1) = \mathring{V}(i+1) \cup \mathcal{V}(\mathring{S}(i)). \quad (\text{C-nest.b})$$

Assuming (C-nest) to hold, clearly the choice (R- λ) satisfies (3.5a), while (3.5b) follows if

$$\mathbb{E}[\lambda_{\ell,j,i+1}|\mathcal{O}_i] = \psi_{\ell,i+1}\hat{\sigma}_{\ell,i+1}(1 - \chi_{\hat{V}(i+1)}^\circ(\ell)) - \phi_{j,i}\hat{\tau}_{j,i}(1 - \chi_{\hat{S}(i)}^\circ(j)), \quad (\ell \in \mathcal{V}(j)).$$

Inserting (R- λ), we see this to be satisfied if for some $\eta_{i+1} \in \mathcal{R}(\mathcal{O}_i; (0, \infty))$ holds

$$\mathbb{E}[\phi_{j,i+1}\hat{\tau}_{j,i+1}\chi_{\hat{S}(i+1)}^\circ(j)|\mathcal{O}_i] = \eta_{i+1} - \phi_{j,i}\hat{\tau}_{j,i}(1 - \chi_{\hat{S}(i)}^\circ(j)) \geq 0, \quad \text{and} \quad (\text{C-step.a})$$

$$\mathbb{E}[\psi_{\ell,i+2}\hat{\sigma}_{\ell,i+2}\chi_{\hat{V}(i+2)}^\circ(\ell)|\mathcal{O}_i] = \eta_{i+1} - \psi_{\ell,i+1}\hat{\sigma}_{\ell,i+1}(1 - \chi_{\hat{V}(i+1)}^\circ(\ell)) \geq 0, \quad (\text{C-step.b})$$

with $j = 1, \dots, m$; $\ell = 1, \dots, n$; and $i \geq -1$, taking

$$\dot{S}(-1) = \{1, \dots, m\}, \quad \text{and} \quad \dot{V}(0) = \{1, \dots, n\}. \quad (\text{C-step.c})$$

If the testing variables $\phi_{j,i+1}$ and $\psi_{\ell,i+2}$ are known, (C-step.a) and (C-step.b) give update rules for $\tau_{j,i+1}$ and $\sigma_{\ell,i+2}$ when $j \in \dot{S}(i+1)$ and, respectively, $\ell \in \dot{V}(i+2)$. To cover $j \in S(i+1) \setminus \dot{S}(i+1)$ and $\ell \in V(i+2) \setminus \dot{V}(i+2)$, for some $\eta_{\tau,i+1}^\perp, \eta_{\sigma,i+1}^\perp \in \mathcal{R}(\mathcal{O}_i; [0, \infty))$, we demand

$$\mathbb{E}[\phi_{j,i+1}\hat{\tau}_{j,i+1}(1 - \chi_{\hat{S}(i+1)}^\circ(j))|\mathcal{O}_i] = \eta_{\tau,i+1}^\perp, \quad \text{and} \quad (\text{C-step.d})$$

$$\mathbb{E}[\psi_{\ell,i+2}\hat{\sigma}_{\ell,i+2}(1 - \chi_{\hat{V}(i+2)}^\circ(\ell))|\mathcal{O}_i] = \eta_{\sigma,i+1}^\perp. \quad (\text{C-step.e})$$

Then

$$\mathbb{E}[\phi_{j,i+1}\hat{\tau}_{j,i+1}] = \mathbb{E}[\eta_{i+1} + \eta_{\tau,i+1}^\perp - \eta_{\tau,i}^\perp], \quad \text{and} \quad (3.6a)$$

$$\mathbb{E}[\psi_{\ell,i+2}\hat{\sigma}_{\ell,i+2}] = \mathbb{E}[\eta_{i+1} + \eta_{\sigma,i+1}^\perp - \eta_{\sigma,i}^\perp]. \quad (3.6b)$$

The condition (C \mathcal{G}) is satisfied if $\mathbb{E}[\phi_{j,i+1}\hat{\tau}_{j,i+1}] = \bar{\eta}_i = \mathbb{E}[\psi_{\ell,i+2}\hat{\sigma}_{\ell,i+2}]$. This follows if

$$\mathbb{E}[\eta_{\tau,i}^\perp - \eta_{\sigma,i}^\perp] = \eta^\perp := \text{constant}. \quad (\text{C-}\eta^\perp)$$

We therefore obtain the following.

Lemma 3.1. *Let us choose Λ_i according to (S- Λ) and (R- λ). If (C-nest), (C-step), and (C- η^\perp) are satisfied, then (C2 ℓ .a)–(C2 ℓ .c) and (C \mathcal{G}) hold, as does (C-cons).*

We will frequently assume for some $\epsilon \in (0, 1)$ that

$$\left. \begin{array}{l} i \mapsto \eta_i \\ i \mapsto \eta_{\tau,i}^\perp \\ i \mapsto \eta_{\sigma,i}^\perp \end{array} \right\} \text{ are non-decreasing, and } \left\{ \begin{array}{l} \epsilon \eta_i \cdot \min_j (\pi_{j,i} - \hat{\pi}_{j,i}) \geq \eta_{\tau,i}^\perp, \\ \eta_i \cdot \min_\ell (\nu_{\ell,i+1} - \hat{\nu}_{\ell,i+1}) \geq \eta_{\sigma,i}^\perp. \end{array} \right. \quad (\text{C-}\eta)$$

This ensures the non-negativity conditions in (C-step.a) and (C-step.b), while simplifying other derivations to follow. In fact, the condition suggests to take either both $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$ as constants, or to take $\eta_{\sigma,i+1}^\perp = \eta_{\tau,i}^\perp =: c\eta_i$ for some $c > 0$ such that the non-negativity conditions in (C-step.a) and (C-step.b) are satisfied. Note that (C- η) guarantees $\bar{\eta}_i \geq \mathbb{E}[\eta_i]$.

If we deterministically take $\dot{V}(i+1) = \emptyset$, then (C-step.e) implies $\eta_{\sigma,i}^\perp \equiv 0$. But then (C-step.b) will be incompatible with (3.6b). Therefore $\dot{V}(i+1)$ has to be chosen randomly to satisfy (C \mathcal{G}). The same holds for $\dot{S}(i)$. Thus algorithms satisfying (C \mathcal{G}) are necessarily doubly-stochastic, randomly updating both the primal and dual variables, or neither.

The alternative (C \mathcal{G}_*) requires $\mathbb{E}[\phi_{j,i+1}\hat{\tau}_{j,i+1}] = \bar{\eta}_{i+1} = \mathbb{E}[\psi_{\ell,i+1}\hat{\sigma}_{\ell,i+1}]$. This holds when

$$\mathbb{E}[\eta_{i+1} + \eta_{\tau,i+1}^\perp - \eta_{\tau,i}^\perp] = \bar{\eta}_i = \mathbb{E}[\eta_i + \eta_{\sigma,i}^\perp - \eta_{\sigma,i-1}^\perp].$$

It does not appear possible to simultaneously satisfy this condition and the non-negativity conditions in (C-step.a) and (C-step.b) for an accelerated method, unless we deterministically take $\dot{V}(i+1) = \emptyset$ for all i . Then (C-step) and (C-nest) imply $V(i+1) = \{1, \dots, n\}$, $S(i) = \dot{S}(i)$, as well as $\eta_{\tau,i}^\perp \equiv 0$, and $\eta_{\sigma,i+1}^\perp = \eta_i$. This choice satisfies (C- η) if $i \mapsto \eta_i$ is non-decreasing and positive. Conversely, choosing $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$ this way, and taking the expectation with respect to O_{i-1} in (C-step.b), we see that $\dot{V}(i+1) = \emptyset$. This says that to satisfy (C \mathcal{G}_*), we need to perform full dual updates. This is akin to most existing primal-dual coordinate descent methods [32, 34, 35]. The algorithms in [36–38] are more closely related to our method. However only [38] provides convergence rates for very limited single-block sampling schemes under the strong assumption that both G and F^* are strongly convex.

The conditions (C-step) now reduce to

$$\phi_{j,i+1}\tau_{j,i+1}\pi_{j,i+1} = \eta_{i+1}, \quad (j = 1, \dots, m), \quad \text{and} \quad (\text{C-step'.a})$$

$$\psi_{\ell,i+1}\sigma_{\ell,i+1} = \eta_{i+1}, \quad (\ell = 1, \dots, n). \quad (\text{C-step'.b})$$

Moreover $\lambda_{\ell,j,i} = \phi_{j,i}\hat{\tau}_{j,i}\chi_{\dot{S}(i)}(j)$. Clearly this satisfies (C-cons) through (3.2). In summary:

Lemma 3.2. *Let us choose Λ_i according to (S-A) and (R- λ). If we ensure (C-step'), take $\eta_{\tau,i}^\perp \equiv 0$ and $\eta_{\sigma,i+1}^\perp = \eta_i$, and force*

$$\dot{S}(i) = S(i), \quad \dot{V}(i+1) = \emptyset, \quad \text{and} \quad V(i+1) = \{1, \dots, n\},$$

then (C-nest) and (C-step) hold, as do (C2'.a)–(C2'.c), (C \mathcal{G}_), and (C-cons). If $i \mapsto \eta_i > 0$ is non-decreasing, then (C- η) holds.*

3.4. Satisfaction of the primal penalty bound (C2'.d)

Split into blocks, the conditions asks for each $j = 1, \dots, m$ the upper bound

$$\mathbb{E}[q_{j,i+2}(\tilde{Y}_j)\|x_j^{i+1} - x_j^i\|^2 + h_{j,i+2}(\tilde{Y}_j)\|x_j^i - \hat{x}_j\|^2] \leq b_{j,i+2}^x(\tilde{Y}_j), \quad (3.7)$$

where

$$\begin{aligned} q_{j,i+2}(\tilde{Y}_j) &:= (\mathbb{E}[\phi_{j,i+1} - \phi_{j,i}(1 + 2\hat{\tau}_{j,i}\tilde{Y}_j)|O_i] \\ &\quad + \alpha_i|\mathbb{E}[\phi_{j,i+1} - \phi_{j,i}(1 + 2\hat{\tau}_{j,i}\tilde{Y}_j)|O_i] - \delta\phi_{j,i})\chi_{S(i)}(j), \end{aligned} \quad (3.8)$$

and

$$\begin{aligned} h_{j,i+2}(\tilde{Y}_j) &:= \mathbb{E}[\phi_{j,i+1} - \phi_{j,i}(1 + 2\hat{\tau}_{j,i}\tilde{Y}_j)|O_{i-1}] \\ &\quad + \alpha_i^{-1}|\mathbb{E}[\phi_{j,i+1} - \phi_{j,i}(1 + 2\hat{\tau}_{j,i}\tilde{Y}_j)|O_i]|. \end{aligned} \quad (3.9)$$

We easily obtain from (3.7) the following lemma.

Lemma 3.3. *Suppose for some $C_x > 0$ either*

$$\|x_j^{i+1} - \hat{x}_j\|^2 \leq C_x, \quad \text{or} \quad (\text{C-xbnd.a})$$

$$h_{j,i+2}(\tilde{Y}_j) \leq 0 \quad \text{and} \quad q_{j,i+2}(\tilde{Y}_j) \leq 0, \quad (\text{C-xbnd.b})$$

for all $j = 1, \dots, m$ and $i \in \mathbb{N}$. Then (C2'.d) holds with $b_{i+2}^x(\tilde{\Gamma}) := \sum_{j=1}^m b_{j,i+2}^x(\tilde{Y}_j)$ for any

$$b_{j,i+2}^x(\tilde{Y}_j) \geq 4C_x\mathbb{E}[\max\{0, q_{j,i+2}(\tilde{Y}_j)\}] + C_x\mathbb{E}[\max\{0, h_{j,i+2}(\tilde{Y}_j)\}]. \quad (3.10)$$

Since Corollary 2.2 involve sums $\sum_{i=0}^{N-1} b_{i+2}^x(\tilde{\Gamma})$, we define for convenience

$$d_{j,N}^x(\tilde{Y}_j) := \sum_{i=0}^{N-1} b_{j,i+2}^x(\tilde{Y}_j). \quad (3.11)$$

We still need to bound $\mathbb{E}[\max\{0, q_{j,i+2}(\tilde{Y}_j)\}]$ and $\mathbb{E}[\max\{0, h_{j,i+2}(\tilde{Y}_j)\}]$. We do this through primal test update rules, constructing next two possibilities.

Example 3.1 (Random primal test updates). For some constant $\rho_j \geq 0$, let us take

$$\phi_{j,i+1} := \phi_{j,i}(1 + 2\tilde{Y}_j\hat{\tau}_{j,i}) + 2\rho_j\pi_{j,i}^{-1}\chi_{S(i)}(j). \quad (\text{R-}\phi\text{rnd})$$

Note that $\phi_{j,i+1} \in \mathcal{R}(\mathcal{O}_i; (0, \infty))$ instead of just $\mathcal{R}(\mathcal{O}_{i+1}; (0, \infty))$, as we have assumed so far. If we set $\rho_j = 0$ and have just a single deterministically updated block, (R- ϕ rnd) is gives the standard update rule (2.3) with the identification $\phi_i = \tau_i^{-2}$. The role of $\rho_j > 0$ is to ensure some (slower) acceleration for non-strongly-convex blocks with $\tilde{Y}_j = 0$. This is necessary for convergence rate estimates.

We compute

$$q_{j,i+2}(\tilde{Y}_j) = (2(1 + \alpha_i)\rho_j\pi_{j,i}^{-1} - \delta\phi_{j,i})\chi_{S(i)}(j) \leq 2(1 + \alpha_i)\rho_j\pi_{j,i}^{-1}\chi_{S(i)}(j), \quad \text{and} \quad (3.12a)$$

$$h_{j,i+2}(\tilde{Y}_j) = 2\rho_j\pi_{j,i}^{-1}\mathbb{E}[\chi_{S(i)}(j)|\mathcal{O}_{i-1}] + 2\alpha_i^{-1}\rho_j\pi_{j,i}^{-1}\mathbb{E}[\chi_{S(i)}(j)|\mathcal{O}_i]. \quad (3.12b)$$

This implies that (C-xbnd.b) holds if $\rho_j = 0$. Taking the expectation, we compute

$$\mathbb{E}[\max\{0, q_{j,i+2}(\tilde{Y}_j)\}] \leq 2(1 + \alpha_i)\rho_j, \quad \text{and} \quad (3.13a)$$

$$\mathbb{E}[\max\{0, h_{j,i+2}(\tilde{Y}_j)\}] = 2(1 + \alpha_i^{-1})\rho_j. \quad (3.13b)$$

Choosing $\alpha_i = 1/2$ and using (3.13) in (3.10), we obtain:

Lemma 3.4. *If (C-xbnd.a) holds, take $\rho_j \geq 0$, otherwise take $\rho_j = 0$, ($j = 1, \dots, m$). If we define $\phi_{j,i+1} \in \mathcal{R}(\mathcal{O}_i; (0, \infty))$ through (R- ϕ rnd), then we may take $b_{j,i+2}^x(\tilde{Y}_j) = 18C_x\rho_j$. In particular (C2'.d) holds, and $d_{j,N}^x(\tilde{Y}_j) = 18C_x\rho_jN$.*

Remark 3.1. In (R- ϕ rnd), we could easily replace $r_{j,i} = \rho_j\pi_{j,i}^{-1}\chi_{S(i)}(j)$ by a split version $r_{j,i} = \dot{\rho}_j\dot{\pi}_{j,i}^{-1}\chi_{\dot{S}(i)}(j) + \rho_j^\perp(\pi_{j,i} - \dot{\pi}_{j,i})^{-1}\chi_{S(i)\setminus\dot{S}(i)}(j)$ without destroying the property $\mathbb{E}[r_{j,i}|\mathcal{O}_{i-1}] = \dot{\rho}_j + \rho_j^\perp =: \rho_j$.

The difficulty with the rule (R- ϕ rnd) is, as we will see in Section 4, that η_{i+1} will depend on the random realisations of $S(i)$ through $\phi_{j,i+1}$. This will require communication in a parallel implementation of the algorithm. We therefore desire a deterministic update rule for η_{i+1} . As we will see, this can be achieved if $\phi_{j,i+1}$ is updated deterministically.

Example 3.2 (Deterministic primal test updates). Let us assume (C-step) and (C- η) to hold, and for some $\rho_j \geq 0$ and $\tilde{Y}_j \in [0, 1)$ take

$$\phi_{j,i+1} := \phi_{j,i} + 2(\tilde{Y}_j\eta_i + \rho_j). \quad (\text{R-}\phi\text{det})$$

Since $\eta_i \in \mathcal{R}(\mathcal{O}_{i-1}; (0, \infty))$, we see that $\phi_{j,i+1} \in \mathcal{R}(\mathcal{O}_{i-1}; (0, \infty))$. In fact, $\phi_{j,i+1}$ is deterministic as long as η_i is chosen deterministically, for example as a function of $\{\phi_{j,i}\}_{i=1}^m$. Since $i \mapsto \eta_{\tau,i}^\perp$ is non-decreasing by (C- η), (C-step) gives

$$\mathbb{E}[\phi_{j,i}\hat{\tau}_{j,i}|\mathcal{O}_{i-1}] = \eta_i + \eta_{\tau,i}^\perp - \eta_{i-1,\tau}^\perp \geq \eta_i. \quad (3.14)$$

Abbreviating $\gamma_{j,i} := \bar{\gamma}_j + \rho_j \eta_i^{-1}$, we can write $\phi_{j,i+1} = \phi_{j,i} + 2\gamma_{j,i}\eta_i$. With this, expansion of (3.9) gives

$$\begin{aligned} h_{j,i+2}(\tilde{\gamma}_j) &= 2\mathbb{E}[\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i} | \mathcal{O}_{i-1}] + 2\alpha_i^{-1}|\mathbb{E}[\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i} | \mathcal{O}_i]| \\ &\leq 2(\gamma_{j,i} - \tilde{\gamma}_j)\eta_i + 2\alpha_i^{-1}|\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}| \\ &\leq 2(1 + \alpha_i^{-1})\rho_j + 2(\bar{\gamma}_j - \tilde{\gamma}_j)\eta_i + 2\alpha_i^{-1}|\bar{\gamma}_j\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}|. \end{aligned}$$

Forcing

$$\alpha_i^{-1}|\bar{\gamma}_j\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}| \leq (\tilde{\gamma}_j - \bar{\gamma}_j)\eta_i, \quad (3.15)$$

and taking the expectation, this gives

$$\mathbb{E}[\max\{0, h_{j,i+2}(\tilde{\gamma}_j)\}] \leq 2(1 + \alpha_i^{-1})\rho_j. \quad (3.16)$$

If $\bar{\gamma}_j\eta_i > \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}$, (3.15) holds when we take

$$\alpha_i = \alpha_{i,1} := \min_j \bar{\gamma}_j / (\tilde{\gamma}_j - \bar{\gamma}_j). \quad (3.17)$$

Otherwise, if $\bar{\gamma}_j\eta_i \leq \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}$, for (3.15) to hold, we need

$$\phi_{j,i}\hat{\tau}_{j,i} \leq \left(\alpha_i \frac{\tilde{\gamma}_j - \bar{\gamma}_j}{\tilde{\gamma}_j} - \frac{\bar{\gamma}_j}{\tilde{\gamma}_j} \right) \eta_i. \quad (3.18)$$

Taking

$$\alpha_i = \alpha_{i,2} := \max_j \left(\frac{\tilde{\gamma}_j \pi_{j,i}^{-1} - \bar{\gamma}_j}{\tilde{\gamma}_j - \bar{\gamma}_j} \right),$$

we see that (3.18) holds if $\phi_{j,i}\hat{\tau}_{j,i} \leq \pi_{j,i}^{-1}\eta_i$. We have to consider the cases $j \in \dot{S}(i)$ and $j \in S(i) \setminus \dot{S}(i)$ separately. The conditions (C-step.a) and (C-step.d) show that

$$\phi_{j,i}\hat{\tau}_{j,i}\pi_{j,i}\chi_{\dot{S}(i)}(j) \leq \eta_i, \quad \text{and} \quad \phi_{j,i}\hat{\tau}_{j,i}(\pi_{j,i} - \pi_{j,i}^*)(1 - \chi_{\dot{S}(i)}(j)) \leq \eta_{\tau,i}^\perp.$$

Using (C-η) in the latter estimate, we verify (3.18) (provided or not that $\pi_{j,i}^* > 0$).

Next, we expand (3.8), obtaining

$$\begin{aligned} q_{j,i+2}(\tilde{\gamma}_j) &= (2\mathbb{E}[\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i} | \mathcal{O}_i] + 2\alpha_i|\mathbb{E}[\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i} | \mathcal{O}_i]| - \delta\phi_{j,i})\chi_{S(i)}(j), \\ &= (2(\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}) + 2\alpha_i|\gamma_{j,i}\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}| - \delta\phi_{j,i})\chi_{S(i)}(j), \\ &\leq (2(1 + \alpha_i)\rho_j + 2(\bar{\gamma}_j\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}) + 2\alpha_i|\bar{\gamma}_j\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}| - \delta\phi_{j,i})\chi_{S(i)}(j). \end{aligned}$$

Again, as η_i and $\phi_{j,i}\tau_{j,i}$ will be increasing, we want

$$2(\bar{\gamma}_j\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}) + 2\alpha_i|\bar{\gamma}_j\eta_i - \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}| \leq \delta\phi_{j,i}, \quad (j \in S(i)). \quad (3.19)$$

Then

$$\mathbb{E}[q_{j,i+2}(\tilde{\gamma}_j)] \leq 2(1 + \alpha_i)\rho_j. \quad (3.20)$$

We only need to consider the case $\bar{\gamma}_j\eta_i > \tilde{\gamma}_j\phi_{j,i}\hat{\tau}_{j,i}$, as (3.19) is trivial in the opposite case. Then α_i is given by (3.17). With this the condition (3.19) expands into

$$\tilde{\gamma}_j = \bar{\gamma}_j = 0 \quad \text{or} \quad \frac{2\tilde{\gamma}_j\bar{\gamma}_j}{\tilde{\gamma}_j - \bar{\gamma}_j}\eta_i \leq \delta\phi_{j,i}, \quad (j \in S(i), i \in \mathbb{N}). \quad (\text{C-}\phi\text{det})$$

In summary:

Lemma 3.5. Suppose (C-step), (C- η), and (C- ϕ det) hold. If (C-xbnd.a) holds, take $\rho_j \geq 0$, otherwise take $\rho_j = 0$, ($j = 1, \dots, m$). Let $\phi_{j,i+1} \in \mathcal{R}(O_{i-1}; (0, \infty))$ be defined by (R- ϕ det). Then we may take

$$b_{j,i+2}^x(\tilde{Y}_j) = \max_{\alpha=\alpha_{i,1}, \alpha_{i,2}} (2(1 + \alpha^{-1})\rho_j C_x + 8(1 + \alpha)\rho_j C_x).$$

In particular, if $\pi_{j,i} \geq \epsilon > 0$ for all $i \in \mathbb{N}$, and some $\epsilon > 0$, then there exists a constant $C_\alpha > 0$ such that (C2'.d) holds with $d_{j,N}^x(\tilde{Y}_j) = \rho_j C_x C_\alpha N$.

Proof. We see from (3.16) and (3.20) that (C-xbnd.b) holds if we take $\rho_j = 0$. Therefore (C-xbnd) always holds. The expression for $b_{j,i+2}^x(\tilde{Y}_j)$ now follows from Lemma 3.3. For the expression of $d_{j,N}^x(\tilde{Y}_j)$, we note that the condition $\pi_{j,i} \geq \epsilon > 0$ bounds $\alpha_{i,2}$. \square

Remark 3.2. In the rule (R- ϕ det), we could replace η_i by $\eta_i + \eta_{\tau,i}^\perp - \eta_{i-1,\tau}^\perp$; cf. (3.14).

3.5. Satisfaction of the dual penalty bound (C2'.e)

To satisfy this bound, we make assumptions similar to Lemma 3.3.

Lemma 3.6. Suppose

$$\mathbb{E}[\psi_{\ell,i+2}|O_i] \geq \mathbb{E}[\psi_{\ell,i+1}|O_i], \quad (\ell = 1, \dots, n). \quad (\text{C-}\psi\text{inc})$$

as well as either

$$\|y_\ell^{i+1} - \hat{y}_\ell\|^2 \leq C_y, \quad \text{or} \quad (\text{C-ybnd.a})$$

$$\mathbb{E}[\psi_{\ell,i+2} - \psi_{\ell,i+1}|O_i] = 0, \quad (j = 1, \dots, m; i \in \mathbb{N}), \quad (\text{C-ybnd.b})$$

Then (C2'.e) holds with $b_{i+2}^y = \sum_{\ell=1}^n b_{\ell,i+2}^y$ where

$$b_{\ell,i+2}^y := 9C_y \mathbb{E}[\psi_{\ell,i+2} - \psi_{\ell,i+1}]. \quad (3.21)$$

For convenience, we also define the sum

$$d_{\ell,N}^y := \sum_{i=0}^{N-1} b_{\ell,i+2}^y = 9C_y \mathbb{E}[\psi_{\ell,N+1} - \psi_{\ell,0}]. \quad (3.22)$$

Proof. (C- ψ inc) implies $|\mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i]| = \mathbb{E}[\Psi_{i+2} - \Psi_{i+1}|O_i] \geq 0$, so (C2'.e) becomes

$$\mathbb{E}[(1 + \beta_i)\|y^{i+1} - y^i\|_{\mathbb{E}[\Psi_{i+2}-\Psi_{i+1}|O_i]}^2 + (1 + \beta_i^{-1})\|y^i - \hat{y}\|_{\mathbb{E}[\Psi_{i+2}-\Psi_i|O_{i-1}]}^2] \leq b_{i+2}^y.$$

In other words, for each block $\ell = 1, \dots, n$ should hold

$$\mathbb{E}[(1 + \beta_i)\|y_\ell^{i+1} - y_\ell^i\|_{\mathbb{E}[\psi_{\ell,i+2}-\psi_{\ell,i+1}|O_i]}^2 + (1 + \beta_i^{-1})\|y_\ell^i - \hat{y}_\ell\|_{\mathbb{E}[\psi_{\ell,i+2}-\psi_{\ell,i}|O_{i-1}]}^2] \leq b_{\ell,i+2}^y.$$

Taking $\beta_i = 1/2$ and estimating (3.23) with (C-ybnd) gives (3.21). \square

3.6. Satisfaction of the positivity condition (C1')

This requires $(1 - \delta)\Psi_{i+1} \geq \Lambda_i \Phi_i^{-1} \Lambda_i^*$, which can be expanded as

$$(1 - \delta) \sum_{\ell=1}^n \psi_{\ell,i+1} Q_\ell \geq \sum_{j=1}^m \sum_{\ell,k=1}^n \lambda_{\ell,j,i} \lambda_{k,j,i} \phi_{j,i}^{-1} Q_\ell K P_j K^* Q_k. \quad (3.23)$$

To go further from here, we require the functions κ_ℓ introduced next. After a general lemma that follows from the properties of the κ_ℓ , we look at specific constructions.

Definition 3.1. Writing $\mathcal{P} := \{P_1, \dots, P_m\}$, and $\mathcal{Q} := \{Q_1, \dots, Q_n\}$, we denote $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, \mathcal{Q})$ if each $\kappa_\ell : [0, \infty)^m \rightarrow [0, \infty)$, ($\ell = 1, \dots, n$), is monotone and we have

(i) (Estimation) The estimate

$$\sum_{j=1}^m \sum_{\ell, k=1}^n z_{\ell, j}^{1/2} z_{k, j}^{1/2} Q_\ell K P_j K^* Q_k \leq \sum_{\ell=1}^n \kappa_\ell(z_{\ell, 1}, \dots, z_{\ell, m}) Q_\ell. \quad (\text{C-}\kappa.\text{a})$$

(ii) (Boundedness) For some $\bar{\kappa} > 0$ the bound

$$\kappa_\ell(z_1, \dots, z_m) \leq \bar{\kappa} \sum_{j=1}^m z_j. \quad (\text{C-}\kappa.\text{b})$$

(iii) (Non-degeneracy) There exists $\underline{\kappa} > 0$ and $\ell^*(j) \in \{1, \dots, n\}$ with

$$\underline{\kappa} z_{j^*} \leq \kappa_{\ell^*(j)}(z_1, \dots, z_m), \quad (j \in \{1, \dots, m\}). \quad (\text{C-}\kappa.\text{c})$$

Lemma 3.7. Let $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, \mathcal{Q})$. The condition (C1') then holds if

$$(1 - \delta)\psi_{\ell, i+1} \geq \kappa_\ell(\dots, \lambda_{\ell, j, i}^2 \phi_{j, i}^{-1}, \dots), \quad (\ell = 1, \dots, n). \quad (\text{C-}\kappa\psi)$$

Proof. Clearly Φ_{i+1} is self-adjoint and positive definite. The remaining condition in (C1') is equivalent to (3.23), which follows from (C- κ .a) with $z_{\ell, j} := \lambda_{\ell, j, i}^2 \phi_{j, i}^{-1}$. \square

Example 3.3 (Simple structural κ). Using Cauchy's inequality, we deduce

$$\sum_{\ell, k=1}^n z_{\ell, j}^{1/2} z_{k, j}^{1/2} Q_\ell K P_j K^* Q_k \leq \sum_{\ell=1}^n z_{\ell, j} a_{\ell, j} Q_\ell, \quad (j = 1, \dots, m), \quad (3.24)$$

for $a_{\ell, j} := \|Q_\ell K P_j\|^2 \cdot \#\mathcal{V}(j)$. Thus (C- κ .a) and (C- κ .b) hold with $\bar{\kappa} = \max_{\ell, j} a_{\ell, j}$ if we take

$$\kappa_\ell(z_1, \dots, z_m) := \sum_{j=1}^m z_j a_{\ell, j}.$$

Clearly κ_ℓ is also monotone. If $\min_j \#\mathcal{V}(j) > 0$, then also (C- κ .c) is satisfied with $\underline{\kappa} = \min_j \max_{\ell \in \mathcal{V}(j)} a_{\ell, j} > 0$ and $\ell^*(j) := \arg \min_{\ell \in \mathcal{V}(j)} a_{\ell, j}$.

Example 3.4 (Worst-case κ). If $\#\mathcal{V}(j)$ is generally large, the previous example may provide very poor estimates. In this case, we may alternatively proceed with $\bar{z}_\ell := \max_j z_{j, \ell}$ as follows:

$$\sum_{j=1}^m \sum_{\ell, k=1}^n z_{\ell, j}^{1/2} z_{k, j}^{1/2} Q_\ell K P_j K^* Q_k \leq \sum_{\ell, k=1}^n \bar{z}_\ell^{1/2} \bar{z}_k^{1/2} Q_\ell K K^* Q_k \leq \sum_{\ell=1}^n \bar{z}_\ell \|K\|^2 Q_\ell.$$

Therefore (C- κ .a) and (C- κ .b) hold with $\bar{\kappa} = \|K\|^2$ for the monotone choice

$$\kappa_\ell(z_1, \dots, z_m) := \|K\|^2 \max\{z_1, \dots, z_m\}.$$

Clearly also $\underline{\kappa} = \bar{\kappa}$ for any choice of $\ell^*(j) \in \{1, \dots, n\}$.

Example 3.5 (Balanced κ). One more option is to choose the minimal κ_ℓ satisfying (C- κ .a) and the balancing condition

$$\kappa_\ell(z_{\ell, 1}, \dots, z_{\ell, m}) = \kappa_k(z_{k, 1}, \dots, z_{k, m}), \quad (\ell, k = 1, \dots, n).$$

This involves more refined use of Cauchy's inequality than the rough estimate (3.24), but tends to perform very well, as we will see in Section 5. This rule uses the data $\{z_{\ell, m}\}$ non-linearly.

3.7. Summary so far

We now summarise our findings so far, starting with writing out the proximal point iteration (PP) explicitly in terms of blocks. We already reformulated it in (3.2). We continue from there, first writing the $\lambda_{\ell,j,i}$ from (R- λ) in operator form as

$$\Lambda_i = K\mathring{T}_i^* \Phi_i^* - \Psi_{i+1} \mathring{\Sigma}_{i+1} K,$$

where $\mathring{T}_i := \sum_{j=1}^m \chi_{\mathring{S}(i)}(j) \hat{\tau}_{j,i} P_j$, and $\mathring{\Psi}_{i+1} := \sum_{j=1}^\ell \chi_{\mathring{V}(i+1)}(\ell) \hat{\sigma}_{\ell,i} Q_\ell$. Also defining $T_i^\perp := T_i - \mathring{T}_i$, and $\Sigma_{i+1}^\perp := \Sigma_{i+1} - \mathring{\Sigma}_{i+1}$, we can therefore rewrite (3.2) non-sequentially as

$$v^{i+1} := \Phi_i^{-1} K^* \mathring{\Sigma}_{i+1}^* \Psi_{i+1}^* (y^{i+1} - y^i) + T_i^\perp K^* y^{i+1}, \quad (3.25a)$$

$$x^{i+1} := (I + T_i \partial G)^{-1} (x^i - \mathring{T}_i K^* y^i - v^{i+1}), \quad (3.25b)$$

$$z^{i+1} := \Psi_{i+1}^{-1} K \mathring{T}_i^* \Phi_i^* (x^{i+1} - x^i) + \Sigma_{i+1}^\perp K x^{i+1}, \quad (3.25c)$$

$$y^{i+1} := (I + \Sigma_{i+1} \partial F^*)^{-1} (y^i + \mathring{\Sigma}_{i+1} K x^i + z^{i+1}). \quad (3.25d)$$

Let us set

$$\Theta_i := \sum_{j \in S(i)} \sum_{\ell \in \mathcal{V}(j)} \theta_{\ell,j,i} Q_\ell K P_j \quad \text{with} \quad \theta_{\ell,j,i+1} := \frac{\tau_{j,i} \phi_{j,i}}{\sigma_{\ell,i+1} \psi_{\ell,i+1}}.$$

Then thanks to (C-nest), we have $\Sigma_{i+1}^\perp \Theta_{i+1} = \Psi_{i+1}^{-1} K \mathring{T}_i^* \Phi_i^*$. Likewise,

$$B_i := \sum_{\ell \in V(i+1)} \sum_{j \in \mathcal{V}^{-1}(\ell)} b_{\ell,j,i} Q_\ell K P_j \quad \text{with} \quad b_{\ell,j,i+1} := \frac{\sigma_{\ell,i+1} \psi_{\ell,i+1}}{\tau_{j,i} \phi_{j,i}},$$

satisfies $T_i^\perp B_{i+1}^* = \Phi_i^{-1} K^* \mathring{\Sigma}_{i+1} \Psi_{i+1}$. Now we can rewrite (3.25) as

$$v^{i+1} := T_i^\perp [B_{i+1}^* (y^{i+1} - y^i) + K^* y^{i+1}], \quad (3.26a)$$

$$x^{i+1} := (I + T_i \partial G)^{-1} (x^i - \mathring{T}_i K^* y^i - v^{i+1}), \quad (3.26b)$$

$$z^{i+1} := \Sigma_{i+1}^\perp [\Theta_{i+1} (x^{i+1} - x^i) + K x^{i+1}], \quad (3.26c)$$

$$y^{i+1} := (I + \Sigma_{i+1} \partial F^*)^{-1} (y^i + \mathring{\Sigma}_{i+1} K x^i + z^{i+1}). \quad (3.26d)$$

Observe how (3.26b) can thanks to (S-G) be split into separate steps with respect to \mathring{T}_i and T_i^\perp , while (C-nest.a) guarantees $z^{i+1} = \Sigma_{i+1}^\perp [\Theta_{i+1} (\mathring{x}^{i+1} - x^i) + K \mathring{x}^{i+1}]$. Therefore, we obtain

$$\mathring{x}^{i+1} := (I + \mathring{T}_i \partial G)^{-1} (x^i - \mathring{T}_i K^* y^i), \quad (3.27a)$$

$$w^{i+1} := \Theta_{i+1} (\mathring{x}^{i+1} - x^i) + \mathring{x}^{i+1}, \quad (3.27b)$$

$$y^{i+1} := (I + \Sigma_{i+1} \partial F^*)^{-1} (y^i + \mathring{\Sigma}_{i+1} K x^i + \Sigma_{i+1}^\perp w^{i+1}), \quad (3.27c)$$

$$v^{i+1} := B_{i+1}^* (y^{i+1} - y^i) + y^{i+1}, \quad (3.27d)$$

$$x^{i+1} := (I + T_i^\perp \partial G)^{-1} (\mathring{x}^{i+1} - T_i^\perp v^{i+1}). \quad (3.27e)$$

In the blockwise case under consideration, in particular the setup of Lemma 3.1 together with (S-G) and (S-F*), the iterations (3.27) easily reduce to Algorithm 1. There we write

$$x_j := P_j x, \quad y_\ell := Q_\ell y, \quad \text{and} \quad K_{\ell,j} := Q_\ell K P_j.$$

In particular (C-step.a) can be written

$$\phi_{j,i+1} \mathbb{E}[\tau_{j,i+1} | j \in \mathring{S}(i+1)] \mathring{\pi}_{j,i+1} = \eta_{i+1} - \phi_{j,i} \hat{\tau}_{j,i} (1 - \chi_{\mathring{S}(i)}(j)).$$

Therefore, if $j \in \dot{S}(i)$, shifting indices down by one, we obtain the formula

$$\tau_{j,i} = \frac{\eta_i - \phi_{j,i-1}\tau_{j,i-1}\chi_{S(i-1)\setminus\dot{S}(i-1)}(j)}{\phi_{j,i}\pi_{j,i}}.$$

Similarly we obtain the other step length formulas in Algorithm 1.

Algorithm 1 Block-stochastic primal-dual method: general form

Require: Convex, proper, lower semi-continuous functions $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ with the separable structures (S-G) and (S-F*). Rules for $\phi_{j,i}$, $\psi_{\ell,i+1}$, η_{i+1} , $\eta_{\tau,i+1}^\perp$, $\eta_{\sigma,i+1}^\perp \in \mathcal{R}(O_i; [0, \infty))$, as well as sampling rules for $\dot{S}(i)$ and $\dot{V}(i+1)$, ($j = 1, \dots, m$; $\ell = 1, \dots, n$; $i \in \mathbb{N}$).

1: Choose initial iterates $x^0 \in X$ and $y^0 \in Y$.

2: **for all** $i \geq 0$ **until** a stopping criterion is satisfied **do**

3: Sample $\dot{S}(i) \subset S(i) \subset \{1, \dots, m\}$ and $\dot{V}(i+1) \subset V(i+1) \subset \{1, \dots, n\}$ subject to (C-nest).

4: For each $j \notin \dot{S}(i)$, set $x_j^{i+1} := x_j^i$.

5: For each $j \in \dot{S}(i)$, compute

$$\tau_{j,i} := \frac{\eta_i - \phi_{j,i-1}\tau_{j,i-1}\chi_{S(i-1)\setminus\dot{S}(i-1)}(j)}{\phi_{j,i}\pi_{j,i}}, \quad \text{and}$$

$$x_j^{i+1} := (I + \tau_{j,i}\partial G_j)^{-1} \left(x_j^i - \tau_{j,i} \sum_{\ell \in \mathcal{V}(j)} K_{\ell,j}^* y_\ell^i \right).$$

6: For each $j \in \dot{S}(i)$ and $\ell \in \mathcal{V}(j)$, set

$$\tilde{w}_{\ell,j}^{i+1} := \theta_{\ell,j,i+1}(x_j^{i+1} - x_j^i) + x_j^{i+1} \quad \text{with} \quad \theta_{\ell,j,i+1} := \frac{\tau_{j,i}\phi_{j,i}}{\sigma_{\ell,i+1}\psi_{\ell,i+1}}.$$

7: For each $\ell \notin V(i+1)$, set $y_\ell^{i+1} := y_\ell^i$.

8: For each $\ell \in \dot{V}(i+1)$, compute

$$\sigma_{j,i+1} := \frac{\eta_i - \psi_{j,i}\sigma_{j,i}\chi_{V(i)\setminus\dot{V}(i)}(j)}{\psi_{j,i+1}\dot{v}_{\ell,i+1}}, \quad \text{and}$$

$$y_\ell^{i+1} := (I + \sigma_{\ell,i+1}\partial F_\ell^*)^{-1} \left(y_\ell^i + \sigma_{\ell,i+1} \sum_{j \in \mathcal{V}^{-1}(\ell)} K_{\ell,j} x_j^i \right).$$

9: For each $\ell \in V(i+1) \setminus \dot{V}(i+1)$ compute

$$\sigma_{j,i+1} := \frac{\eta_{\sigma,i}^\perp}{\psi_{j,i+1}(v_{\ell,i+1} - \dot{v}_{\ell,i+1})}, \quad \text{and}$$

$$y_\ell^{i+1} := (I + \sigma_{\ell,i+1}\partial F_\ell^*)^{-1} \left(y_\ell^i + \sigma_{\ell,i+1} \sum_{j \in \mathcal{V}^{-1}(\ell)} K_{\ell,j} \tilde{w}_{\ell,j}^{i+1} \right).$$

10: For each $\ell \in \dot{V}(i+1)$ and $j \in \mathcal{V}^{-1}(\ell)$, set

$$\tilde{v}_{\ell,j}^{i+1} := b_{\ell,j,i+1}(y_\ell^{i+1} - y_\ell^i) + y_\ell^i \quad \text{with} \quad b_{\ell,j,i+1} := \frac{\sigma_{\ell,i+1}\psi_{\ell,i+1}}{\tau_{j,i}\phi_{j,i}}.$$

11: For each $j \in S(i) \setminus \dot{S}(i)$, compute

$$\tau_{j,i} := \frac{\eta_{\tau,i}^\perp}{\phi_{j,i}(\pi_{j,i} - \pi_{j,i})}, \quad \text{and}$$

$$x_j^{i+1} := (I + \tau_{j,i}\partial G_j)^{-1} \left(x_j^i - \tau_{j,i} \sum_{\ell \in \mathcal{V}(j)} K_{\ell,j}^* \tilde{v}_{\ell,j}^{i+1} \right).$$

12: **end for**

One way to simplify the algorithm in concept, is to alternate between the two x - y and y - x update directions through the *random* alternating choice of $\dot{S}(i) = \emptyset$ and $\dot{V}(i+1) = \emptyset$. We will discuss this in more detail in Section 4.10. Alternatively, if $\dot{S}(i) = S(i)$, the last two steps of (3.27) vanish, giving $x^{i+1} := \dot{x}^{i+1}$. If this choice is

deterministic, then also $\mathring{V}(i+1) = \emptyset$, so we are in the full dual updates setting of Lemma 3.2. The result is Algorithm 2.

Algorithm 2 Block-stochastic primal-dual method: full dual updates

Require: Convex, proper, lower semi-continuous functions $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ with the separable structures (S-G) and (S-F*). Rules for $\phi_{j,i}, \psi_{\ell,i+1}, \eta_{i+1} \in \mathcal{R}(O_i; (0, \infty))$, as well as a sampling rule for the set $S(i)$, ($j = 1, \dots, m; \ell = 1, \dots, n; i \in \mathbb{N}$).

- 1: Choose initial iterates $x^0 \in X$ and $y^0 \in Y$.
- 2: **for all** $i \geq 0$ **until** a stopping criterion is satisfied **do**
- 3: Select random $S(i) \subset \{1, \dots, m\}$.
- 4: For each $j \notin S(i)$, set $x_j^{i+1} := x_j^i$.
- 5: For each $j \in S(i)$, with $\tau_{j,i} := \eta_i \pi_{j,i}^{-1} \phi_{j,i}^{-1}$, compute

$$x_j^{i+1} := (I + \tau_{j,i} \partial G_j)^{-1} \left(x_j^i - \tau_{j,i} \sum_{\ell \in \mathcal{V}(j)} K_{\ell,j}^* y_\ell^i \right).$$

- 6: For each $j \in S(i)$ set

$$\bar{x}_j^{i+1} := \theta_{j,i+1} (x_j^{i+1} - x_j^i) + x_j^{i+1} \quad \text{with} \quad \theta_{j,i+1} := \frac{\eta_i}{\pi_{j,i} \eta_{i+1}}.$$

- 7: For each $\ell \in \{1, \dots, n\}$ using $\sigma_{\ell,i+1} := \eta_{i+1} \psi_{\ell,i+1}^{-1}$, compute

$$y_\ell^{i+1} := (I + \sigma_{\ell,i+1} \partial F_\ell^*)^{-1} \left(y_\ell^i + \sigma_{\ell,i+1} \sum_{j \in \mathcal{V}^{-1}(\ell)} K_{\ell,j} \bar{x}_j^{i+1} \right).$$

- 8: **end for**
-

We have not yet specified η_i and $\psi_{\ell,i}$. We have also not fixed $\phi_{j,i}$, giving the options in Example 3.1 and Example 3.2. We will return to these choices in the next section, but now summarise our results so far by specialising Corollary 2.2.

Proposition 3.1. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, \mathcal{Q})$. Then the conditions (C0), (C1'), and (C2') hold when we do the following for each $i \in \mathbb{N}$.*

- (i) Sample $\mathring{S}(i) \subset S(i) \subset \{1, \dots, m\}$ and $\mathring{V}(i+1) \subset V(i+1) \subset \{1, \dots, n\}$ subject to (C-nest).
- (ii) Define Φ_{i+1} through (S- $\Phi\Psi$), satisfying (C-xbnd) for $\tilde{\gamma}_j \geq 0$ to be specified.
- (iii) Define Ψ_{i+1} through (S- $\Phi\Psi$), satisfying (C- ψ inc), (C-ybnd), and (C- $\kappa\psi$).
- (iv) Take T_i and Σ_i of the form (S- $T\Sigma$) with the blockwise step lengths satisfying (C-step).
- (v) Define Λ_i through (S- Λ) and (R- λ).
- (vi) Either (Lemma 3.1 or Lemma 3.2)
 - (a) Take $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$ satisfying (C- η^\perp) and (C-step). In this case (C \mathcal{G}) holds; or
 - (b) Satisfy (C-step'), forcing $S(i) = \mathring{S}(i)$, $\mathring{V}(i+1) = \emptyset$, $V(i+1) = \{1, \dots, n\}$. In this case (C \mathcal{G}_*) holds, as do (C-nest) and (C-step) with $\eta_{\tau,i}^\perp \equiv 0$, and $\eta_{\sigma,i+1}^\perp = \eta_i$.

Let then G and F^* have the separable structures (S-G) and (S-F*). For each $j = 1, \dots, m$, suppose G_j is (strongly) convex with corresponding factor $\gamma_j \geq 0$, and pick $\tilde{\gamma}_j \in [0, \gamma_j]$. Then there exists $C_0 > 0$ such that the iterates of (PP) satisfy

$$\delta \sum_{k=1}^m \frac{1}{\mathbb{E}[\phi_{k,N}^{-1}]} \cdot \mathbb{E}[\|x_k^N - \hat{x}_k\|^2] + \tilde{g}_N \leq C_0 + \sum_{j=1}^m d_{j,N}^x(\tilde{\gamma}_j) + \sum_{\ell=1}^n d_{\ell,N}^y, \quad (3.30)$$

where $d_{j,N}^x(\tilde{y}_j)$ is defined in (3.11), $d_{\ell,N}^y$ is defined in (3.22), and we set

$$\tilde{g}_N := \begin{cases} \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N), & (\text{C}\mathcal{G}) \text{ holds and } \tilde{y}_j \leq \gamma_j/2 \text{ for all } j, \\ \zeta_{*,N} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}), & (\text{C}\mathcal{G}_*) \text{ holds and } \tilde{y}_j \leq \gamma_j/2 \text{ for all } j, \\ 0, & \text{otherwise.} \end{cases}$$

Here ζ_N and the ergodic variables \tilde{x}_N and \tilde{y}_N in (2.15), and the gap \mathcal{G} in (2.16). The alternatives $\zeta_{*,N}$, $\tilde{x}_{*,N}$ and $\tilde{y}_{*,N}$ are defined in (2.19).

Proof. We have proved all of the conditions (C0), (C1'), (C2') and (C \mathcal{G}), alternatively (C \mathcal{G}_*), in Lemmas 3.1 to 3.3, 3.6 and 3.7. Only the estimate (3.30) demands further verification.

In Corollary 2.2, we have assumed that either $\tilde{\Gamma} = \Gamma$ or $\tilde{\Gamma} = \Gamma/2$, that is $\tilde{y}_j \in \{\gamma_j, \gamma_j/2\}$. However, G_j is (strongly) convex with factor γ'_j for any $\gamma'_j \in [0, \gamma_j]$, so we may relax this assumption to $0 \leq \tilde{y}_j \leq \gamma_j$ with the gap estimates holding when $\tilde{y}_j \leq \gamma_j/2$.

Setting $C_0 := \frac{1}{2} \|u^0 - \hat{u}\|_{Z_0 L_0}^2$, Corollary 2.2 thus shows

$$\delta \mathbb{E}[\|x^N - \hat{x}\|_{\Phi_N}^2] + \tilde{g}_N \leq C_0 + \sum_{i=0}^{N-1} (b_{i+2}^x(\tilde{y}_j) + b_{i+2}^y).$$

By Hölder's inequality

$$\mathbb{E}[\|x^N - \hat{x}\|_{\Phi_N}^2] = \sum_{k=1}^m \mathbb{E}[\phi_{k,N} \|x_k^N - \hat{x}_k\|^2] \geq \sum_{k=1}^m \mathbb{E}[\|x_k^N - \hat{x}_k\|]^2 / \mathbb{E}[\phi_{k,N}^{-1}].$$

The estimate (3.30) is now immediate. \square

Specialised to Algorithms 1 and 2, we obtain the following corollaries.

Corollary 3.1. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, Q)$. Suppose the primal bound (C-xbnd), and the dual test conditions (C- ψ inc), (C-ybnd), and (C- $\kappa\psi$) hold along with (C- η^\perp), (C- η). Then the iterates of Algorithm 1 satisfy (3.30) with $\tilde{g}_N = \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N)$ when $\tilde{y}_j \leq \gamma_j/2$ for all j , and $\tilde{g}_N = 0$ otherwise.*

Proof. Algorithm 1 satisfies the structural assumptions (S- $\Phi\Psi$), (S- $T\Sigma$), and (S- Λ), the conditions (C-nest) and (R- λ), as well as the alternative condition (a) of Proposition 3.1, provided the non-negativity conditions in (C-step) are satisfied. They are indeed ensured by us assuming (C- η). The remaining conditions of Proposition 3.1 we have also assumed. \square

Corollary 3.2. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, Q)$. Suppose (C-xbnd) and the dual conditions (C- ψ inc), (C-ybnd), and (C- $\kappa\psi$) hold. Then the iterates of Algorithm 2 satisfy (3.30) with $\tilde{g}_N = \zeta_{*,N} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N})$ when $\tilde{y}_j \leq \gamma_j/2$ for all j , and $\tilde{g}_N = 0$ otherwise.*

Proof. Algorithm 2 satisfies the structural assumptions (S- $\Phi\Psi$), (S- $T\Sigma$), (S- Λ) and (R- λ), as well as (b) of Proposition 3.1. Its remaining conditions we have assumed. \square

4. Dual tests and penalty bounds for block-proximal methods

We now need to satisfy the conditions of Corollaries 3.1 and 3.2. This involves choosing update rules for η_{i+1} , $\eta_{\tau,i+1}^\perp$, $\eta_{\sigma,i+1}^\perp$, $\phi_{j,i+1}$ and $\psi_{\ell,i+1}$. Specifically, for both Corollaries, we have to verify the primal bound (C-xbnd). For Corollary 3.1, we moreover need (C- η^\perp), (C- η), and the non-negativity conditions in (C-step.a) and (C-step.b). At the

same time, to obtain good convergence rates, we need to make $d_{j,N}^x(\tilde{y}_j)$ and $d_{\ell,N}^y = \mathbb{E}[\psi_{\ell,N+1} - \psi_{\ell,0}]$ small in (3.30).

We concentrate on the deterministic primal test update rule of Example 3.2, which also provide estimates on $d_{j,N}^x(\tilde{y}_j)$ if the conditions of 3.5 are satisfied. In addition to the verifications above, we need to verify (C- ϕ det). To satisfy (C-xbnd), we have to take $\rho_j = 0$ unless the bound (C-xbnd.a) holds for the specific problem under consideration.

We begin with general assumptions, after which in Section 4.2 we calculate expectation bounds on $\phi_{j,i}$. In Sections 4.3 to 4.8 we give useful choice of η_i and $\psi_{\ell,i}$ that finally yield specific convergence results. We finish the section with choices for $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$ in Section 4.9, and sampling patterns in Section 4.10. The difficulty of (fully) extending our estimates to the random primal test update of Example 3.1, we discuss in Remark 4.2.

4.1. Assumptions and simplifications

Throughout this section, we assume for simplicity that the probabilities stay constant between iterations,

$$\pi_{j,i} \equiv \pi_j > 0, \quad \text{and} \quad v_{\ell,i} \equiv v_\ell. \quad (\text{R-}\pi v)$$

Then (C-nest) shows that

$$\pi_{j,i} \equiv \pi_j > 0, \quad \text{and} \quad v_{\ell,i} \equiv v_\ell > 0.$$

We assume (C- η) to hold. As we recall from Lemma 3.2, this is the case for Algorithm 2 if

$$i \mapsto \eta_i > 0 \quad \text{is non-decreasing.}$$

Finally, aside from the non-negativity conditions that will be verified through the choice of $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$, we note that (C-step) holds in Algorithms 1 and 2. It is therefore assumed.

4.2. Estimates for deterministic primal test updates

We consider the deterministic primal test updates of Example 3.2. To start with, from (R- ϕ det), we compute

$$\phi_{j,N} = \phi_{j,N-1} + 2(\tilde{y}_j \eta_{N-1} + \rho_j) = \phi_{j,0} + 2\rho_j N + 2\tilde{y}_j \sum_{i=0}^{N-1} \eta_i. \quad (4.1)$$

The following lemma lists the fundamental properties that this update rule satisfies.

Lemma 4.1. *Suppose (C- η), (C- ϕ det), and (R- πv) hold. If (C-xbnd.a) holds with the constant $C_x \geq 0$, take $\rho_j \geq 0$, otherwise take $\rho_j = 0$, supposing $\tilde{y}_j + \rho_j > 0$, ($j = 1, \dots, m$). Define $\phi_{j,i+1}$ according to Example 3.2. Suppose $\eta_i \geq b_j(i+1)^p$, for some $p, b_j > 0$. Then for some $c_j > 0$, and $C_\alpha > 0$ holds*

$$\phi_{j,N} \in \mathcal{R}(O_{N-1}; (0, \infty)), \quad (\text{C-}\phi\text{bnd.a})$$

$$\mathbb{E}[\phi_{j,N}] = \phi_{j,0} + 2\rho_j N + 2\tilde{y}_j \sum_{i=0}^{N-1} \mathbb{E}[\eta_i], \quad \text{and} \quad (\text{C-}\phi\text{bnd.b})$$

$$\mathbb{E}[\phi_{j,N}^{-1}] \leq c_j N^{-1}, \quad (N \geq 1). \quad (\text{C-}\phi\text{bnd.c})$$

Moreover, the primal test bound (C-xbnd) holds, and with $C_\alpha = 18$ we have

$$d_{j,N}^x(\tilde{y}_j) = \rho_j C_x C_\alpha N. \quad (\text{C-}\phi\text{bnd.d})$$

Suppose moreover that $\eta_i \geq b_j \min_j \phi_{j,i}^p$, for some $p, b_j > 0$. Then for some $\tilde{c}_j \geq 0$ holds

$$\frac{1}{\mathbb{E}[\phi_{j,N}^{-1}]} \geq \tilde{c}_j N^{p+1}, \quad (N \geq 4). \quad (\text{C-}\phi\text{bnd.e})$$

Proof. The conditions of Lemma 3.5 are guaranteed by our assumptions. It directly proves (C- $\phi\text{bnd.a}$) and (C- $\phi\text{bnd.d}$). Since we assume $i \mapsto \eta_i$ to be increasing, clearly $\phi_{j,N} \geq 2N\tilde{\rho}_j$ for $\tilde{\rho}_j := \rho_j + \tilde{\gamma}_j\eta_0 > 0$. Then $\phi_{j,N}^{-1} \leq \frac{1}{2\tilde{\rho}_j N}$. Taking the expectation proves (C- $\phi\text{bnd.c}$), while (C- $\phi\text{bnd.b}$) is immediate from (4.1). Clearly (C- $\phi\text{bnd.e}$) holds if $\tilde{\gamma}_j = 0$, so assume $\tilde{\gamma}_j > 0$. Under our assumption on η_i , Lemma B.1 shows for some $B_j > 0$ that $\phi_{j,N}^{-1} \leq \frac{1}{B_j N^{1+p}}$. Taking the expectation proves (C- $\phi\text{bnd.e}$) for a $\tilde{c}_j := B_j/\tilde{\gamma}_j$. \square

Remark 4.1. From (4.1), we see that $\eta_i \geq b_j(i+1)^p$ if $\eta_i \geq \tilde{b}_j \min_j \phi_{j,i}^p$, for some $\tilde{b}_j > 0$.

Remark 4.2. Propositions 4.1 and 4.2 to follow, will generalise to any update rule for $\phi_{j,i+1}$ that satisfies (C- ϕbnd). The conditions (C- $\phi\text{bnd.a}$)–(C- $\phi\text{bnd.d}$) can easily be shown for the random primal test updates of Example 3.1. The estimate (C- $\phi\text{bnd.e}$) however is challenging, with any derivation likely dependent on the exact sampling patterns employed. The estimate (C- $\phi\text{bnd.e}$) is, however, only needed to estimate $1/\mathbb{E}[\phi_{j,N}^{-1}]$ from below in the first term of the general estimate (3.30), and therefore only affects convergence of the iterates, not the gap. Hence the estimates in the upcoming Propositions 4.1 and 4.2 on the *ergodic duality gap*, but not the iterates, do hold for the random primal test update rule of Example 3.1.

4.3. Dual bounds—a first attempt

We now need to satisfy the conditions (C- ψinc), (C- ybnd), and (C- $\kappa\psi$) on the dual updates. By the construction of $\lambda_{\ell,j,i}$ in (R- λ), the step length condition (C- step), and the constant probability assumption (R- $\pi\nu$), we have

$$\lambda_{\ell,j,i} \leq \eta_i (\pi_j^{-1} \chi_{\hat{S}(i)}(j) + \nu_\ell^{-1} \chi_{\hat{V}(i+1)}(\ell)) =: \eta_i \hat{\mu}_{\ell,j,i}, \quad (\ell \in \mathcal{V}(j)).$$

Therefore (C- $\kappa\psi$) holds if we take

$$\psi_{\ell,i+1} := \frac{\eta_i^2}{1-\delta} \kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots). \quad (4.2)$$

With this, the condition (C- ψinc) would require for all $\ell = 1, \dots, n$ that

$$\mathbb{E}[\eta_{i+1}^2 \kappa_\ell(\dots, \hat{\mu}_{\ell,j,i+1}^2 \phi_{j,i+1}^{-1}, \dots) | \mathcal{O}_i] \geq \mathbb{E}[\eta_i^2 \kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_i] \quad (4.3)$$

Since $\eta_{i+1} \in \mathcal{R}(\mathcal{O}_i; (0, \infty))$, and $\eta_i \in \mathcal{R}(\mathcal{O}_i; (0, \infty))$, we can take η_i and η_{i+1} outside the expectations in (4.3). A first idea would then be to take η_{i+1} as the smallest number satisfying (4.3) for all ℓ . In the deterministic case, the resulting rule will telescope, and reduce to the one that we will follow. In the stochastic case, we have however observed numerical instability, and have also been unable to prove convergence. Therefore, we have to study “less optimal” rules.

4.4. Worst-case conditions

For a random variable $p \in \mathcal{R}(\Omega; \mathbb{R})$ on the probability space $(\Omega, \mathcal{O}, \mathbb{P})$, let us define the conditional worst-case realisation with respect to the σ -algebra $\mathcal{O}' \subset \mathcal{O}$ as the random variable $\mathbb{W}[p | \mathcal{O}'] \in \mathcal{R}(\mathcal{O}'; \mathbb{R})$ defined by

$$p \leq \mathbb{W}[p | \mathcal{O}'] \leq q \quad \mathbb{P}\text{-a.e.} \quad \text{for all } q \in \mathcal{R}(\mathcal{O}'; \mathbb{R}) \quad \text{s.t. } p \leq q \quad \mathbb{P}\text{-a.e.}$$

We also write $\mathbb{W}[p] := \mathbb{W}[p|\mathcal{O}']$ when $\mathcal{O}' = \{\Omega, \emptyset\}$ is the trivial σ -algebra.

Following the derivation of (4.2), the condition (C- $\kappa\psi$) will hold if

$$\psi_{\ell,i+1} \geq \frac{\eta_i^2}{1-\delta} \mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_{i-1}]. \quad (4.4)$$

Accordingly, we take

$$\eta_i := \min_{\ell=1,\dots,n} \sqrt{\frac{(1-\delta)\psi_{\ell,i+1}}{\mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_{i-1}]}}. \quad (\text{R-}\eta)$$

By the construction of \mathbb{W} , we get $\eta_i \in \mathcal{R}(\mathcal{O}_{i-1}; (0, \infty))$ provided also $\psi_{i+1} \in \mathcal{R}(\mathcal{O}_{i-1}; (0, \infty))$.

It is our task in the rest of this section to experiment with different choices of $\psi_{\ell,i+1}$, satisfying (C- ψ inc) and (C- ψ bnd). Before this we establish the following important fact.

Lemma 4.2. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, \mathcal{Q})$. Suppose (R- πv) holds, and that both $i \mapsto \phi_{j,i}$ and $i \mapsto \psi_{\ell,i}$ are non-decreasing for all $j = 1, \dots, m$ and $\ell = 1, \dots, n$. Then $i \mapsto \eta_i$ defined in (R- η) is non-decreasing.*

Proof. We fix $\ell \in \{1, \dots, n\}$. The condition (R- πv) implies that $(\hat{\mu}_{\ell,1,i}, \dots, \hat{\mu}_{\ell,m,i})$ are independently identically distributed for all $i \in \mathbb{N}$. Since $\phi_{j,i} \in \mathcal{R}(\mathcal{O}_{i-1}; (0, \infty))$, we can for some random $(\hat{\mu}_1, \dots, \hat{\mu}_m)$ on a probability space $(\mathbb{P}_\mu, \Omega_\mu, \mathcal{O}_\mu)$, distinct from $(\mathbb{P}, \Omega, \mathcal{O})$, write

$$\mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_{i-1}] \sim \mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_j^2 \phi_{j,i}^{-1}, \dots)], \quad (i \in \mathbb{N}),$$

where \sim stands for “identically distributed”. Since $i \mapsto \phi_{j,i}$ is non-decreasing and κ_ℓ monotone, this implies

$$\mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_{i-1}] \geq \mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i+1}^2 \phi_{j,i+1}^{-1}, \dots) | \mathcal{O}_i], \quad \mathbb{P}\text{-a.e.}$$

Since $i \mapsto \psi_{\ell,i}$ is also non-decreasing, the claim follows. \square

4.5. Partial strong convexity: Bounded ψ

In addition to the assumptions in Section 4.1, from now on we assume (C- ϕ bnd) to hold. As we have seen, this is the case for the deterministic primal test update rule of both Example 3.2. For the random primal test update rule of (3.1), the rest of the conditions hold, but we have not been able to verify (C- ϕ bnd.e). This has the implication that only the gap estimates hold.

As a first option for $\psi_{\ell,i}$, let us take $\psi_{\ell,i} \equiv \psi_{\ell,0}$. Then both (C- ψ inc) and (C- ψ bnd.b) clearly hold, and $d_{\ell,N}^y \equiv 0$. Moreover Lemma 4.2 shows that $i \mapsto \eta_i$ is non-decreasing, as we have required in Section 4.1. To obtain convergence rates, we still need to estimate the primal penalty $d_{j,N}^x(\tilde{y}_j)$ as well as ζ_N and $\zeta_{*,N}$. Presently the dual penalty $d_{\ell,N}^y = \mathbb{E}[\psi_{\ell,N} - \psi_{\ell,0}] \equiv 0$.

With $\underline{\psi}_0 := \min_{\ell=1,\dots,n} \psi_{\ell,0}$, we compute

$$\eta_i \geq \sqrt{\frac{(1-\delta)\underline{\psi}_0}{\max_{\ell=1,\dots,n} \mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_{i-1}]}}.$$

Let us define

$$\mathbb{w}_j := \max_{\ell \in \mathcal{V}(j)} \mathbb{W}[\hat{\mu}_{\ell,j,i} | \mathcal{O}_{i-1}] = \max_{\ell \in \mathcal{V}(j)} \mathbb{W}[\hat{\pi}_j^{-1} \chi_{\hat{S}(i)}(j) + \hat{v}_\ell^{-1} \chi_{\hat{V}(i+1)}(\ell) \mid \mathcal{O}_{i-1}],$$

which is independent of $i \geq 0$; see the proof of Lemma 4.2. Since $\hat{\mu}_{\ell,j,i} = 0$ for $\ell \notin \mathcal{V}(j)$, using (C- κ .b) and $\phi_{j,i} \in \mathcal{R}(\mathcal{O}_{i-1}; (0, \infty))$, we further get

$$\eta_i \geq \sqrt{\frac{(1-\delta)\psi_0}{\bar{\kappa} \max_{\ell \in \mathcal{V}(j)} \mathbb{W}[\sum_{j=1}^n \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1} | \mathcal{O}_{i-1}]}} \geq \sqrt{\frac{1}{\sum_{j=1}^n b_j^{-1} \phi_{j,i}^{-1}}} \quad (4.5)$$

for $b_j := (1-\delta)\psi_0/(\bar{\kappa}\mathbb{W}_j^2)$. Jensen's inequality thus gives

$$\mathbb{E}[\eta_i] \geq \left(\mathbb{E} \left[\sqrt{\sum_{j=1}^n b_j^{-1} \phi_{j,i}^{-1}} \right] \right)^{-1} \geq \left(\sum_{j=1}^n b_j^{-1} \mathbb{E}[\phi_{j,i}^{-1}] \right)^{-1/2}. \quad (4.6)$$

Using (C- ϕ bnd.c), it follows

$$\mathbb{E}[\eta_i] \geq C_\eta i^{1/2} \quad \text{for } C_\eta := \left(\sum_{j=1}^m b_j^{-1} c_j \right)^{-1/2}. \quad (4.7)$$

We recall (C- η) guaranteeing $\bar{\eta}_i \geq \mathbb{E}[\eta_i]$. Consequently, we estimate ζ_N from (2.14) by

$$\begin{aligned} \zeta_N &= \sum_{i=0}^{N-1} \bar{\eta}_i \geq \sum_{i=0}^{N-1} \mathbb{E}[\eta_i] \geq C_\eta \sum_{i=0}^{N-1} i^{1/2} \geq C_\eta \int_0^{N-2} x^{1/2} dx \geq \frac{2C_\eta}{3} (N-2)^{3/2} \\ &= \frac{2C_\eta}{3} N^{3/2} \left(\frac{N-2}{N} \right)^{3/2} = \frac{C_\eta}{\sqrt{18}} N^{3/2}, \quad (N \geq 4). \end{aligned} \quad (4.8)$$

Similarly $\zeta_{*,N}$ defined in (2.18) satisfies

$$\zeta_{*,N} \geq \sum_{i=1}^{N-1} \mathbb{E}[\eta_i] \geq \frac{2C_\eta}{3} ((N-2)^{3/2} - 1) \geq \frac{(2^{3/2} - 1)C_\eta}{2^{3/2} \sqrt{18}} N^{3/2} \geq \frac{C_\eta}{2\sqrt{18}} N^{3/2}, \quad (N \geq 4). \quad (4.9)$$

For the deterministic primal test update rule of Example 3.2, by Lemma 3.5, we still need to satisfy (C- ϕ det), that is $2\tilde{\gamma}_j \bar{\gamma}_j \eta_i \leq \delta \phi_{j,i} (\tilde{\gamma}_j - \bar{\gamma}_j)$ for $j \in S(i)$ and $i \in \mathbb{N}$. The non-degeneracy assumption (C- κ .c) applied in (R- η) gives $\eta_i \leq \sqrt{(1-\delta)\psi_{\ell,0}\phi_{j,i}/(\bar{\kappa}\mathbb{W}_j)}$. Since under both rules Examples 3.1 and 3.2, $\phi_{j,i}$ is increasing in i , it therefore suffices to choose $\bar{\gamma}_j \geq 0$ to satisfy

$$\tilde{\gamma}_j = \bar{\gamma}_j = 0 \quad \text{or} \quad \frac{2\tilde{\gamma}_j \bar{\gamma}_j}{\tilde{\gamma}_j - \bar{\gamma}_j} \sqrt{\frac{1-\delta}{\bar{\kappa}\mathbb{W}_j}} \leq \delta \psi_{\ell,0}^{-1/2} \phi_{j,0}^{1/2}. \quad (\text{C-}\phi\text{det}')$$

These findings can be summarised as:

Proposition 4.1. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, \mathcal{Q})$. Pick $\rho_j \geq 0$ and $\tilde{\gamma}_j \in [0, \gamma_j]$, ($j = 1, \dots, m$). In Algorithm 1 or Algorithm 2, take*

- (i) *the probabilities $\pi_{j,i} \equiv \pi_j$ and (in Algorithm 1) $\mathring{v}_{\ell,i} \equiv \mathring{v}_\ell$ constant over iterations,*
- (ii) *η_i according to (R- η), and (in Algorithm 1) $\eta_{\tau,i}^\perp, \eta_{\sigma,i}^\perp > 0$ satisfying (C- η) and (C- η^\perp),*
- (iii) *$\phi_{j,0} > 0$ by free choice, and $\phi_{j,i}$ for $i \geq 1$ following Example 3.2, taking $0 \leq \bar{\gamma}_j < \tilde{\gamma}_j$ or $\bar{\gamma}_j = 0$, and satisfying (C- ϕ det'),*
- (iv) *$\psi_{\ell,i} := \psi_{\ell,0}$ for some fixed $\psi_{\ell,0} > 0$, ($\ell = 1, \dots, n$).*

Suppose for each $j = 1, \dots, m$ that $\rho_j + \tilde{\gamma}_j > 0$ and either $\rho_j = 0$ or (C-xbnd.a) holds with the constant C_x . Let \tilde{c}_k be the constant provided by Lemma 4.1. Then

$$\sum_{k=1}^m \delta \tilde{c}_k \tilde{\gamma}_k \mathbb{E} [\|x_k^N - \hat{x}_k\|^2] + \frac{C_\eta}{\sqrt{18}} g_N \leq \frac{C_0 + C_x C_\alpha (\sum_{j=1}^m \rho_j) N}{N^{3/2}}, \quad (N \geq 4), \quad (4.10)$$

where

$$g_N := \begin{cases} \mathcal{G}(\tilde{x}_N, \tilde{y}_N), & \text{Algorithm 1, } \tilde{\gamma}_j \leq \gamma_j/2 \text{ for all } j, \\ \frac{1}{2} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}), & \text{Algorithm 2, } \tilde{\gamma}_j \leq \gamma_j/2 \text{ for all } j, \\ 0, & \text{otherwise.} \end{cases}$$

Proof. By (4.5), we have $\eta_i \geq b_j^{1/2} \min_j \phi_{j,i}^{1/2}$ for each $j = 1, \dots, m$, so (C- ϕ bnd.e) holds with $p = 1/2$. Therefore, (4.10) is immediate from (C- ϕ bnd.e), Corollaries 3.1 and 3.2, and the estimates (4.8) and (4.9), whose assumptions follow from Lemmas 4.1 and 4.2 and (C- ϕ det'). \square

Remark 4.3 (Practical parameter initialisation). In practise, we take $\tau_{j,0}$, η_i , and $\delta \in (0, 1)$ as the free step length parameters. Then (C-step.a) and (C-step.b) give $\phi_{j,0} = \eta_0 / (\tau_{j,0} \tilde{\tau}_{j,0})$. As $\psi_{\ell,0}$ we take a value reaching the maximum in (R- η). In practise, we take $\eta_0 = 1 / \min_j \tau_{j,0}$. This choice appears to work well, and is consistent with the basic algorithm (2.2) corresponding to $\phi_{j,0} = \tau_{j,0}^{-2}$. For the deterministic update rule (R- ϕ det), we use (C- ϕ det') to bound $\tilde{\gamma}_j$.

4.6. Partial strong convexity: Increasing ψ

In (R- η), let us take $\psi_{\ell,i+1} := \psi_{\ell,0} \eta_i$. Then

$$\eta_i = \min_{\ell=1,\dots,n} \frac{(1-\delta)\psi_{\ell,0}}{\mathbb{W}[\kappa_\ell(\dots, \hat{\mu}_{\ell,j,i}^2 \phi_{j,i}^{-1}, \dots) | \mathcal{O}_{i-1}]}. \quad (\text{R-}\eta^2)$$

Adapting Lemma 4.2, we see that $i \mapsto \eta_i$ is non-decreasing. Thus $i \mapsto \psi_{\ell,i}$ is also increasing, so (C- ψ inc) holds. As (C-ybnd.b) does not hold, we need to assume (C-ybnd.a). To obtain convergence rates, we need to estimate both the primal and the dual penalties $d_{j,N}^x(\tilde{\gamma}_j)$ and $d_{\ell,N}^y$, as well as ζ_N and $\mathbb{E}[\psi_{k,N}]$.

Similarly to the derivation of (4.7), we deduce with the help of (C- ϕ bnd.c) that

$$\mathbb{E}[\eta_i] \geq C_\eta^2 i. \quad (4.11)$$

Consequently

$$\begin{aligned} \zeta_N &= \sum_{i=0}^{N-1} \bar{\eta}_i \geq \sum_{i=0}^{N-1} \mathbb{E}[\eta_i] \geq C_\eta^2 \sum_{i=0}^{N-1} i \geq C_\eta^2 \int_0^{N-2} x \, dx \geq \frac{C_\eta^2}{2} (N-2)^2 \\ &= \frac{C_\eta^2}{2} N^2 \left(\frac{N-2}{N} \right)^2 \geq \frac{C_\eta^2}{8} N^2, \quad (N \geq 4). \end{aligned} \quad (4.12)$$

Similarly

$$\zeta_{*,N} \geq \sum_{i=1}^{N-1} \mathbb{E}[\eta_i] \geq \frac{C_\eta^2}{2} ((N-2)^2 - 1) = \frac{3}{8} C_\eta^2 N^2, \quad (N \geq 4). \quad (4.13)$$

We still need to bound $\psi_{\ell,N+1}$ to bound $d_{\ell,N}^y$. To do this, we assume the existence of some j^* with $\gamma_{j^*} = 0$. With the help (C- κ .c), we then deduce from (R- η^2) that

$$\eta_i \leq \frac{(1-\delta)\psi_{\ell^*(j^*),0}}{\underline{\kappa} \mathbb{W}_{\ell^*(j^*)}^2} \phi_{j^*,i}.$$

Since $\gamma_{j^*} = 0$, a referral to (C- ϕ bnd.b) shows that

$$\mathbb{E}[\phi_{j^*,N}] = \phi_{j^*,0} + N\rho_{j^*}.$$

Consequently

$$\mathbb{E}[d_{y,\ell}^N] = \psi_{\ell,0}(\mathbb{E}[\eta_N] - 1) \leq \psi_{\ell,0} \left(\frac{(1-\delta)\psi_{\ell^*(j^*),0}}{\underline{\kappa}\mathbb{W}_{j^*}^2} \mathbb{E}[\phi_{j^*,N}] - 1 \right) \leq \psi_{\ell,0}(C_{\eta,*}N + \delta_*) \quad (4.14)$$

for

$$C_{\eta,*} := \frac{(1-\delta)\psi_{\ell^*(j^*),0}\rho_{j^*}}{\underline{\kappa}\mathbb{W}_{j^*}} \quad \text{and} \quad \delta_* := \frac{(1-\delta)\psi_{\ell^*(j^*),0}\phi_{j^*,0}}{\underline{\kappa}\mathbb{W}_{j^*}} - 1. \quad (4.15)$$

For the deterministic primal test update rule of Example 3.2, we still need to satisfy (C- ϕ det). Similarly to the derivation of (C- ϕ det'), we obtain for $\bar{\gamma}_j \geq 0$ the condition

$$\tilde{\gamma}_j = \bar{\gamma}_j = 0 \quad \text{or} \quad \frac{2\tilde{\gamma}_j\bar{\gamma}_j}{\tilde{\gamma}_j - \bar{\gamma}_j} \frac{(1-\delta)\psi_{\ell,0}}{\underline{\kappa}\mathbb{W}_j} \leq \delta, \quad (\ell \in \mathcal{V}(j)). \quad (\text{C-}\phi\text{det}'')$$

In summary:

Proposition 4.2. *Let $\delta \in (0,1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, \mathcal{Q})$. Pick $\rho_j \geq 0$ and $\tilde{\gamma}_j \in [0, \gamma_j]$, ($j = 1, \dots, m$). In Algorithm 1 or Algorithm 2, take*

- (i) *the probabilities $\hat{\pi}_{j,i} \equiv \hat{\pi}_j$ and (in Algorithm 1) $\hat{v}_{\ell,i} \equiv \hat{v}_\ell$ constant over iterations,*
- (ii) *η_i according to (R- η^2), and (in Algorithm 1) $\eta_{\tau,i}^\perp, \eta_{\sigma,i}^\perp > 0$ satisfying (C- η) and (C- η^\perp),*
- (iii) *$\phi_{j,0} > 0$ by free choice, and $\phi_{j,i}$ for $i \geq 1$ following Example 3.2, taking $0 \leq \bar{\gamma}_j < \tilde{\gamma}_j$ or $\bar{\gamma}_j = 0$, and satisfying (C- ϕ det'),*
- (iv) *$\psi_{\ell,i} := \eta_i\psi_{\ell,0}$ for some fixed $\psi_{\ell,0} > 0$, ($\ell = 1, \dots, n$).*

Suppose for each $j = 1, \dots, m$ that $\rho_j + \bar{\gamma}_j > 0$ and either $\rho_j = 0$ or (C-xbnd.a) holds with the constant C_x . Also assume that $\bar{\gamma}_{j^*} = 0$ for some $j^* \in \{1, \dots, m\}$, and that (C-ybnd.a) holds with the corresponding constant C_y . Let \tilde{c}_k be the constant provided by Lemma 4.1. Then

$$\sum_{k=1}^m \delta \tilde{c}_k \bar{\gamma}_k \mathbb{E}[\|x_k^N - \hat{x}_k\|^2] + \frac{C_\eta}{8} g_N \leq \frac{C_0 + C_x C_\alpha (\sum_{j=1}^m \rho_j) N + 9C_y \sum_{\ell=1}^n \psi_{\ell,0} (C_{\eta,*}N + \delta_*)}{N^2}$$

for $N \geq 4$ with

$$g_N := \begin{cases} \mathcal{G}(\tilde{x}_N, \tilde{y}_N), & \text{Algorithm 1, } \tilde{\gamma}_j \leq \gamma_j/2 \text{ for all } j, \\ \frac{3}{4} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}), & \text{Algorithm 2, } \tilde{\gamma}_j \leq \gamma_j/2 \text{ for all } j, \\ 0, & \text{otherwise.} \end{cases}$$

Proof. Similarly to the derivation of (4.5), by (R- η^2), $\eta_i \geq b_j \min_j \phi_{j,i}$, so (C- ϕ bnd.e) holds with $p = 1$. Therefore, the claim is immediate from (C- ϕ bnd.e), Corollaries 3.1 and 3.2, and the estimates (4.12)–(4.14), whose assumptions are provided by Lemmas 4.1 and 4.2 and (C- ϕ det''). \square

Remark 4.4. Note from (4.15) and the estimates of Proposition (4.2) that the factors ρ_j for such j that $\gamma_j = 0$ are very important for the convergence rate, and should therefore be chosen small. That is, we should not try to accelerate non-strongly-convex blocks very much, although some acceleration is necessary to obtain any estimates on the strongly convex blocks. As we will next see, without any acceleration or strong convexity at all, it is still however possible to obtain $O(1/N)$ convergence of the ergodic duality gap.

4.7. Unaccelerated algorithm

If $\rho_j = 0$ and $\tilde{y}_j = 0$ for all $j = 1, \dots, m$, then $\phi_{j,i} \equiv \phi_{j,0}$. Consequently (R- η) shows that $\eta_i \equiv \eta_0$. Recalling ζ_N from (2.14), we see that $\zeta_N = N\eta_0$. Likewise $\zeta_{*,N}$ from (2.18) satisfies $\zeta_{*,N} = (N-1)\eta_0$. Inserting this information into (3.30) in Proposition 3.1, we immediately obtain the following result.

Proposition 4.3. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, Q)$. In Algorithm 1 or 2, take*

- (i) $\phi_{j,i} \equiv \phi_{j,0} > 0$ constant between iterations,
- (ii) the probabilities $\pi_{j,i} \equiv \pi_j$ and (in Algorithm 1) $v_{\ell,i} \equiv v_\ell$ constant over iterations,
- (iii) $\psi_{\ell,i} \equiv \psi_{\ell,0}$ for some fixed $\psi_{\ell,0} > 0$, ($\ell = 1, \dots, n$), and
- (iv) $\eta_i \equiv \eta_0$, and (in Algorithm 1) $\eta_\tau^\perp, \eta_\sigma^\perp > 0$ satisfying (C- η).

Then

- (I) The iterates of Algorithm 1 satisfy $\mathcal{G}(\tilde{x}_N, \tilde{y}_N) \leq C_0 \eta_0^{-1}/N$, ($N \geq 1$).
- (II) The iterates of Algorithm 2 satisfy $\mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}) \leq C_0 \eta_0^{-1}/(N-1)$, ($N \geq 2$).

Remark 4.5. The obvious advantage of this unaccelerated algorithm is that in a parallel implementation no communication between different processors is necessary for the formation of η_i , which stays constant even with the random primal test update rule of Example 3.1.

4.8. Full primal strong convexity

Can we derive an $O(1/N^2)$ algorithm if G is full strongly convex? We still concentrate on the deterministic primal test updates of Example 3.2. Further, we follow the route of constant $\psi_{\ell,i} \equiv \psi_{\ell,0}$ in Section 4.5, as we seek to eliminate the penalties $d_{\ell,N}^y$ that any other choice would include in the convergence rates.

To eliminate the penalty $d_{j,N}^x(\tilde{y}_j)$, we take $\rho_j = 0$ and suppose $\gamma := \min_j \tilde{y}_j > 0$. This ensures (C-xbnd) as well as $\rho_j + \tilde{y}_j > 0$. The primal test update rule (R- ϕ det) then gives

$$\phi_{j,N} \geq \underline{\phi}_0 + \underline{\gamma} \sum_{i=0}^{N-1} \eta_i \geq \underline{\phi}_0 + \underline{\gamma} \sum_{i=0}^{N-1} \eta_i \quad \text{with} \quad \underline{\phi}_0 := \min_j \phi_{j,0} > 0.$$

Continuing from (4.5), therefore

$$\eta_N^2 \geq \underline{b}\underline{\phi}_0 + \underline{b}\underline{\gamma} \sum_{i=0}^{N-1} \eta_i \quad \text{with} \quad \underline{b} := \min_j b_j.$$

Otherwise written this says $\eta_N^2 \geq \tilde{\eta}_N^2$, where

$$\tilde{\eta}_N^2 = \underline{b}\underline{\phi}_0 + \underline{b}\underline{\gamma} \sum_{i=0}^{N-1} \tilde{\eta}_i = \tilde{\eta}_{N-1}^2 + c^2 \underline{\gamma} \tilde{\eta}_{N-1} = \tilde{\eta}_{N-1}^2 + \underline{b}\underline{\gamma} \tilde{\eta}_{N-1}^{-1}.$$

This implies by the estimates in [23] for the acceleration rule (2.3) that for some $Q_\eta > 0$ holds $\eta_i \geq \tilde{\eta}_i \geq Q_\eta i$. Replacing C_η by Q_η , repeating (4.12), (4.13) and (C- ϕ det'), and finally inserting the fact that now $\rho_j = 0$, we deduce:

Proposition 4.4. *Let $\delta \in (0, 1)$ and $(\kappa_1, \dots, \kappa_n) \in \mathcal{K}(K, \mathcal{P}, Q)$. Assume $\min_j \gamma_j > 0$, and pick $0 < \tilde{\gamma}_j \leq \gamma_j$, ($j = 1, \dots, m$). In Algorithm 1 or Algorithm 2, take*

- (i) the probabilities $\pi_{j,i} \equiv \pi_j$ and (in Algorithm 1) $\mathring{v}_{\ell,i} \equiv \mathring{v}_\ell$ constant over iterations,
- (ii) η_i according to (R- η), and (in Algorithm 1) $\eta_{\tau,i}^\perp, \eta_{\sigma,i}^\perp > 0$ satisfying (C- η) and (C- η^\perp),
- (iii) $\phi_{j,0} > 0$ by free choice, and $\phi_{j,i+1} := \phi_{j,i}(1 + 2\tilde{\gamma}_j\tau_{j,i})$, ($i \geq 1$), for some fixed $\tilde{\gamma}_j \in (0, \tilde{\gamma}_j)$, ($j = 1, \dots, m$), and
- (iv) $\psi_{\ell,i} := \psi_{\ell,0}$ for some fixed $\psi_{\ell,0} > 0$, ($\ell = 1, \dots, n$), satisfying (C- $\phi\text{det}'$).

Let \tilde{c}_k be the constant provided by Lemma 4.1. Then

$$\sum_{k=1}^m \delta \tilde{c}_k \tilde{\gamma}_k \mathbb{E} [\|x_k^N - \hat{x}_k\|^2] + g_N \leq \frac{8C_0}{Q_\eta N^2}, \quad (N \geq 4),$$

where

$$g_N := \begin{cases} \mathcal{G}(\tilde{x}_N, \tilde{y}_N), & \text{Algorithm 1, } \tilde{\gamma}_j \leq \gamma_j/2 \text{ for all } j, \\ \frac{3}{4} \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}), & \text{Algorithm 2, } \tilde{\gamma}_j \leq \gamma_j/2 \text{ for all } j, \\ 0, & \text{otherwise.} \end{cases}$$

Remark 4.6 (Linear rates under full primal-dual strong convexity). If both G and F^* are strongly convex, then it is possible to derive linear rates using the $\phi_{j,i+1}$ update rule of either Example 3.1 or Example 3.2, however fixing τ_i to a constant. This will cause $\mathbb{E}[\phi_{j,i}]$ to grow exponentially. Thanks to (C-step) exponential growth will also be the case for $\mathbb{E}[\eta_i]$. Through (4.4), also $\mathbb{E}[\psi_{\ell,i+1}]$ will grow exponentially. To counteract this, throughout the entire proof, starting from Theorems 2.1 and 2.2 we need to carry the strong convexity of F^* through the derivations similarly to how the strong convexity of G is carried in $\tilde{\Gamma}$ within $\Xi_{i+1}(\tilde{\Gamma})$ and $\Delta_{i+1}(\tilde{\Gamma})$.

4.9. Choices for $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$

We have not yet specified how exactly to choose $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$ in Algorithm 1, merely requiring the satisfaction of (C- η) and (C- η^\perp).

Example 4.1 (Constant $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$). We can take $\eta_{\tau,i}^\perp \equiv \eta_\tau^\perp$ and $\eta_{\sigma,i}^\perp \equiv \eta_\sigma^\perp$ for some $\eta_\sigma^\perp, \eta_\tau^\perp > 0$. This satisfies (C- η^\perp). Since our constructions of $i \mapsto \eta_i$ are increasing, and we assume fixed probabilities (R- πv), the condition (C- η) is satisfied for some $\epsilon \in (0, 1)$ if

$$\eta_0 \cdot \min_j (\pi_j - \pi_j) > \eta_\tau^\perp, \quad \text{and} \quad \eta_0 \cdot \min_\ell (v_\ell - \mathring{v}_\ell) > \eta_\sigma^\perp.$$

Example 4.2 (Proportional $\eta_{\tau,i}^\perp$ and $\eta_{\sigma,i}^\perp$). For some $\alpha \in (0, 1)$ let us take $\eta_{\tau,i}^\perp := \eta_\tau^\perp$ and $\eta_{\sigma,i}^\perp := \alpha \eta_i$. With the fixed probabilities (R- πv), this choice satisfies (C- η^\perp). The condition (C- η) holds for some $\epsilon \in (0, 1)$ if moreover

$$\min_j (\pi_j - \pi_j) > \alpha, \quad \text{and} \quad \min_\ell (v_\ell - \mathring{v}_\ell) \geq \alpha.$$

4.10. Sampling patterns

Since we for simplicity make the fixed probability assumption (R- πv), the only fully deterministic sampling patterns allowed are to consistently take $\mathring{S}(i) = \{1, \dots, m\}$ and $\mathring{V}(i+1) = \emptyset$, or alternatively $\mathring{S}(i) = \emptyset$ and $\mathring{V}(i+1) = \{1, \dots, n\}$. Regarding stochastic algorithms, let us first consider a few options for sampling $S(i)$ in Algorithm 2.

Example 4.3 (Independent unchanging probabilities). If all the blocks $\{1, \dots, m\}$ are chosen independently of each other, we have $\mathbb{P}(\{j, k\} \subset S(i)) = \pi_j \pi_k$ for $j \neq k$, where $\pi_j \in (0, 1]$.

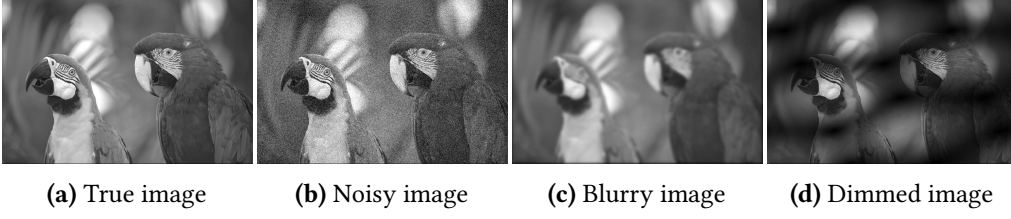


Figure 1: We use sample image (b) for TGV² denoising, (c) for TV deblurring, and (d) for TV undimming experiments.

Example 4.4 (Fixed number of random blocks). If we have a fixed number M of processors, we might choose a subset $S(i) \subset \{1, \dots, m\}$ such that $\#S(i) = M$.

The next example gives the simplest way to satisfy (C-nest.a) for Algorithm 1.

Example 4.5 (Alternating x-y and y-x steps). Let us randomly alternate between $\dot{S}(i) = \emptyset$ and $\dot{V}(i+1) = \emptyset$, choosing the non-empty set by a suitable sampling rule, such as those in Example 4.3 or Example 4.4. That is, with some probability \mathbb{p}_x , we choose to take an x-y step that omits lines 11 and 10 in Algorithm 1, and with probability $1 - \mathbb{p}_x$, an y-x step that omits the lines 6 and 9. If $\tilde{\pi}_j = \mathbb{P}[j \in \dot{S} | \dot{S} \neq \emptyset]$, and $\tilde{v}_\ell = \mathbb{P}[\ell \in \dot{V} | \dot{V} \neq \emptyset]$ denote the probabilities of the rule used to sample $\dot{S} = \dot{S}(i)$ and $\dot{V} = \dot{V}(i+1)$ when non-empty, then (C-nest) gives

$$\begin{aligned} \pi_j &= \mathbb{p}_x \tilde{\pi}_j, & \pi_j &= \mathbb{p}_x \tilde{\pi}_j + (1 - \mathbb{p}_x) \mathbb{P}[j \in \mathcal{V}^{-1}(\dot{V}) | \dot{V} \neq \emptyset], \\ v_\ell &= (1 - \mathbb{p}_x) \tilde{v}_\ell, & v_\ell &= (1 - \mathbb{p}_x) \tilde{v}_\ell + \mathbb{p}_x \mathbb{P}[\ell \in \mathcal{V}(\dot{S}) | \dot{S} \neq \emptyset]. \end{aligned}$$

To compute π_j and v_ℓ we thus need to know \mathcal{V} and the exact sampling pattern.

Remark 4.7. Based on Example 4.5, we can derive an algorithm where the only randomness comes from alternating between full x-y and y-x steps.

5. Numerical experience

We now apply several variants of the proposed algorithms to image processing problems. We generally consider discretisations, as our methods are formulated in Hilbert spaces, but the space of functions of bounded variation—where image processing problems are typically formulated—is only a Banach space. Our specific example problems will be TGV² denoising, TV deblurring, and TV undimming. In the latter, we solve

$$\min_{u \in \text{BV}(\Omega)} \frac{1}{2} \|f - \gamma \cdot u\|^2 + \alpha \text{TV}(u),$$

for a dimming mask $\gamma : \Omega \rightarrow \mathbb{R}$. In TGV² denoising and TV deblurring, we likewise use the L^2 -squared fidelity term, modelling Gaussian noise in the discretised setting.

We present the corrupt and ground-truth images in Figure 1, with values in the range $[0, 255]$. We use the images both at the original resolution of 768×512 , and scaled down to 192×128 pixels. To the noisy high-resolution test image in Figure 1b, we have added Gaussian noise with standard deviation 29.6 (12dB). In the downscaled image, this becomes 6.15 (25.7dB). The image in Figure 1c we have distorted with Gaussian blur of standard deviation 4. To avoid inverse crimes, we have added Gaussian noise of standard deviation 2.5. The dimmed image in Figure 1d, we have distorted by multiplying the image with a sinusoidal mask γ . The details of our construction beyond that seen in Figure 1c can be found in the source code [that will be archived per EPSRC regulations](#)

Table 1: Algorithm variants. The letters indicate how to read names like A-PRBO (P: only the primal variable x is randomly updated; R: the update rule for ϕ is the random one, and I , the increasing η rule of Proposition 4.2 is used; O: “Balanced” κ from Example 3.5).

Letter:	1st	2nd	3rd	4th
	Randomisation	ϕ rule	η and ψ rules	κ choice
A-	D: Deterministic P: Primal only B: Primal & Dual	R: Random, Ex. 3.1 D: Determ., Ex. 3.2 C: Constant	B: Bounded, Pr. 4.1 I: Increasing, Pr. 4.2	O: Balanc., Ex. 3.5 M: Max., Ex. 3.4

when the final version of the manuscript is submitted. Again, we have added the small amount of noise to the blurry image.

Besides the basic unaccelerated PDHGM (2.2)—note that our example problems are not strongly convex and hence the basic PDHGM cannot be accelerated—we evaluate our algorithms against the relaxed PDHGM of [44, 46], denoted in our results as ‘Relax’. In our precursor work [3], we have evaluated these two algorithms against the mixed-rate method of [47], and the adaptive PDHGM of [48]. To keep our tables and figures easily legible, we also do not include the algorithms of [3] in our evaluations. It is worth noting that even in the two-block case, the algorithms presented in this paper will not reduce to those of that paper: our rules for $\sigma_{\ell,i}$ are very different from the rules for the single σ_i therein.

We define abbreviations of our algorithm variants in Table 1. We do not report the results or apply all variants to all example problems. This would not be informative. We generally only consider the deterministic variants, as our problems are not large enough to benefit from being split on a computer cluster, where the benefits of the stochastic approaches would be apparent. We demonstrate the performance of the stochastic variants on TGV² denoising only.

To rely on Propositions 4.1 and 4.2 for convergence, we still need to satisfy (C-ybnd.a) and (C-xbnd.a), or take $\rho_j = 0$. The bound C_y in (C-ybnd) is easily calculated, as in all of our example problems, the functional F^* will restrict the dual variable to lie in a ball of known size. The primal variable, on the other hand, is not explicitly bounded. It is however possible to prove data-based conservative bounds on the optimal solution, see, e.g., [49, Appendix A]. We can therefore add an artificial bound to the problem to force all iterates to be bounded, replacing G by $\tilde{G}(x) := G(x) + \delta_{B(0, C_x)}(x)$. In practise, to avoid figuring out the exact magnitude of C_x , we update it dynamically, so as to avoid the constraint ever becoming active. It therefore does not affect the algorithm itself at all. In [49] a “pseudo duality gap” based on this idea was introduced. It is motivated by the fact that the real duality gap is also in practise infinite in numerical TGV² reconstructions. We will also use this type of dynamic duality gaps in our reporting: we take the bound C_x as the maximum over all iterations of all tested algorithms, and report the duality gap for the problem with \tilde{G} replacing G .

For each algorithm, we report the pseudo-duality gap, distance to a target solution, and function value. The target solution \hat{u} we compute by taking one million iterations of the basic PDHGM (2.2). In the calculation of the final duality gaps comparing each algorithm, we then take as C_x the maximum over all evaluations of all the algorithms. This makes the results fully comparable. We always report the pseudo-duality gap in decibels $10 \log_{10}(\text{gap}^2/\text{gap}_0^2)$ relative to the initial iterate. Similarly, we report the distance to the target solution \hat{u} in decibels $10 \log_{10}(\|u^i - \hat{u}\|^2/\|\hat{u}\|^2)$, and the primal objective value $\text{val}(x) := G(x) + F(Kx)$ relative to the target as $10 \log_{10}((\text{val}(x) - \text{val}(\hat{x}))^2/\text{val}(\hat{x})^2)$. Our computations were performed in Matlab+C-MEX on a MacBook Pro with 16GB RAM and a 2.8 GHz Intel Core i5 CPU.

5.1. TGV² denoising

In this problem, for regularisation parameters $\alpha, \beta > 0$, we have $x = (v, w)$ and $y = (\phi, \psi)$, with

$$G(x) = G_0(v), \quad K = \begin{pmatrix} \nabla & -I \\ 0 & \mathcal{E} \end{pmatrix}, \quad \text{and} \quad F^*(y) = \delta_{B(0, \alpha)}(\phi) + \delta_{B(0, \beta)}(\psi),$$

where the balls are pointwise in L^∞ , and \mathcal{E} the symmetrised gradient. Since there is no further spatial non-uniformity in this problem, it is natural to take as our projections $P_1x = v$, $P_2x = w$, $Q_1y = \phi$, and $Q_2y = \psi$. It is then not difficult to calculate the optimal κ_ℓ of Example 3.5, so we use only the ‘xxxO’ variants of the algorithms in Table 1.

As the regularisation parameters (β, α) , we choose (4.4, 4) for the downsampled image. For the original image we scale these parameters by $(0.25^{-2}, 0.25^{-1})$ corresponding to the image downscaling factor [50]. Since G is not strongly convex with respect to w , we have $\tilde{\gamma}_2 = 0$. For v we take $\tilde{\gamma}_1 = 1/2$, corresponding to the gap versions of our convergence estimates.

We take $\delta = 0.01$, and parametrise the standard PDHGM with $\sigma_0 = 1.9/\|K\|$ and $\tau_0 \approx 0.52/\|K\|$ solved from $\tau_0\sigma_0 = (1 - \delta)\|K\|^2$. These are values that typically work well. For forward-differences discretisation of TGV² with cell width $h = 1$, we have $\|K\|^2 \leq 11.4$ [49]. For the ‘Relax’ method from [46], we use the same σ_0 and τ_0 , as well as the value 1.5 for the inertial ρ parameter. For the increasing- ψ ‘xxIx’ variants of our algorithms, we take $\rho_1 = \rho_2 = 5$, $\tau_{1,0} = \tau_0$, and $\tau_{2,0} = 3\tau_0$. For the bounded- ψ ‘xxBx’ variants we take $\rho_1 = \rho_2 = 5$, $\tau_{1,0} = \tau_0$, and $\tau_{2,0} = 8\tau_0$. For both methods we also take $\eta_0 = 1/\tau_{0,1}$. These parametrisations force $\phi_{1,0} = 1/\tau_{1,0}^2$, and keep the initial step length $\tau_{1,0}$ for v consistent with the basic PDHGM. This justifies our algorithm comparisons using just a single set of parameters. To get an overview of the stability of convergence with respect to initialisation, we experiment with both initialisations $v^0 = 0$ and $v^0 = f$ the noisy image. The remaining variables w^0 , ϕ^0 , and ψ^0 we always initialise as zero.

The results for the deterministic variants of our algorithms are in Table 2 and Figure 2. For each algorithm we display the first 5000 iterations in a logarithmic fashion. To reduce computational overheads, we compute the reported quantities only every 10 iterations. To reduce the effects of other processes occasionally slowing down the computer, the CPU times reported are based on the average $\text{iteration_time} = \text{total_time}/\text{total_iterations}$, excluding time spent initialising the algorithm.

Our first observation is that the variants ‘xDxx’ based on the deterministic ϕ rule perform better than the “random” ϕ rule ‘xRxx’. Presently, with no randomisation, the only difference between the rules is the value of $\tilde{\gamma}$. The value 0.0105 from (C- $\phi\text{det}'$) and the value 0.0090 from (C- $\phi\text{det}''$) appear to give better performance than the maximal value $\tilde{\gamma}_1 = 0.5$. Generally, the A-DDBO seems to have the best asymptotic performance, with A-DRBO close. A-DDIO has good initial performance, although especially on the higher resolution image, the PDHGM and ‘Relax’ perform initially the best. Overall, however, the question of the best performer seems to be a rather fair competition between ‘Relax’ and A-DDBO.

5.2. TGV² denoising with stochastic algorithm variants

We also tested a few stochastic variants of our algorithms. We used the alternating sampling based on Example 4.4 with $M = 1$ and, when appropriate, Example 4.5. We took all probabilities equal to 0.5, that is $\mathbf{p}_x = \tilde{\pi}_1 = \tilde{\pi}_2 = \tilde{v}_1 = \tilde{v}_2 = 0.5$. In the doubly-stochastic ‘Bxxx’ variants of the algorithms, we have taken $\eta_{\tau,i}^\perp = \eta_{\sigma,i}^\perp = 0.9 \cdot 0.5\eta_i$.

The results are in Figure 3. To conserve space, we have only included a few descriptive algorithm variants. On the x axis, to better describe to the amount of actual work

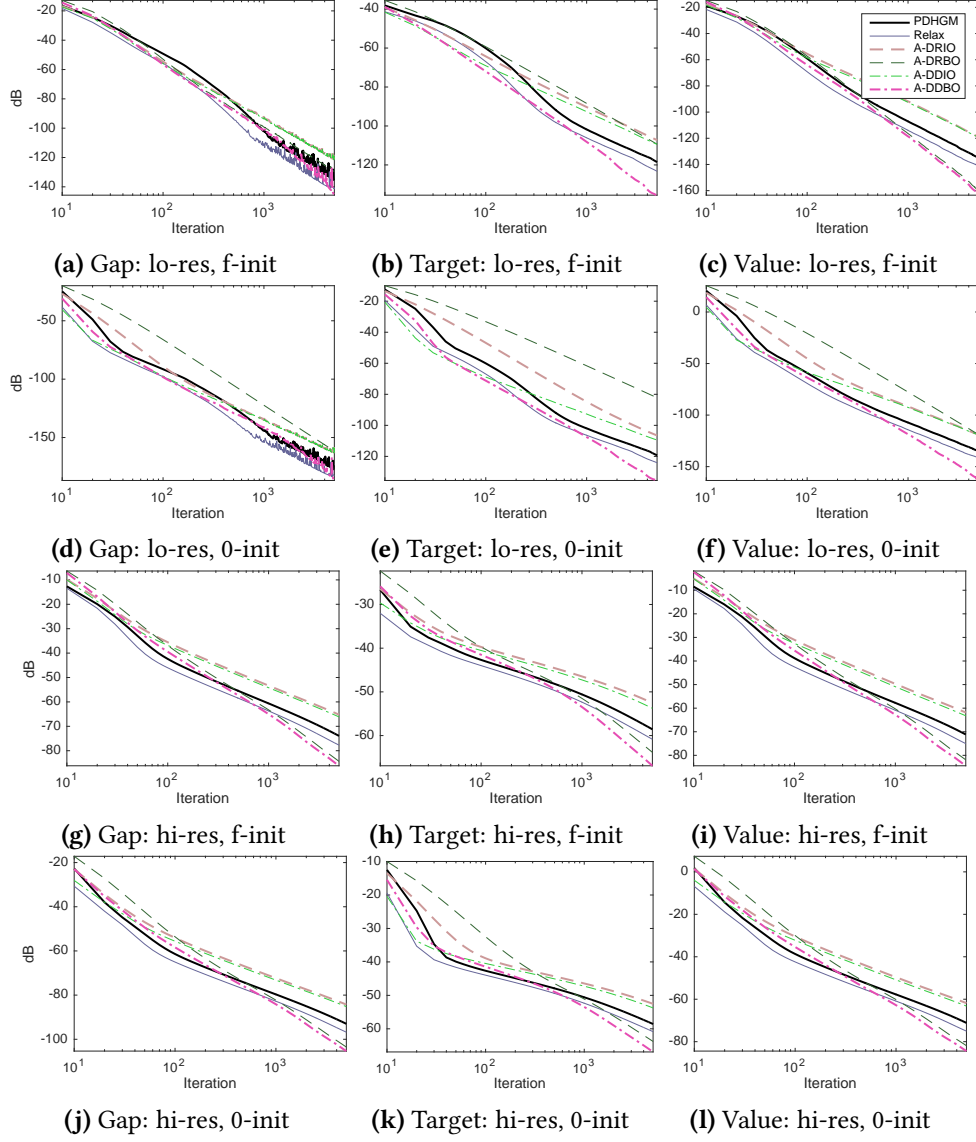


Figure 2: TGV² denoising performance of fully deterministic variants of our algorithms with pixelwise step lengths, 5000 iterations, high (hi-res) and low (lo-res) resolution images. Two different initialisations: $x^0 = 0$ (0-init) and $x^0 = (f, 0)$ (f-init). The plots are logarithmic.

performed by the stochastic methods, the “iteration” count refers to the *expected* number of full primal-dual updates. For all the displayed stochastic variants, with the present choice of probabilities, the expected number of full updates in each iteration is 0.75.

We run each algorithm 50 times, and plot for each iteration the 90% confidence interval according to Student’s t -distribution. Towards the 5000th iteration, these generally become very narrow, indicating reliability of the random method. Overall, the full-dual-update ‘Pxxx’ variants perform better than the doubly-stochastic ‘Bxxx’ variants. In particular, A-PDBO has performance comparable to or even better than the PDHGM.

5.3. TV deblurring

We now want to remove the blur in Figure 1c. We use TV parameter $\alpha = 2.55$ for the high resolution image and the scaled parameter $\alpha = 2.55 * 0.15$ for the low resolution image. We parametrise the PDHGM and ‘Relax’ algorithms exactly as for TGV² denoising above,

Table 2: TGV² denoising performance: CPU time and number of iterations (at a resolution of 10) to reach given duality gap, distance to target, or primal objective value.

low resolution / f-init							low resolution / 0-init						
Method	gap ≤ -60 dB	tgt ≤ -60 dB	val ≤ -60 dB	iter	time		gap ≤ -60 dB	tgt ≤ -60 dB	val ≤ -60 dB	iter	time	iter	time
PDHGM	190	1.29s	100	0.67s	110	0.74s	30	0.21s	100	0.72s	110	0.79s	
Relax	130	1.21s	70	0.64s	70	0.64s	20	0.20s	70	0.71s	70	0.71s	
A-DRIO	140	0.83s	80	0.47s	140	0.83s	40	0.26s	230	1.55s	180	1.22s	
A-DRBO	140	0.83s	110	0.65s	120	0.71s	80	0.54s	890	6.07s	500	3.41s	
A-DDIO	130	0.78s	50	0.29s	110	0.66s	20	0.14s	50	0.36s	110	0.80s	
A-DDBO	120	0.70s	50	0.29s	90	0.53s	30	0.19s	50	0.32s	90	0.58s	
high resolution / f-init							high resolution / 0-init						
Method	gap ≤ -50 dB	tgt ≤ -50 dB	val ≤ -50 dB	iter	time		gap ≤ -50 dB	tgt ≤ -50 dB	val ≤ -50 dB	iter	time	iter	time
PDHGM	250	32.17s	870	112.26s	370	47.67s	50	6.31s	870	111.83s	370	47.49s	
Relax	170	29.29s	580	100.34s	250	43.15s	40	6.93s	580	102.89s	250	44.25s	
A-DRIO	640	83.20s	2740	356.64s	1040	135.28s	70	9.17s	2750	365.52s	1050	139.48s	
A-DRBO	300	38.89s	790	102.63s	410	53.20s	80	10.56s	860	114.81s	420	56.00s	
A-DDIO	570	77.73s	2130	290.84s	900	122.81s	60	7.37s	2140	267.29s	900	112.34s	
A-DDBO	260	34.23s	600	79.16s	340	44.80s	60	7.85s	600	79.67s	340	45.09s	

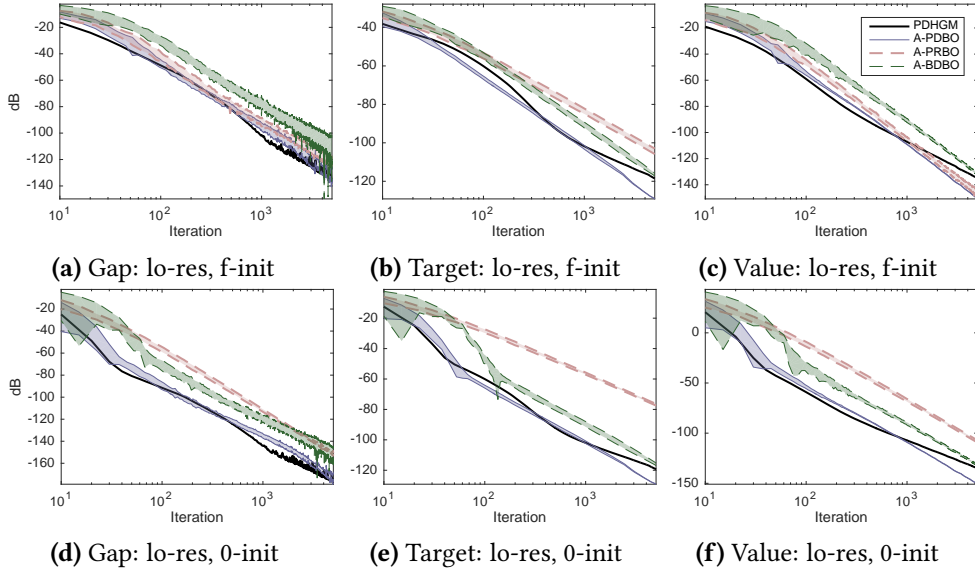


Figure 3: TGV² denoising performance of stochastic variants of our algorithms: 5000 iterations, low resolution images, initialisation both by zero and noisy data. Logarithmic plots with iteration counts scaled by the fraction of variables updated on average. We plot for each iteration the 90% confidence interval according to the t -distribution over 50 random runs.

with the natural difference of using the estimate $8 \geq \|K\|^2$ for $K = \nabla$ [51]. We write the forward blur operator as elementwise multiplication by factors $a = (a_1, \dots, a_m)$ in the discrete Fourier basis; that is $G(x) = \frac{1}{2} \|f - \mathcal{F}^*(a\mathcal{F}x)\|^2$ for \mathcal{F} the discrete Fourier transform. We then take as P_j the projection to the j :th Fourier component, and as Q_ℓ the projection to the ℓ :th pixel. This is to say that *each dual pixel and each primal Fourier component have their own step length parameter*. We initialise this as $\tau_{j,0} = \tau_0 / (\lambda + (1 - \lambda)\gamma_j)$, where the componentwise factor of strong convexity $\gamma_j = |a_j|^2$. For the bounded- ψ ‘xxBx’ algorithm variants we take $\lambda = 0.01$, and for the increasing- ψ ‘xxIx’ variants $\lambda = 0.1$.

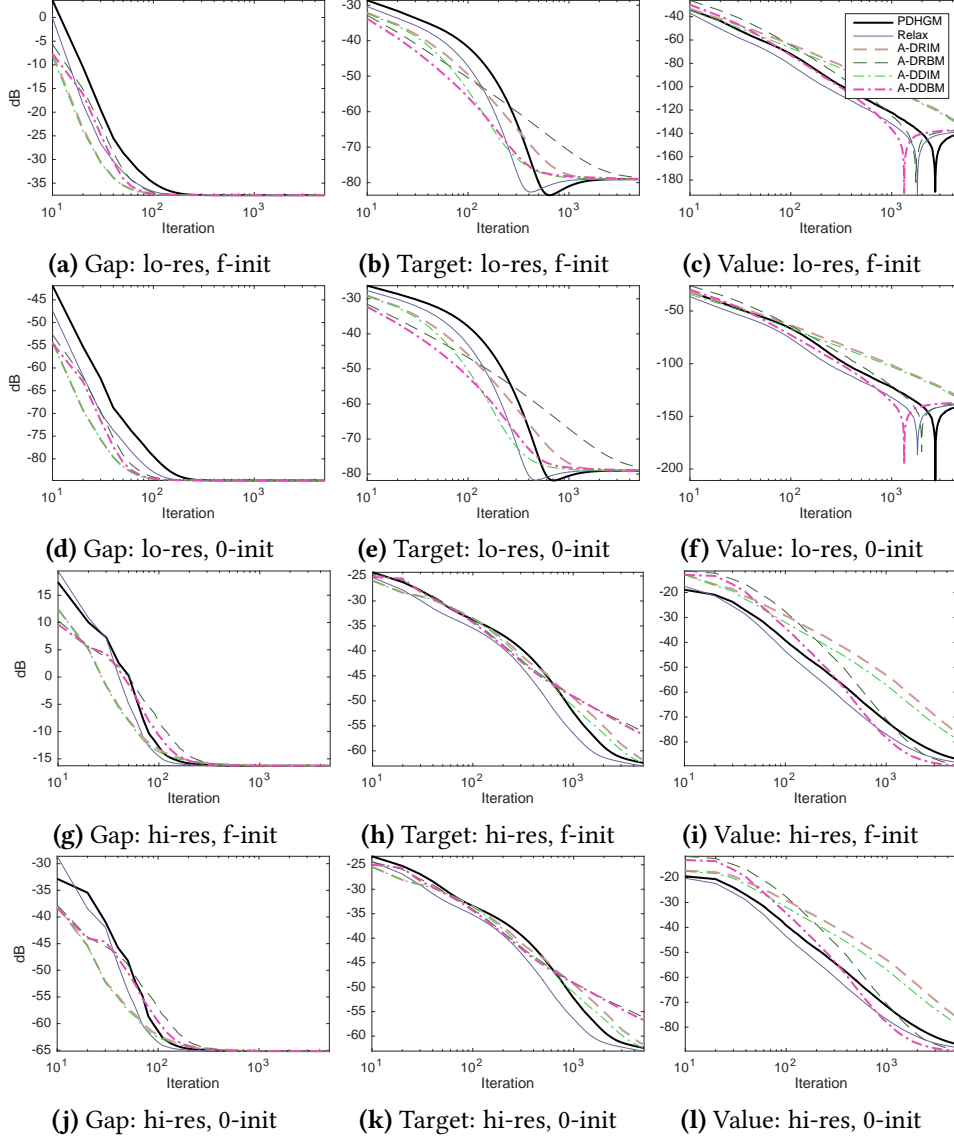


Figure 4: TV deblurring performance of fully deterministic variants of our algorithms with pixelwise step lengths, first 5000 iterations, high (hi-res) and low (lo-res) resolution images. Two different initialisations: $x^0 = 0$ (0-init) and $x^0 = f$ (f-init). The plots are logarithmic.

We only experiment with deterministic algorithms, as we do not expect small-scale randomisation to be beneficial. We also use the maximal κ ‘xxxM’ variants, as a more optimal κ would be very difficult to compute. The results are in Table 3 and Figure 4. Similarly to A-DDBO in our TGV² denoising experiments, A-DDBM performs reliably well, indeed better than the PDHGM or ‘Relax’. However, in many cases, A-DRBM and A-DDIM are even faster.

5.4. TV undimming

For TV undimming, our setup is exactly the same as TV deblurring, with the natural difference that the projection P_j are no longer to the Fourier basis, but to individual image pixels. The results are in Figure 5, and Table 4. They tell roughly the same story as TV deblurring, with A-DDBM performing well and reliably.

Table 3: TV deblurring performance: CPU time and number of iterations (at a resolution of 10) to reach given duality gap, distance to target, or primal objective value.

low resolution / f-init							low resolution / 0-init						
Method	gap ≤ -30 dB		tgt ≤ -60 dB		val ≤ -60 dB		Method	gap ≤ -60 dB		tgt ≤ -60 dB		val ≤ -60 dB	
	iter	time	iter	time	iter	time		iter	time	iter	time	iter	time
PDHGM	60	0.38s	280	1.78s	60	0.38s		30	0.18s	330	2.05s	70	0.43s
Relax	40	0.23s	190	1.14s	40	0.23s		20	0.11s	220	1.30s	50	0.29s
A-DRIM	30	0.22s	220	1.67s	80	0.60s		20	0.14s	280	2.08s	80	0.59s
A-DRBM	50	0.38s	310	2.42s	90	0.70s		20	0.14s	490	3.58s	90	0.65s
A-DDIM	30	0.22s	140	1.07s	70	0.53s		20	0.14s	170	1.25s	70	0.51s
A-DDBM	40	0.32s	140	1.12s	60	0.48s		20	0.15s	180	1.37s	60	0.45s

high resolution / f-init							high resolution / 0-init						
Method	gap ≤ -5 dB		tgt ≤ -40 dB		val ≤ -40 dB		Method	gap ≤ -50 dB		tgt ≤ -40 dB		val ≤ -40 dB	
	iter	time	iter	time	iter	time		iter	time	iter	time	iter	time
PDHGM	70	5.93s	330	28.25s	110	9.36s		60	5.04s	330	28.12s	110	9.31s
Relax	60	5.46s	220	20.25s	90	8.23s		50	4.32s	220	19.30s	90	7.84s
A-DRIM	40	4.44s	280	31.77s	310	35.19s		30	3.27s	280	31.41s	320	35.92s
A-DRBM	80	8.97s	240	27.14s	220	24.87s		60	6.48s	240	26.27s	220	24.07s
A-DDIM	40	4.43s	260	29.43s	230	26.02s		30	3.17s	260	28.35s	230	25.06s
A-DDBM	70	7.84s	230	26.00s	150	16.92s		50	5.56s	230	25.98s	150	16.90s

Table 4: TV undimming performance: CPU time and number of iterations (at a resolution of 10) to reach given duality gap, distance to target, or primal objective value.

low resolution / f-init							low resolution / 0-init						
Method	gap ≤ -80 dB		tgt ≤ -60 dB		val ≤ -60 dB		Method	gap ≤ -80 dB		tgt ≤ -60 dB		val ≤ -60 dB	
	iter	time	iter	time	iter	time		iter	time	iter	time	iter	time
PDHGM	110	0.30s	200	0.54s	120	0.32s		70	0.18s	200	0.51s	120	0.30s
Relax	70	0.16s	130	0.30s	80	0.18s		50	0.16s	130	0.41s	80	0.25s
A-DRIM	50	0.13s	150	0.39s	80	0.21s		30	0.10s	160	0.57s	80	0.28s
A-DRBM	40	0.10s	170	0.45s	60	0.16s		20	0.05s	170	0.47s	60	0.16s
A-DDIM	50	0.13s	100	0.25s	60	0.15s		30	0.08s	110	0.30s	60	0.16s
A-DDBM	30	0.07s	70	0.18s	40	0.10s		20	0.05s	70	0.18s	40	0.10s

high resolution / f-init							high resolution / 0-init						
Method	gap ≤ -80 dB		tgt ≤ -60 dB		val ≤ -60 dB		Method	gap ≤ -80 dB		tgt ≤ -60 dB		val ≤ -60 dB	
	iter	time	iter	time	iter	time		iter	time	iter	time	iter	time
PDHGM	170	5.75s	290	9.83s	210	7.11s		100	3.41s	300	10.31s	210	7.21s
Relax	110	4.42s	200	8.07s	140	5.64s		70	3.03s	200	8.73s	140	6.10s
A-DRIM	320	13.44s	750	31.56s	630	26.50s		80	3.52s	760	33.82s	640	28.48s
A-DRBM	240	9.87s	370	15.24s	380	15.65s		90	3.95s	370	16.39s	380	16.84s
A-DDIM	240	10.03s	570	23.88s	420	17.58s		70	3.05s	580	25.57s	430	18.94s
A-DDBM	140	5.84s	230	9.61s	200	8.35s		60	2.63s	230	10.22s	200	8.88s

Conclusions

We have derived from abstract theory several accelerated block-proximal primal-dual methods, both stochastic and deterministic. So far, we have primarily concentrated on applying them deterministically, taking advantage of blockwise—indeed pixelwise—factors of strong convexity, to obtain improved performance compared to standard methods. In future work, it will be interesting to evaluate the methods on real large scale problems to other state-of-the-art stochastic optimisation methods. Moreover, interesting questions include heuristics and other mechanisms for optimal initialisation of the pixelwise parameters.

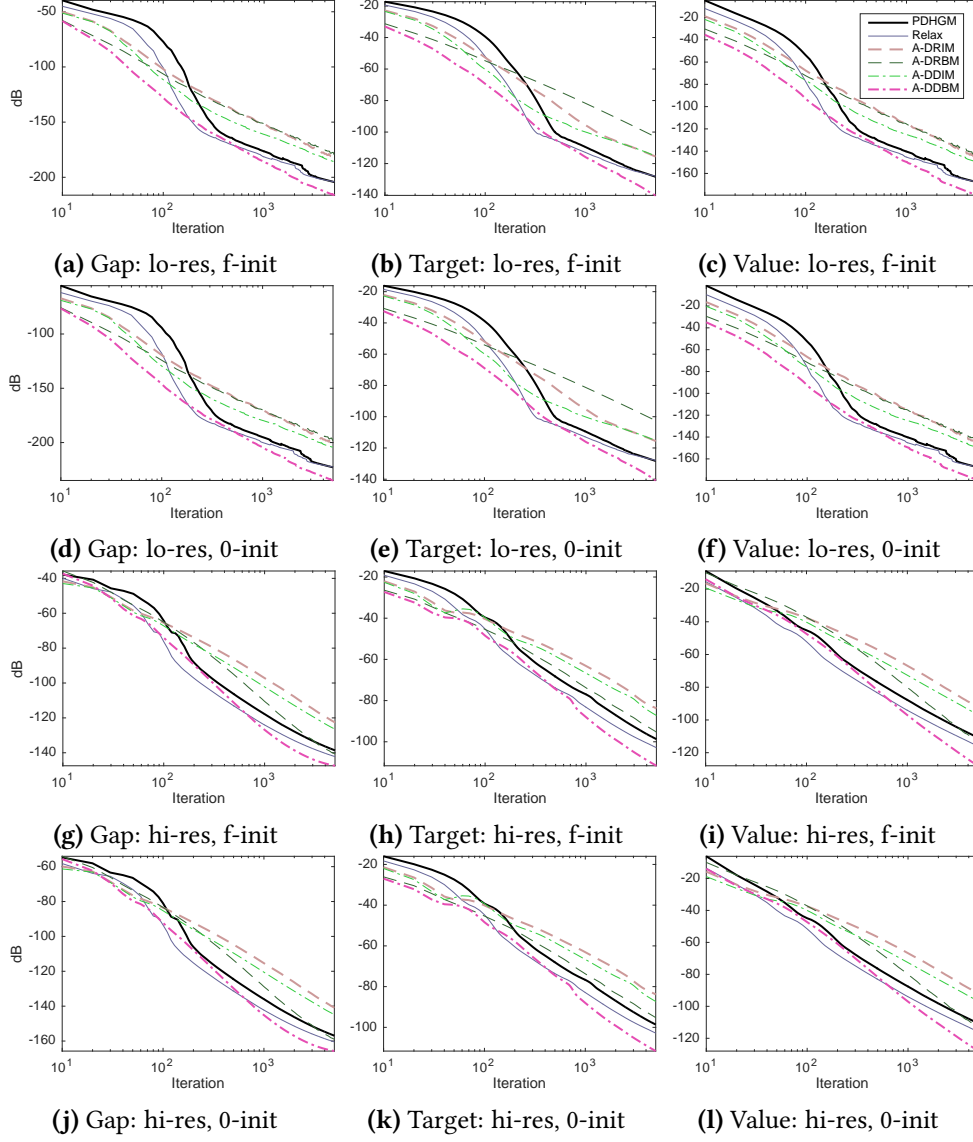


Figure 5: TV undimming performance of fully deterministic variants of our algorithms with pixelwise step lengths, 5000 iterations, high (hi-res) and low (lo-res) resolution images. Two different initialisations: $x^0 = 0$ (0-init) and $x^0 = f$ (f-init). The plots are logarithmic.

Acknowledgements

The author would like to thank Peter Richtárik and Olivier Fercoq for several fruitful discussions, and for introducing him to stochastic optimisation. Moreover, the support of the EPSRC grant EP/M00483X/1 “Efficient computational tools for inverse imaging problems” is acknowledged during the initial two months of the research.

A data statement for the EPSRC

The codes will be archived once the final version of the paper is submitted. The sample photo is from the free Kodak image suite, at the time of writing available online at <http://r0k.us/graphics/kodak/>.

A. Proofs of the general estimates

Here, we prove Theorems 2.1 to 2.3. In fact, as the technique streamlines the proof, we consider instead of (PP) for convex $V_{i+1} : X \times Y \rightarrow \overline{\mathbb{R}}$ the general iteration

$$0 \in Z_{i+1}W_{i+1}H(u^{i+1}) + \partial V_{i+1}(u^{i+1}). \quad (\text{B-PP})$$

Our motivation is that under (C0), $Z_{i+1}L_{i+1}$ is self-adjoint, which has the effect

$$Z_{i+1}L_{i+1}(u^{i+1} - u^i) = \nabla V_{i+1}(u^{i+1}) \quad \text{for} \quad V_{i+1}(u) := \frac{1}{2}\|u - u^i\|_{Z_{i+1}L_{i+1}}^2. \quad (\text{A.1})$$

In (B-PP), we therefore take V_i to be an arbitrary convex function satisfying for each $i \in \mathbb{N}$ for some $\Delta_{i+1} \in \mathcal{L}(X \times Y; X \times Y)$ and $\delta_i \geq 0$ the condition

$$V_{i+1}(u) + \delta_{i+1} \leq V_i(u) + \langle p, u^i - u \rangle + \frac{1}{2}\|u - u^i\|_{Z_i\Xi_i(\Gamma) + \Delta_{i+1}(\Gamma)}^2, \quad (u \in X \times Y, p \in \partial V_i(u^i)). \quad (\text{B-C0})$$

Example A.1. For the choice (A.1), using (2.10) in (B-C0), we obtain the requirement

$$0 \leq \frac{1}{2}\|u - u^{i-1}\|_{Z_iL_i}^2 + \langle p, u^i - u \rangle - \frac{1}{2}\|u - u^i\|_{Z_iL_i}^2, \quad (u \in X \times Y), \quad (\text{A.2})$$

where $p = L_i^*Z_i^*(u^i - u^{i-1})$. We observe for self-adjoint M the identity

$$\langle u^i - u^{i-1}, u^i - u \rangle_M = \frac{1}{2}\|u^i - u^{i-1}\|_M^2 - \frac{1}{2}\|u^{i-1} - u\|_M^2 + \frac{1}{2}\|u^i - u\|_M^2.$$

Applying this to $M = Z_iL_i$, we verify (A.2) for $\delta_i := \frac{1}{2}\|u^i - u^{i-1}\|_{Z_iL_i}^2$. By (C0), $\delta_i \geq 0$.

Example A.2 (Bregman distances, cf. [52]). Suppose that V is a Bregman distance, that is, for some convex function J and some $p_i \in \partial J(u^i)$ to be fixed, we take

$$V_{i+1}(u) = J(u) - \langle p_i, u - u^i \rangle - J(u^i).$$

Then (B-C0) is satisfied $\Delta_{i+1}(\Gamma) = Z_i\Xi_i(\Gamma)$ if

$$J(u) - \langle p_i, u - u^{i+1} \rangle - J(u^{i+1}) \leq J(u) - \langle p_{i-1}, u - u^i \rangle - J(u^i) + \langle p, u^{i+1} - u \rangle. \quad (\text{A.3})$$

We have $p \in \partial V_i(u^i)$ if and only if $p \in \partial J(u^i) - p_{i-1}$. Taking $p_i := p + p_{i-1}$, (A.3) becomes

$$\langle p + p_i, u^{i+1} - u \rangle - J(u^{i+1}) \leq -\langle p_i, u - u^i \rangle - J(u^i) + \langle p, u^{i+1} - u \rangle.$$

In other words

$$J(u^i) + \langle p_i, u^{i+1} - u^i \rangle \leq J(u^{i+1}).$$

This automatically holds by the convexity of J . Since there is no acceleration, Theorem A.1 will not give any kind of convergence rates for this choice of V_i . We at most obtain the convergence of the ergodic gap from Lemma A.1.

For better rates, we need more structure, and to incorporate acceleration parameters in V_i . Normally in the algorithm (2.2) & (2.3), we would for $c = \|K\|^2/(1 - \delta)$ have

$$Z_{i+1}L_{i+1} = \begin{pmatrix} \tau_i^{-2}I & -\tau_i^{-1}K^* \\ -\tau_i^{-1}K & cI \end{pmatrix}$$

This suggests to use

$$V_{i+1}(u) = \tau_i^{-2}V_{x,i+1}(x) + cV_{y,i+1}(y) - \tau_i^{-1}\langle K(x - x^i), y - y^i \rangle$$

for $V_{x,i+1}$ and $V_{y,i+1}$ the Bregman distances corresponding to some J_x and J_y at x^i and y^i . If J_x and J_y possess sufficient strong convexity, V_{i+1} will also be strongly convex. In that case, we could by adapting (C1) and the analysis of Example A.2, get convergence similar to [52].

We have the following abstract convergence estimate. To keep the necessary setup short, and to avoid introducing additional notation, we do not represent the result in the full generality of Banach spaces. Such a generalisation is, however, immediate.

Theorem A.1. *Let us be given $K \in \mathcal{L}(X; Y)$, and convex, proper, lower semicontinuous functionals $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ on Hilbert spaces X and Y , satisfying (G-PM) and (F*-PM) for some $0 \leq \Gamma \in \mathcal{L}(X; X)$. Suppose (B-PP) is solvable, and that (B-C0) is satisfied for each $i \in \mathbb{N}$ for some operators $T_i, \Phi_i \in \mathcal{L}(X; X)$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{L}(Y; Y)$, with $\Phi_i T_i \in \mathcal{T}$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{S}$. Then the iterates $u^i = (x^i, y^i)$ of (B-PP) satisfy*

$$V_{N+1}(\hat{u}) + \sum_{i=0}^{N-1} \delta_{i+1} \leq V_0(\hat{u}) + \frac{1}{2} \sum_{i=0}^{N-1} \|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\Gamma)}^2, \quad (N \geq 1) \quad (\text{A.4})$$

Proof. We take $p_{i+1} \in H(u^{i+1})$ such that $-Z_{i+1}W_{i+1}p_{i+1} \in \partial V_{i+1}(u^{i+1})$, guaranteed to exist by the iteration (B-PP) being by assumption solvable. Using the expansion

$$Z_{i+1}W_{i+1} = \begin{pmatrix} \Phi_i T_i & 0 \\ 0 & \Psi_{i+1} \Sigma_{i+1} \end{pmatrix},$$

and the fact that $0 \in H(\hat{u})$, we deduce

$$\begin{aligned} \langle p_{i+1}, W_{i+1}^* Z_{i+1}^* (u^{i+1} - \hat{u}) \rangle &\subset \langle H(u^{i+1}) - H(\hat{u}), W_{i+1}^* Z_{i+1}^* (u^{i+1} - \hat{u}) \rangle \\ &= \langle \partial G(x^{i+1}) - \partial G(\hat{x}), T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle \\ &\quad + \langle \partial F^*(y^{i+1}) - \partial F^*(\hat{y}), \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle \\ &\quad + \langle K^*(y^{i+1} - \hat{y}), T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle \\ &\quad - \langle K(x^{i+1} - \hat{x}), \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle. \end{aligned}$$

An application of (G-PM) and (F*-PM) consequently gives

$$\begin{aligned} \langle p_{i+1}, W_{i+1}^* Z_{i+1}^* (u^{i+1} - \hat{u}) \rangle &\geq \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma}^2 \\ &\quad + \langle \Phi_i T_i K^* (y^{i+1} - \hat{y}), x^{i+1} - \hat{x} \rangle \\ &\quad - \langle \Psi_{i+1} \Sigma_{i+1} K (x^{i+1} - \hat{x}), y^{i+1} - \hat{y} \rangle. \quad (\text{A.5}) \\ &= \frac{1}{2} \|u^{i+1} - \hat{u}\|_{Z_{i+1} \Xi_{i+1}(\Gamma)}^2. \end{aligned}$$

Next, (B-C0) at $u = \hat{u}$ gives

$$V_{i+1}(\hat{u}) - V_{i+2}(\hat{u}) - \delta_{i+1} + \frac{1}{2} \|u^{i+1} - \hat{u}\|_{Z_{i+1} \Xi_{i+1}(\Gamma) + \Delta_{i+2}(\Gamma)}^2 \geq \langle Z_{i+1}W_{i+1}p_{i+1}, u^{i+1} - \hat{u} \rangle. \quad (\text{A.6})$$

Combining (A.5) with (A.6), we thus deduce

$$V_{i+2}(\hat{u}) + \delta_{i+1} \leq V_{i+1}(\hat{u}) + \frac{1}{2} \|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\Gamma)}^2. \quad (\text{A.7})$$

Summing (A.7) over $i = 0, \dots, N-1$, we obtain (A.4). \square

Proof of Theorem 2.1. Insert V_i from (A.1) into (A.4), and use Theorem A.1 and Example A.1. \square

Lemma A.1. *Let X and Y be Hilbert spaces, $K \in \mathcal{L}(X; Y)$, and $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ be convex, proper, lower semicontinuous. Take $0 \leq \Gamma \in \mathcal{L}(X; X)$. Suppose (B-PP) is solvable, and that (B-C0) is satisfied for each $i \in \mathbb{N}$ for some $T_i, \Phi_i \in \mathcal{L}(X; X)$ and*

$\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{L}(Y; Y)$, with $\Phi_i T_i \in \mathcal{T}$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{S}$. Then the iterates $u^i = (x^i, y^i)$ of (B-PP) satisfy

$$V_{N+1}(\hat{u}) + \sum_{i=0}^{N-1} (\mathcal{G}'_{i+1} + \delta_{i+1}) \leq V_0(\hat{u}) + \frac{1}{2} \sum_{i=0}^{N-1} \|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\Gamma/2)}^2, \quad (N \geq 1) \quad (\text{A.8})$$

for

$$\begin{aligned} \mathcal{G}'_{i+1} := & \langle \partial G(x^{i+1}), T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle - \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma/2}^2 \\ & + \langle \partial F^*(y^{i+1}), \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle \\ & - \langle \hat{y}, (KT_i^* \Phi_i^* - \Psi_{i+1} \Sigma_{i+1} K) \hat{x} \rangle - \langle y^{i+1}, \Psi_{i+1} \Sigma_{i+1} K \hat{x} \rangle + \langle \hat{y}, KT_i^* \Phi_i^* x^{i+1} \rangle. \end{aligned}$$

Proof. Similarly to the proof of Theorem 2.1, we take $p_{i+i} \in H(u^{i+1})$ such that $-Z_{i+1} p_{i+i} \in \partial V_{i+1}(u^{i+1})$, guaranteed to exist by the iteration (B-PP) being by assumption solvable. Then

$$\begin{aligned} \langle Z_{i+1} W_{i+1} p_{i+i}, u^{i+1} - \hat{u} \rangle &= \langle \partial G(x^{i+1}), T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle \\ &\quad + \langle \partial F^*(y^{i+1}), \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle \\ &\quad + \langle K^* y^{i+1}, T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle - \langle K x^{i+1}, \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle \\ &= g_{i+1} + \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma/2}^2 \\ &\quad + \langle y^{i+1}, (KT_i^* \Phi_i^* - \Psi_{i+1} \Sigma_{i+1} K) x^{i+1} \rangle \\ &\quad - \langle y^{i+1}, KT_i^* \Phi_i^* \hat{x} \rangle + \langle \hat{y}, \Psi_{i+1} \Sigma_{i+1} K x^{i+1} \rangle. \end{aligned}$$

for

$$g_{i+1} := \langle \partial G(x^{i+1}), T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle - \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma/2}^2 + \langle \partial F^*(y^{i+1}), \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle.$$

A little bit of reorganisation and referral to the expression for Δ_{i+1} in (2.9) gives

$$\begin{aligned} \langle Z_{i+1} W_{i+1} p_{i+i}, u^{i+1} - \hat{u} \rangle &\geq g_{i+1} + \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma/2}^2 \\ &\quad + \langle y^{i+1} - \hat{y}, (KT_i^* \Phi_i^* - \Psi_{i+1} \Sigma_{i+1} K) (x^{i+1} - \hat{x}) \rangle \\ &\quad - \langle \hat{y}, (KT_i^* \Phi_i^* - \Psi_{i+1} \Sigma_{i+1} K) \hat{x} \rangle \\ &\quad - \langle y^{i+1}, \Psi_{i+1} \Sigma_{i+1} K \hat{x} \rangle + \langle \hat{y}, KT_i^* \Phi_i^* x^{i+1} \rangle \\ &= \mathcal{G}'_{i+1} + \frac{1}{2} \|u^{i+1} - \hat{u}\|_{\Xi_{i+1}(\Gamma/2)}^2. \end{aligned}$$

Combining (A.6) and (A.9), we obtain

$$V_{i+2}(\hat{u}) + \mathcal{G}'_{i+1} + \delta_{i+1} \leq V_{i+1}(\hat{u}) + \frac{1}{2} \|u^{i+1} - \hat{u}\|_{\Delta_{i+2}(\Gamma/2)}^2.$$

Summing this for $i = 0, \dots, N-1$ gives (A.8). \square

Proof of Theorem 2.2. We use Lemma A.1. We already verified (B-C0) in Example A.1, so it remains to derive (2.17). Using (C-G), (G-EC), and (F*-EC), we compute

$$\begin{aligned} \sum_{i=0}^{N-1} \mathbb{E}[\mathcal{G}'_{i+1}] &:= \sum_{i=0}^{N-1} \mathbb{E} \left[\langle \partial G(x^{i+1}), T_i^* \Phi_i^* (x^{i+1} - \hat{x}) \rangle - \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma/2}^2 \right] \\ &\quad + \mathbb{E} \left[\langle \partial F^*(y^{i+1}), \Sigma_{i+1}^* \Psi_{i+1}^* (y^{i+1} - \hat{y}) \rangle \right] - \zeta_N \langle \tilde{y}_N, K \hat{x} \rangle + \zeta_N \langle \hat{y}, K \tilde{x}_N \rangle \\ &\geq \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N) - \zeta_N \langle \tilde{y}_N, K \hat{x} \rangle + \zeta_N \langle \hat{y}, K \tilde{x}_N \rangle. \end{aligned}$$

We therefore obtain (2.17) by taking the expectation in (A.8). \square

Proof of Theorem 2.3. Using (G-PM) and (OC), we deduce

$$\mathcal{G}'_1 \geq \langle \partial F^*(y^1), \Sigma_1^* \Psi_1^*(y^1 - \hat{y}) \rangle + \langle \hat{y}, \Psi_1 \Sigma_1 K \hat{x} \rangle - \langle y^1, \Psi_1 \Sigma_1 K \hat{x} \rangle.$$

Likewise (F*-PM) and (OC) give

$$\begin{aligned} \mathcal{G}'_N &\geq \langle \partial G(x^N), T_{N-1}^* \Phi_{N-1}^*(x^N - \hat{x}) \rangle - \|x^N - \hat{x}\|_{\Phi_{N-1} T_{N-1} \Gamma/2}^2 \\ &\quad - \langle \hat{y}, K T_{N-1}^* \Phi_{N-1}^* \hat{x} \rangle + \langle \hat{y}, K T_{N-1}^* \Phi_{N-1}^* x^N \rangle. \end{aligned}$$

Shifting indices of y^i by one compared to \mathcal{G}'_{i+1} , we define

$$\begin{aligned} \mathcal{G}'_{*,i+1} &:= \langle \partial G(x^{i+1}), T_i^* \Phi_i^*(x^{i+1} - \hat{x}) \rangle - \|x^{i+1} - \hat{x}\|_{\Phi_i T_i \Gamma/2}^2 \\ &\quad + \langle \partial F^*(y^i), \Sigma_i^* \Psi_i^*(y^i - \hat{y}) \rangle \\ &\quad - \langle \hat{y}, (K T_i^* \Phi_i^* - \Psi_i \Sigma_i K) \hat{x} \rangle - \langle y^i, \Psi_i \Sigma_i K \hat{x} \rangle + \langle \hat{y}, K T_i^* \Phi_i^* x^{i+1} \rangle, \end{aligned}$$

Correspondingly reorganising terms, we observe

$$\sum_{i=0}^{N-1} \mathcal{G}'_{i+1} = \mathcal{G}'_1 + \sum_{i=1}^{N-2} \mathcal{G}'_{i+1} + \mathcal{G}'_N \geq \sum_{i=1}^{N-1} \mathcal{G}'_{*,i+1}.$$

We now estimate $\sum_{i=1}^{N-1} \mathbb{E}[\mathcal{G}'_{*,i+1}]$ analogously to the proof of Theorem 2.2. \square

B. An inequality

We needed the following for convergence rates in Section 4.2.

Lemma B.1. *Suppose $\phi_N \geq \phi_0 + b \sum_{i=0}^{N-1} (i+1)^p$ for each $N \geq 0$ for some constants $p \geq 0$ and $\phi_0, b > 0$. Then $\phi_N \geq \phi_0 + C N^{p+1}$ for some constant $C = C(b, \phi_0, p) > 0$.*

Proof. We calculate

$$\phi_N \geq \phi_0 + b \sum_{i=1}^N i^p \geq \phi_0 + b \int_2^N x^p dx \geq \phi_0 + p^{-1} b (N^{p+1} - 2).$$

The lower bound $\phi_N \geq \phi_0$ for $0 \leq N \leq 2$, and suitably choice of $C > 0$ verify the claim. \square

References

- [1] J. Bolte, S. Sabach and M. Teboulle, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Mathematical Programming **146** (2013), 459–494, doi:[10.1007/s10107-013-0701-9](https://doi.org/10.1007/s10107-013-0701-9).
- [2] T. Möllenhoff, E. Strekalovskiy, M. Moeller and D. Cremers, *The primal-dual hybrid gradient method for semiconvex splittings*, SIAM Journal on Imaging Sciences **8** (2015), 827–857, doi:[10.1137/140976601](https://doi.org/10.1137/140976601).
- [3] T. Valkonen and T. Pock, *Acceleration of the PDHGM on partially strongly convex functions* (2015), submitted, [arXiv:1511.06566](https://arxiv.org/abs/1511.06566).
URL <http://iki.fi/tuomov/mathematics/cpaccel.pdf>
- [4] P. Ochs, Y. Chen, T. Brox and T. Pock, *iPiano: Inertial proximal algorithm for non-convex optimization* (2014), preprint, [arXiv:1404.4805](https://arxiv.org/abs/1404.4805).

- [5] L. Rudin, S. Osher and E. Fatemi, *Nonlinear total variation based noise removal algorithms*, Physica D **60** (1992), 259–268.
- [6] K. Bredies, K. Kunisch and T. Pock, *Total generalized variation*, SIAM Journal on Imaging Sciences **3** (2011), 492–526, doi:[10.1137/090769521](https://doi.org/10.1137/090769521).
- [7] A. Chambolle and P.-L. Lions, *Image recovery via total variation minimization and related problems*, Numerische Mathematik **76** (1997), 167–188, doi:[10.1007/s002110050258](https://doi.org/10.1007/s002110050258).
- [8] I. Daubechies, M. Defrise and C. De Mol, *An iterative thresholding algorithm for linear inverse problems with a sparsity constraint*, Communications on Pure and Applied Mathematics **57** (2004), 1413–1457, doi:[10.1002/cpa.20042](https://doi.org/10.1002/cpa.20042).
- [9] I. Loris and C. Verhoeven, *On a generalization of the iterative soft-thresholding algorithm for the case of non-separable penalty*, Inverse Problems **27** (2011), 125007, doi:[10.1088/0266-5611/27/12/125007](https://doi.org/10.1088/0266-5611/27/12/125007).
- [10] A. Beck and M. Teboulle, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM Journal on Imaging Sciences **2** (2009), 183–202, doi:[10.1137/080716542](https://doi.org/10.1137/080716542).
- [11] S. Wright, *Coordinate descent algorithms*, Mathematical Programming **151** (2015), 3–34, doi:[10.1007/s10107-015-0892-3](https://doi.org/10.1007/s10107-015-0892-3).
- [12] Y. Nesterov, *Efficiency of coordinate descent methods on huge-scale optimization problems*, SIAM Journal on Optimization **22** (2012), 341–362, doi:[10.1137/100802001](https://doi.org/10.1137/100802001).
- [13] P. Richtárik and M. Takáč, *Parallel coordinate descent methods for big data optimization*, Mathematical Programming (2015), 1–52, doi:[10.1007/s10107-015-0901-6](https://doi.org/10.1007/s10107-015-0901-6).
- [14] O. Fercoq and P. Richtárik, *Accelerated, parallel and proximal coordinate descent* (2013), preprint, [arXiv:1312.5799](https://arxiv.org/abs/1312.5799).
- [15] P. Richtárik and M. Takáč, *Distributed coordinate descent method for learning with big data* (2013), [arXiv:1310.2059](https://arxiv.org/abs/1310.2059).
- [16] Z. Qu, P. Richtárik and T. Zhang, *Randomized dual coordinate ascent with arbitrary sampling* (2014), preprint, [arXiv:1411.5873](https://arxiv.org/abs/1411.5873).
- [17] P. Zhao and T. Zhang, *Stochastic optimization with importance sampling* (2014), preprint, [arXiv:1401.2753](https://arxiv.org/abs/1401.2753).
- [18] S. Shalev-Shwartz and T. Zhang, *Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization*, Mathematical Programming **155** (2014), 105–145, doi:[10.1007/s10107-014-0839-0](https://doi.org/10.1007/s10107-014-0839-0).
- [19] Z. Qu, P. Richtárik, M. Takáč and O. Fercoq, *SDNA: stochastic dual newton ascent for empirical risk minimization* (2015), [arXiv:1502.02268](https://arxiv.org/abs/1502.02268).
- [20] D. Csiba, Z. Qu and P. Richtárik, *Stochastic dual coordinate ascent with adaptive probabilities*, preprint, [arXiv:1502.08053](https://arxiv.org/abs/1502.08053).
- [21] P. L. Combettes and J.-C. Pesquet, *Stochastic forward-backward and primal-dual approximation algorithms with application to online image restoration* (2016), [arXiv:1602.08021](https://arxiv.org/abs/1602.08021).

- [22] Z. Peng, Y. Xu, M. Yan and W. Yin, *ARock: an algorithmic framework for asynchronous parallel coordinate updates* (2015), uCLA CAM Report 15-37.
URL <ftp://ftp.math.ucla.edu/pub/camreport/cam15-37.pdf>
- [23] A. Chambolle and T. Pock, *A first-order primal-dual algorithm for convex problems with applications to imaging*, Journal of Mathematical Imaging and Vision **40** (2011), 120–145, doi:[10.1007/s10851-010-0251-1](https://doi.org/10.1007/s10851-010-0251-1).
- [24] T. Pock, D. Cremers, H. Bischof and A. Chambolle, *An algorithm for minimizing the mumford-shah functional*, in: *12th IEEE Conference on Computer Vision*, 2009, 1133–1140, doi:[10.1109/ICCV.2009.5459348](https://doi.org/10.1109/ICCV.2009.5459348).
- [25] E. Esser, X. Zhang and T. F. Chan, *A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science*, SIAM Journal on Imaging Sciences **3** (2010), 1015–1046, doi:[10.1137/09076934X](https://doi.org/10.1137/09076934X).
- [26] D. Gabay, *Applications of the method of multipliers to variational inequalities*, in: *Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems*, volume 15, Edited by M. Fortin and R. Glowinski, North-Holland 1983, 299–331.
- [27] J. Douglas, Jim and J. Rachford, H. H., *On the numerical solution of heat conduction problems in two and three space variables*, Transactions of the American Mathematical Society **82** (1956), pp. 421–439, doi:[10.2307/1993056](https://doi.org/10.2307/1993056).
- [28] K. Bredies and H. P. Sun, *Preconditioned Douglas–Rachford algorithms for TV- and TGV-regularized variational imaging problems*, Journal of Mathematical Imaging and Vision **52** (2015), 317–344, doi:[10.1007/s10851-015-0564-1](https://doi.org/10.1007/s10851-015-0564-1).
- [29] T. Goldstein and S. Osher, *The split bregman method for l1-regularized problems*, SIAM Journal on Imaging Sciences **2** (2009), 323–343, doi:[10.1137/080725891](https://doi.org/10.1137/080725891).
- [30] W. Yin, S. Osher, D. Goldfarb and J. Darbon, *Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing*, SIAM Journal on Imaging Sciences **1** (2008), 143–168, doi:[10.1137/070703983](https://doi.org/10.1137/070703983).
- [31] S. Setzer, *Operator splittings, Bregman methods and frame shrinkage in image processing*, International Journal of Computer Vision **92** (2011), 265–280, doi:[10.1007/s11263-010-0357-3](https://doi.org/10.1007/s11263-010-0357-3).
- [32] T. Suzuki, *Stochastic dual coordinate ascent with alternating direction multiplier method* (2013), preprint, [arXiv:1311.0622](https://arxiv.org/abs/1311.0622).
- [33] Y. Zhang and L. Xiao, *Stochastic primal-dual coordinate method for regularized empirical risk minimization* (2014), [arXiv:1409.3257](https://arxiv.org/abs/1409.3257).
- [34] O. Fercoq and P. Bianchi, *A coordinate descent primal-dual algorithm with large step size and possibly non separable functions* (2015), [arXiv:1508.04625](https://arxiv.org/abs/1508.04625).
- [35] P. Bianchi, W. Hachem and F. Iutzeler, *A stochastic coordinate descent primal-dual algorithm and applications to large-scale composite optimization*, preprint, [arXiv:1407.0898](https://arxiv.org/abs/1407.0898).
- [36] Z. Peng, T. Wu, Y. Xu, M. Yan and W. Yin, *Coordinate friendly structures, algorithms and applications* (2016), [arXiv:1601.00863](https://arxiv.org/abs/1601.00863).

- [37] J.-C. Pesquet and A. Repetti, *A class of randomized primal-dual algorithms for distributed optimization* (2014), [arXiv:1406.6404](#).
- [38] A. W. Yu, Q. Lin and T. Yang, *Doubly stochastic primal-dual coordinate method for empirical risk minimization and bilinear saddle-point problem* (2015), [arXiv:1508.03390](#).
- [39] C. Chen, B. He, Y. Ye and X. Yuan, *The direct extension of ADMM for multi-block convex minimization problems is not necessarily convergent*, *Mathematical Programming* **155** (2014), 57–79, doi:[10.1007/s10107-014-0826-5](#).
- [40] A. S. Lewis and S. Zhang, *Partial smoothness, tilt stability, and generalized hessians*, *SIAM Journal on Optimization* **23** (2013), 74–94, doi:[10.1137/110852103](#).
- [41] A. S. Lewis, *Active sets, nonsmoothness, and sensitivity*, *SIAM Journal on Optimization* **13** (2002), 702–725, doi:[10.1137/S1052623401387623](#).
- [42] J. Liang, J. Fadili and G. Peyré, *Local linear convergence of forward-backward under partial smoothness*, *Advances in Neural Information Processing Systems* **27** (2014), 1970–1978.
URL <http://papers.nips.cc/paper/5260-local-linear-convergence-of-forward-backward-under-partial-smoothness.pdf>
- [43] T. Pock and A. Chambolle, *Diagonal preconditioning for first order primal-dual algorithms in convex optimization*, in: *Computer Vision (ICCV), 2011 IEEE International Conference on*, 2011, 1762–1769, doi:[10.1109/ICCV.2011.6126441](#).
- [44] B. He and X. Yuan, *Convergence analysis of primal-dual algorithms for a saddle-point problem: From contraction perspective*, *SIAM Journal on Imaging Sciences* **5** (2012), 119–149, doi:[10.1137/100814494](#).
- [45] A. N. Shiriaev, *Probability*, Graduate Texts in Mathematics, Springer 1996.
- [46] A. Chambolle and T. Pock, *On the ergodic convergence rates of a first-order primal-dual algorithm*, *Mathematical Programming* (2015), 1–35, doi:[10.1007/s10107-015-0957-3](#).
- [47] Y. Chen, G. Lan and Y. Ouyang, *Optimal primal-dual methods for a class of saddle point problems*, *SIAM Journal on Optimization* **24** (2014), 1779–1814, doi:[10.1137/130919362](#).
- [48] T. Goldstein, M. Li and X. Yuan, *Adaptive primal-dual splitting methods for statistical learning and image processing*, *Advances in Neural Information Processing Systems* **28** (2015), 2080–2088.
- [49] T. Valkonen, K. Bredies and F. Knoll, *Total generalised variation in diffusion tensor imaging*, *SIAM Journal on Imaging Sciences* **6** (2013), 487–525, doi:[10.1137/120867172](#).
URL <http://iki.fi/tuomov/mathematics/dtireg.pdf>
- [50] J. C. de Los Reyes, C.-B. Schönlieb and T. Valkonen, *Bilevel parameter learning for higher-order total variation regularisation models*, *Journal of Mathematical Imaging and Vision* (2016), doi:[10.1007/s10851-016-0662-8](#), published online, [arXiv:1508.07243](#).
URL http://iki.fi/tuomov/mathematics/tgv_learn.pdf
- [51] A. Chambolle, *An algorithm for mean curvature motion*, *Interfaces and Free Boundaries* **6** (2004), 195.

- [52] T. Hohage and C. Homann, *A generalization of the Chambolle-Pock algorithm to Banach spaces with applications to inverse problems* (2014), preprint, [arXiv:1412.0126](https://arxiv.org/abs/1412.0126).