

First-order primal-dual methods for nonsmooth non-convex optimisation

Tuomo Valkonen

Abstract We provide an overview of primal-dual algorithms for nonsmooth and non-convex-concave saddle-point problems. This flows around a new analysis of such methods, using Bregman divergences to formulate simplified conditions for convergence.

1 Introduction

Interesting imaging problems can often be written in the general form

$$\min_{x \in X} \max_{y \in Y} F(x) + K(x, y) - G_*(y), \quad (\text{S})$$

where X and Y are Banach spaces, $K \in C^1(X, Y)$, and $F : X \rightarrow \overline{\mathbb{R}}$ and $G_* : Y \rightarrow \overline{\mathbb{R}}$ are convex, proper, lower semicontinuous functions with G_* the **preconjugate** of some $G : Y^* \rightarrow \overline{\mathbb{R}}$, meaning $G = (G_*)^*$. The functions F and G_* may be nonsmooth. In this chapter, we provide an overview of proximal-type primal-dual algorithms for this class of problems together with a simplified analysis, based on Bregman divergences.

> Notation, conventions, and basic convex analysis

As is standard in optimisation, all vector/Banach/Hilbert spaces in this chapter are over the real field without it being explicitly mentioned. For basic definitions of convex analysis, such as the (pre)conjugate and the subdifferential, see the [glossary](#) at the end of the chapter or textbooks such as [37, 56, 25, 31].

Tuomo Valkonen

Center for Mathematical Modeling, Escuela Politécnica Nacional, Quito, Ecuador *and* Department of Mathematics and Statistics, University of Helsinki, Finland; e-mail: tuomo.valkonen@iki.fi

A common instance of (S) is when $K(x, y) = \langle Ax|y \rangle$ for a linear operator $A \in \mathbb{L}(X; Y^*)$ with $\langle \cdot | \cdot \rangle : Y^* \times Y \rightarrow \mathbb{R}$ denoting the dual product. Then (S) arises from writing G in terms of its (pre)conjugate G_* in

$$\min_{x \in X} F(x) + G(Ax). \quad (1)$$

We now discuss sample imaging and inverse problems of the types (S) and (1), and then outline our approach to solving them in the rest of the chapter.

1.1 Sample problems

Optimisation problems of the type (1) can effectively model linear *inverse problems*; typically one would attempt to minimise the sum of a data-term and a regulariser,

$$\min_{x \in X} \Phi(z - Tx) + G(Ax), \quad (2)$$

where

- $T : \mathbb{L}(X; \mathbb{R}^n)$ is a forward operator, mapping our unknown x into a finite number of measurements.
- Φ models noise ν in the data $z \in \mathbb{R}^n$; for normal-distributed noise, $\Phi(z) = \frac{1}{2} \|z\|^2$;
- $G \circ A$ is a typically nonsmooth regularisation term that models our prior assumptions on what a good solution to the ill-posed problem $z = Tx + \nu$ should be; in imaging, what “looks good”.

For conventional total variation regularisation on a domain $\Omega \subset \mathbb{R}^m$ one would take $G(y^*) = \alpha \|y^*\|_{\mathcal{M}(\Omega; \mathbb{R}^m)}$ the Radon norm of the measure $y^* \in \mathcal{M}(\Omega; \mathbb{R}^m)$ weighted by the regularisation parameter $\alpha > 0$, and $A = D \in \mathbb{L}(\text{BV}(\Omega); \mathcal{M}(\Omega; \mathbb{R}^m))$ the [distributional derivative](#) [1]. Simple examples of a *linear* forward operator T include:

- the identity for denoising [58],
- a convolution operation for deblurring or deconvolution [67],
- a subsampling operator for inpainting [59],
- the Fourier transform for magnetic resonance imaging (MRI) [51, 46], and
- the Radon transform for computational (CT) or positron emission tomography (PET) [52].

The last two examples would frequently be combined with subsampling for reconstruction from limited data.

In many important problems T is, however, nonlinear:

- a pointwise application of $(r, \varphi) \mapsto r e^{-i\varphi}$ for phase and amplitude reconstruction for velocity-encoded magnetic resonance imaging [62],
- a pointwise application of $u \mapsto s_0 - s e^{-\langle u, b \rangle}$ to model the Stejskal–Tanner equation in diffusion tensor imaging [62, 41], or

- the solution operator of nonlinear partial differential equation (PDE) for several forms of tomography from magnetic and electric to acoustic and optical [51, 52, 2, 42, 39, 60, 61, 44].

In the last example, the PDE governs the physics of measurement, typically relating boundary measurements and excitations to interior data. The methods we study in this chapter are applied to electrical impedance tomography in [40, 48].

How to fit a nonlinear forward operator T into the framework (S) that requires both F and G_* to be convex? If the noise model $\Phi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is convex, proper, and lower semicontinuous, we can write (2) using the Fenchel conjugate Φ^* and $K_{TA}(x, (y_1, y_2)) := \langle z - T(x)|y_1 \rangle + \langle Ax|y_2 \rangle$ as

$$\min_{x \in X} \max_{(y_1, y_2) \in \mathbb{R}^n \times Y} K_{TA}(x, (y_1, y_2)) - \Phi^*(y_1) - G_*(y_2). \quad (3)$$

This is of the form (S) for the functions $\tilde{F} \equiv 0$ and $\tilde{G}_*(y_1, y_2) := \Phi^*(y_1) - G_*(y_2)$. Even for linear T , although (2) is readily of the form (1) and hence (S), this reformulation may allow expressing (2) in the form (S) with both \tilde{F} and \tilde{G}_* “prox-simple”. We will make this concept, important for the effective realisation of algorithms, more precise in Section 3.

Finally, fully general K in (S) was shown in [24] to be useful for highly nonsmooth and nonconvex problems, such as the [34]. Indeed, the “0-function”

$$|t|_0 := \begin{cases} 0, & t = 0, \\ 1, & t \neq 0, \end{cases}$$

can be written

$$|t|_0 = \sup_{s \in \mathbb{R}} \rho(st) \quad \text{for} \quad \rho(t) = 2t - t^2.$$

For the (anisotropic) Potts model this is applied pixelwise on a discretised image gradient computed for an $n_1 \times n_2$ image by $\nabla_h : \mathbb{R}^{n_1 n_2} \rightarrow \mathbb{R}^{2 \times n_1 n_2}$ [24]:

$$\min_{x \in \mathbb{R}^{n_1 n_2}} \max_{y \in \mathbb{R}^{2 \times n_1 n_2}} \frac{1}{2} \|b - x\|_2^2 + \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \rho(\langle [\nabla_h x]_{ij}, y_{ij} \rangle), \quad (4)$$

where $b \in \mathbb{R}^{n_1 n_2}$ is the image to be segmented.

1.2 Outline

We introduce in Section 3 methods for (S) inspired by the *primal-dual proximal splitting* (PDPS) of [17, 55] for bilinear K , commonly known as the *Chambolle–Pock method*. We work in Banach spaces, as was done in [38]. To be able to define proximal-type methods in Banach spaces, in Section 2, we introduce and recall the crucial properties of so-called *Bregman divergences*.

Our main reason for working with Bregman divergences is, however, not the generality of Banach spaces. Rather, they provide a powerful proof tool to deal with the general K in (S). This approach allows us in Section 4 to significantly simplify and better explain the original convergence proofs and conditions of [17, 62, 23, 24, 48]. Without additional effort, they also allow us to present block-adapted methods like those in [66, 64, 48].

Our overall approach and the internal organisation of Section 4 centres around the following three main ingredients of the convergence proof:

- (i) A **three-point identity**, satisfied by all Bregman divergences (shown in Section 2 and employed in Section 4.1),
- (ii) **(Semi-)ellipticity** of the algorithm-defining Bregman divergences (concept defined in Section 2, specific Bregman divergence in Section 3, and its ellipticity verified in Sections 4.2 and 4.3 through several examples), and
- (iii) A **non-smooth second-order growth** condition around a solution of (S) (treated in Sections 4.4 and 4.5).

With these basic ingredients, we then prove convergence in Sections 4.6 and 4.7. In the present overview, with focus on key concepts and aiming to avoid technical complications, we only cover, weak, strong, and linear convergence of iterates, and the convergence of gap functionals when K is convex-concave.

In Section 5 we improve the basic method by adding dependencies to earlier iterates, a form of inertia. This is needed to develop an effective algorithm for K not affine in y , including the aforementioned formulation of the Potts segmentation model. We finish in Section 6 with pointers to alternative methods and further extensions.

2 Bregman divergences

The norm and inner product in a (real) Hilbert space X satisfy the three-point identity

$$\langle x - y, x - z \rangle_X = \frac{1}{2} \|x - y\|_X^2 - \frac{1}{2} \|y - z\|_X^2 + \frac{1}{2} \|x - z\|_X^2 \quad (x, y, z \in X). \quad (5)$$

This is crucial for convergence proofs of optimisation methods [63], so we would like to have something similar in Banach spaces—or other more general spaces. Towards this end, we let $J : X \rightarrow \mathbb{R}$ be a Gâteaux-differentiable function.¹ Then one can define the asymmetric *Bregman divergence*

$$B_J(z, x) := J(z) - J(x) - \langle DJ(x) | z - x \rangle_X \quad (x, z \in X). \quad (6)$$

¹ The differentiability assumption is for notational and presentational simplicity; otherwise we would need to write the Bregman divergence as $B_J^p(z, x) := J(z) - J(x) - \langle p | z - x \rangle_X$ for some subdifferential p of J , and define explicit updates of this subdifferential in algorithms.

This function is non-negative *if and only if*² the *generating function* J is convex; it is not in general a true distance, as it can happen that $B_J(x, z) = 0$ although $x \neq z$.

Writing D_1 for the Gâteaux derivative with respect to the first parameter, we have

$$D_1 B_J(x, z) = DJ(z) - DJ(x). \quad (7)$$

Moreover, the Bregman divergence satisfies for any $\bar{x} \in X$ the *three-point identity*

$$\begin{aligned} \langle D_1 B_J(x, z) | x - \bar{x} \rangle_X &= \langle DJ(x) - DJ(z) | x - \bar{x} \rangle_X \\ &= B_J(\bar{x}, x) - B_J(\bar{x}, z) + B_J(x, z). \end{aligned} \quad (8)$$

Indeed, writing the right-hand side out, we have

$$\begin{aligned} B_J(\bar{x}, x) - B_J(\bar{x}, z) + B_J(x, z) &= [J(\bar{x}) - J(x) - \langle DJ(x) | \bar{x} - x \rangle_X] \\ &\quad - [J(\bar{x}) - J(z) - \langle DJ(z) | \bar{x} - z \rangle_X] \\ &\quad + [J(x) - J(z) - \langle DJ(z) | x - z \rangle_X], \end{aligned}$$

which immediately gives the three-point identity.

Example 1 In a Hilbert space X , the *standard generating function* $J = N_X := \frac{1}{2} \|\cdot\|_X^2$ yields $B_J(z, x) = \frac{1}{2} \|z - x\|_X^2$, so (8) recovers (5).

We will frequently require B_J to be *non-negative* or *semi-elliptic* ($\gamma = 0$) or *elliptic* ($\gamma > 0$) within some $\Omega \subset X$. These notions mean that

$$B_J(z, x) \geq \frac{\gamma}{2} \|z - x\|_X^2 \quad (x, z \in \Omega). \quad (9)$$

Equivalently, this defines J to be (γ -strongly) *subdifferentiable* within Ω . When $\Omega = X$, we simply call B_J (semi-)elliptic and J (γ -strongly) subdifferentiable.³

We will in [Section 5](#) also need a Cauchy inequality for Bregman divergences. We base this on strong subdifferentiability and the smoothness property (10) in the next lemma. The latter holding with $\Omega = X$ implies that DJ is L -Lipschitz, and in Hilbert spaces is equivalent to this property; see [5, Theorem 18.15] or [63, Appendix C].

Lemma 1 *Suppose $J : X \rightarrow \mathbb{R}$ is Gâteaux-differentiable and γ -strongly subdifferentiable within Ω , and satisfies for some $L > 0$ the subdifferential smoothness*

$$\frac{1}{2L} \|DJ(x) - DJ(y)\|_{X^*}^2 \leq J(x) - J(y) - \langle DJ(y) | x - y \rangle \quad (x, y \in \Omega). \quad (10)$$

Then, for any $\alpha > 0$,

$$|\langle D_1 B_J(x, y) | z - x \rangle| \leq \frac{L}{\alpha} B_J(x, y) + \frac{\alpha}{\gamma} B_J(z, x) \quad (x, y, z \in \Omega).$$

² For the entirely algebraic proof of the “only if”, see [37, Theorem 4.1.1].

³ In Banach spaces strong subdifferentiability is implied by strong convexity, as defined without subdifferentials. In Hilbert spaces the two properties are equivalent.

Proof By Cauchy's inequality and (7),

$$|\langle D_1 B_J(x, y) | z - x \rangle| \leq \frac{1}{2\alpha} \|DJ(x) - DJ(y)\|_{X^*}^2 + \frac{\alpha}{2} \|z - x\|_X^2.$$

By the strong convexity, $\frac{\gamma}{2} \|z - x\|_X^2 \leq B_J(z, x)$, and by the smoothness property (10), $\frac{1}{2L} \|DJ(x) - DJ(y)\|_{X^*}^2 \leq B_J(x, y)$. Together these estimates yield the claim. \square

3 Primal-dual proximal splitting

We now formulate a basic version of our primal-dual method. Later in Section 5 we improve the algorithm to be more effective when K is not affine in y .

> Notation

Throughout the manuscript, we combine the primal and dual variables x and y into variables involving the letter u :

$$u = (x, y), \quad u^k = (x^k, y^k), \quad \hat{u} = (\hat{x}, \hat{y}), \quad \text{etc.}$$

3.1 Optimality conditions and proximal points

We define the *Lagrangian* as

$$\mathcal{L}(x, y) := F(x) + K(x, y) - G_*(y).$$

A *saddle point* $\hat{u} = (\hat{x}, \hat{y})$ of the problem (S) satisfies, by definition

$$\mathcal{L}(\hat{x}, y) \leq \mathcal{L}(\hat{x}, \hat{y}) \leq \mathcal{L}(x, \hat{y}) \quad \text{for all } u = (x, y) \in X \times Y.$$

Writing $D_x K$ and $D_y K$ for the Gâteaux derivatives of K with respect to the two variables, if K is convex-concave, basic results in convex analysis [31, 5] show that

$$-D_x K(\hat{x}, \hat{y}) \in \partial F(\hat{x}) \quad \text{and} \quad D_y K(\hat{x}, \hat{y}) \in \partial G_*(\hat{y}) \quad (11)$$

is necessary and sufficient for \hat{u} to be saddle point. If K is C^1 , the theory of generalised subdifferentials of Clarke [22] still indicates⁴ the necessity of (11).

We can alternatively write (11) as

⁴ The Fermat-rule $0 \in \partial_C [F + K(\cdot, \hat{y})](\hat{x})$ holds. Since F is convex and $K(\cdot, \hat{y})$ is C^1 , \hat{x} is a regular point of both, so also the subdifferential sum rule holds. We argue $G_* + K(\hat{y}, \cdot)$ similarly.

$$0 \in H(\hat{u}) := \left(\begin{array}{l} \partial F(\hat{x}) + D_x K(\hat{x}, \hat{y}) \\ \partial G_*(\hat{y}) - D_y K(\hat{x}, \hat{y}) \end{array} \right). \quad (12)$$

If X and Y were Hilbert spaces, we could in principle use the classical *proximal point method* [49, 57] to solve (12): given step length parameters $\tau_k > 0$, iteratively solve u^{k+1} from

$$0 \in H(u^{k+1}) + \tau_k^{-1}(u^{k+1} - u^k). \quad (13)$$

If K were bilinear, H would be a so-called monotone operator and convergence of iterates would follow from [57]. In practise the steps of the method are too expensive to realise as the primal and dual iterates x^{k+1} and y^{k+1} are coupled: generally, one cannot solve one before the other.

Fortunately, the iterates can be decoupled by introducing a *preconditioner* that switches $D_x K(x^{k+1}, y^{k+1})$ on the first line of $H(u^{k+1})$ to $D_x K(x^k, y^k)$. This gives rise to the *primal-dual proximal splitting* (PDPS), introduced in [17, 55] for bilinear $K(x, y) = \langle Ax|y \rangle$. That the PDPS is actually a preconditioned proximal point method was first observed in [36]. In the following, we describe its extension from [62, 23, 24] to general K and the general problem (S). To simplify the proofs and concepts in them, we work with Bregman divergences, at no cost in Banach spaces.

3.2 Algorithm formulation

Given Gâteaux-differentiable functions $J_X : X \rightarrow \overline{\mathbb{R}}$ and $J_Y : Y \rightarrow \overline{\mathbb{R}}$ with the corresponding Bregman divergences $B_X := B_{J_X}$ and $B_Y := B_{J_Y}$, we define

$$J^0(x, y) := J_X(x) + J_Y(y) - K(x, y). \quad (14)$$

Introducing the short-hand notation $B^0 := B_{J^0}$, we propose to solve (12) through the iterative solution of

$$0 \in H(u^{k+1}) + D_1 B^0(u^{k+1}, u^k) \quad (15)$$

for u^{k+1} . Inserting (12) and (7) for $J = J^0$ as defined in (14), we expand and rearrange this implicitly defined method as:

Primal-dual Bregman-proximal splitting (PDBS)

Iteratively over $k \in \mathbb{N}$, solve for x^{k+1} and y^{k+1} :

$$\begin{aligned} DJ_X(x^k) - D_x K(x^k, y^k) &\in DJ_X(x^{k+1}) + \partial F(x^{k+1}) \quad \text{and} \\ DJ_Y(y^k) - D_y K(x^k, y^k) &\in DJ_Y(y^{k+1}) + \partial G_*(y^{k+1}) - 2D_y K(x^{k+1}, y^{k+1}). \end{aligned} \quad (16)$$

We readily obtain x^{k+1} if the inverse of $DJ_X + \tau \partial F$ has an analytical closed-form expression. In this case we say that F is *prox-simple* with respect to J_X . For y^{k+1} , the same is true if K is affine in y and G_* is prox-simple with respect to J_Y . If, however, K

is not affine in y , it is practically unlikely that $\partial G_* - 2D_y K(x^{k+1}, \cdot)$ would be prox-simple. We will therefore improve the method for general K in Section 5, after first studying fundamental ideas behind convergence proofs in the following Section 4.

If X and Y are Hilbert spaces with $J_X = \tau^{-1}N_X$ and $J_Y = \sigma^{-1}N_Y$ the standard generating functions divided by some step length parameters $\tau, \sigma > 0$, (16) becomes

Primal–dual proximal splitting (PDPS)

Iterate over $k \in \mathbb{N}$:

$$\begin{aligned} x^{k+1} &:= \text{prox}_{\tau F}(x^k - \tau \nabla_x K(x^k, y^k)), \\ y^{k+1} &:= \text{prox}_{\sigma[G_* - 2K(x^{k+1}, \cdot)]}(y^k - \sigma \nabla_y K(x^k, y^k)). \end{aligned} \quad (17)$$

The *proximal map* is defined as

$$\text{prox}_{\tau F}(x) := (I + \tau \partial F)^{-1}(x) = \arg \min_{\tilde{x} \in X} \left(\tau F(\tilde{x}) + \frac{1}{2} \|\tilde{x} - x\|_X^2 \right).$$

When this map has an analytical closed-form expression, we say that F is *prox-simple* (without reference to J_X). In finite dimensions, several worked out proximal maps may be found online [21] or in the book [6]. Some extend directly to Hilbert spaces or by superposition to L^2 .

Remark 1 For K affine in y , i.e., $K(x, y) = \langle A(x)|y \rangle$ for some differentiable $A : X \rightarrow Y^*$, the dual update of (17) reduces to

$$\begin{aligned} y^{k+1} &= \text{prox}_{\sigma G_*}(y^k + \sigma[2\nabla_y K(x^{k+1}, y^k) - \nabla_y K(x^k, y^k)]) \\ &= \text{prox}_{\sigma G_*}(y^k + \sigma[2\nabla A(x^{k+1}) - \nabla A(x^k)]). \end{aligned}$$

This corresponds to the “linearised” variant of the NL-PDPS of [62]. The “exact” variant, studied in further detail in [23], updates

$$y^{k+1} := \text{prox}_{\sigma G_*}(y^k + \sigma \nabla_y K(2x^{k+1} - x^k, y^k)).$$

If K is bilinear the two variants are the exactly same PDPS of [17]. For K not affine in y , the method is neither the generalised PDPS of [24] nor the version for convex-concave K from [35].

3.3 Block-adaptation

We now derive a version of the PDPS (16) adapted to the structure of

$$F(x) = \sum_{j=1}^m F_j(x_j) \quad \text{and} \quad G_*(y) = \sum_{\ell=1}^n G_{\ell^*}(y_\ell),$$

where $x = (x_1, \dots, x_m)$ and $y = (y_1, \dots, y_n)$ in the (for simplicity) Hilbert spaces $X = \prod_{j=1}^m X_j$ and $Y = \prod_{\ell=1}^n Y_\ell$, and $F_j : X_j \rightarrow \overline{\mathbb{R}}$ and $G_{\ell^*} : Y_\ell \rightarrow \overline{\mathbb{R}}$ are convex, proper, and lower semicontinuous.

For some “blockwise” step length parameters $\tau_j, \sigma_\ell > 0$ we take

$$J_X(x) = \sum_{j=1}^m \tau_j^{-1} N_{X_j}(x_j) \quad \text{and} \quad J_Y(y) = \sum_{\ell=1}^n \sigma_\ell^{-1} N_{Y_\ell}(y_\ell)$$

If K is now affine in y , observing [Remark 1](#), (16) readily transforms into:

Block-adapted PDPS for K affine in y

Iteratively over $k \in \mathbb{N}$, for all $j = 1, \dots, m$ and $\ell = 1, \dots, n$, update:

$$\begin{aligned} x_j^{k+1} &:= \text{prox}_{\tau_j F_j}(x_j^k - \tau_j \nabla_{x_j} K(x^k, y^k)), \\ y_\ell^{k+1} &:= \text{prox}_{\sigma_\ell G_{\ell^*}}(y_\ell^k + \sigma_\ell [2\nabla_{y_\ell} K(x^{k+1}, y^k) - \nabla_{y_\ell} K(x^k, y^k)]). \end{aligned} \quad (18)$$

The idea is that the blockwise step length parameters adapt the algorithm to the structure of the problem. We will return their choices in the examples of [Section 4.3](#).

> Performance gains

Correct adaptation of the blockwise step length parameters to the specific problem structure can yield significant performance gains compared to not exploiting the block structure [[54](#), [40](#), [48](#)].

Remark 2 For bilinear K , (18) is the “diagonally preconditioned” method of [[54](#)], or an unaccelerated non-stochastic variant of the methods in [[64](#)]. For K affine in y , (18) differs from the methods in [[48](#)] by placing the over-relaxation in the dual step outside K , compare [Remark 1](#).

Recall the saddle-point formulation (3) for inverse problems with nonlinear forward operators. We can now adapt step lengths to the constituent dual blocks:

Example 2 Let $A_1 \in C^1(X; Y_1^*)$ and $A_2 \in \mathbb{L}(X; Y_2^*)$, and suppose the convex functions $G_1 : Y_1^* \rightarrow \overline{\mathbb{R}}$ and $G_2 : Y_2^* \rightarrow \overline{\mathbb{R}}$ have the preconjuguates G_{1^*} and G_{2^*} . Then we can write the problem

$$\min_{x \in X} G_1(A_1(x)) + G_2(A_2 x) + F(x).$$

in the form (S) with $G_*(y_1, y_2) = G_{1*}(y_1) + G_{2*}(y_2)$ and $K(x, y) = \langle A_1(x)|y_1 \rangle + \langle A_2x|y_2 \rangle$. The algorithm (18) specialises as

$$\begin{aligned} x^{k+1} &:= \text{prox}_{\tau F}(x^k - \tau[\nabla A_1(x^k)^* y_1 + A_2^* y_2]), \\ y_1^{k+1} &:= \text{prox}_{\sigma_1 G_{1*}}(y_1^k + \sigma_1[2A_1(x^{k+1}) - A_1(x^k)]), \\ y_2^{k+1} &:= \text{prox}_{\sigma_2 G_{2*}}(y_2^k + \sigma_2[A_2(2x^{k+1} - x^k)]) \end{aligned}$$

for some step length parameters $\tau, \sigma_1, \sigma_2 > 0$. We return to their choices and the local neighbourhood of convergence in Examples 8 and 17 after developing the necessary convergence theory.

4 Convergence theory

We now seek to understand when the basic version (15) of the PDBS convergences. The organisation of this section centres around the **three main ingredients** of the convergence proof, as discussed in the Introduction:

- (i) the three-point identity (8) employed in the general-purpose estimate of Section 4.1,
- (ii) the (semi-)ellipticity of the algorithm-generating Bregman divergences B_{J_0} for J^0 as in (14), verified for several examples in Sections 4.2 and 4.3, and
- (iii) a second-order growth condition on (S), verified for several examples in Sections 4.4 and 4.5.

With these basic ingredients, we then prove various convergence results in Sections 4.6 and 4.7. The usefulness of both (ii) and (iii) will become apparent from the fundamental estimates and examples of the next Section 4.1.

4.1 A fundamental estimate

We start with a simple estimate applicable to general methods of the form

$$0 \in H(u^{k+1}) + D_1 B(u^{k+1}, u^k) \tag{BP}$$

for some set-valued $H : U \rightrightarrows U^*$ and a Bregman divergence $B := B_J$ generated by some Gâteaux-differentiable $J : U \rightarrow \mathbb{R}$. We analyse (BP) following the “testing” ideas introduced in [63], extending them to the Bregman–Banach space setting, however in a simplified constant-metric setting that cannot model accelerated methods. The *generic gap functional* $\mathcal{G}(u^{k+1}, \bar{u})$ in the next result models any function value differences available from H . Its non-negativity will provide the basis for the aforementioned second-order growth conditions of Sections 4.4 and 4.5. We provide an example and interpretation after the theorem.

Theorem 1 *On a Banach space U , let $H : U \rightrightarrows U^*$, and let $B := B_J$ be generated by a Gâteaux-differentiable $J : U \rightarrow \mathbb{R}$. Suppose (BP) is solvable for $\{u^{k+1}\}_{k \in \mathbb{N}}$ given an initial iterate $u^0 \in U$. Let $N \geq 1$. If for all $k = 0, \dots, N-1$, for some $\bar{u} \in U$ and $\mathcal{G}(u^{k+1}, \bar{u}) \in \mathbb{R}$ the fundamental condition*

$$\langle h^{k+1} | u^{k+1} - \bar{u} \rangle \geq \mathcal{G}(u^{k+1}, \bar{u}) \quad (h^{k+1} \in H(u^{k+1})) \quad (\text{C})$$

holds, then so do the quantitative Δ -Féjer monotonicity

$$B(\bar{u}, u^{k+1}) + B(u^{k+1}, u^k) + \mathcal{G}(u^{k+1}, \bar{u}) \leq B(\bar{u}, u^k) \quad (\text{F})$$

and the descent inequality

$$B(\bar{u}, u^N) + \sum_{k=0}^{N-1} B(u^{k+1}, u^k) + \sum_{k=0}^{N-1} \mathcal{G}(u^{k+1}, \bar{u}) \leq B(\bar{u}, u^0). \quad (\text{D})$$

Proof We can write (BP) as

$$0 = h^{k+1} + D_1 B(u^{k+1}, u^k) \quad \text{for some } h^{k+1} \in H(u^{k+1}). \quad (\text{19})$$

Testing (19) by applying $\langle \cdot | u^{k+1} - \bar{u} \rangle$ we obtain

$$0 = \langle h^{k+1} + D_1 B(u^{k+1}, u^k) | u^{k+1} - \bar{u} \rangle.$$

We use the three-point identity (8) to transform this into

$$B(\bar{u}, u^k) = \langle h^{k+1} | u^{k+1} - \bar{u} \rangle + B(\bar{u}, u^{k+1}) + B(u^{k+1}, u^k).$$

Inserting (C), we obtain (F). Summing the latter over $k = 0, \dots, N-1$ yields (D). \square

Example 3 If $H = \partial F$ for a convex function F , then by the definition of the convex subdifferential, (C) holds with the gap functional

$$\mathcal{G}(u, \bar{u}) = F(u) - F(\bar{u}).$$

If we take \bar{u} is a minimiser of F , then the gap functional is non-negative and indeed positive if u is also not minimiser. This is why it is called a gap functional.

Consider then for some step length parameter $\tau > 0$ the proximal point method (13) in a Hilbert space X , that is, taking $B = \tau^{-1}N_X$,

$$u^{k+1} := \text{prox}_{\tau F}(x^k), \quad \text{equivalently } 0 \in \partial F(u^{k+1}) + \tau(u^{k+1} - u^k).$$

Then (D) reads

$$\frac{1}{2\tau} \|u^N - \bar{u}\|_X^2 + \sum_{k=0}^{N-1} \frac{1}{2} \|u^{k+1} - u^k\|_X^2 + \sum_{k=0}^{N-1} \tau(F(u^{k+1}) - F(\bar{u})) \leq \frac{1}{2} \|\bar{u} - u^0\|_X^2. \quad (\text{20})$$

With \bar{u} a minimiser, this clearly forces $F(u^N) \rightarrow F(\bar{u})$ as $N \rightarrow \infty$, suggesting why we call (D) the “descent inequality”.

If our problem is non-convex, then we can try to locally ensure second-order growth by imposing $\mathcal{G}(u^{k+1}, \bar{u}) \geq 0$. Verifying this for the PDBS will be the topic of Sections 4.4 and 4.5. If B is not given by the standard generating function N_X on a Hilbert spaces X , then to get from (D) an estimate like (20) on norms, we can assume the ellipticity or at least semi-ellipticity of the overall Bregman divergence B . Verifying this for $B = B_{J^0}$ with J^0 given in (14) is our next topic.

4.2 Ellipticity of the Bregman divergences

As just discussed, for Theorem 1 to provide estimates that we can use to prove the convergence of the PDBS, we need at least the semi-ellipticity of B^0 generated by J^0 given in (14). Deriving simple conditions that ensure such semi-ellipticity or ellipticity is the topic of the present subsection. To do this, we need the “basic” Bregman divergences B_X and B_Y on both spaces X and Y to be elliptic:

➤ Standing assumption

In this subsection, we assume that B_X is τ^{-1} -elliptic and B_Y is σ^{-1} -elliptic for some $\tau, \sigma > 0$. This is true for the Hilbert-space PDPS (17) where τ and σ are the primal and dual step length parameters.

The examples that follow the next general lemma will provide improved estimates.

Lemma 2 *Suppose $K \in C^1(X \times Y)$ is Lipschitz-continuously differentiable with the factor L_{DK} in a convex subdomain $\Omega \subset X \times Y$. Then for $u, u' \in \Omega$,*

$$B_K(u', u) \leq \frac{L_{DK}}{2} \|u' - u\|_{X \times Y}^2. \quad (21)$$

Consequently, if B_X is τ^{-1} -elliptic and B_Y is σ^{-1} -elliptic and $1 \geq \max\{\tau, \sigma\}L_{DK}$, then B^0 is semi-elliptic (elliptic if the inequality is strict) within Ω .

Proof By definition, $B_K(u', u) = K(u') - K(u) - \langle DK(u) | u' - u \rangle$. Using the mean value equality in \mathbb{R} with the chain rule and the Cauchy–Schwarz inequality, we get

$$B_K(u', u) = \int_0^1 \langle DK(u + t(u' - u)) - DK(u) | u' - u \rangle dt \leq \int_0^1 t L_{DK} \|u' - u\|_{X \times Y}^2 dt.$$

Calculating the last integral yields (21).

For the (semi-)ellipticity, we need $B^0(u, u') \geq \frac{\varepsilon}{2} \|u - u'\|_{X \times Y}^2$ for some $\varepsilon > 0$ ($\varepsilon = 0$) and all $u, u' \in \Omega$. Since B_X and B_Y are τ^{-1} - and σ^{-1} -elliptic, we have

$$\begin{aligned} B^0(u', u) &= B_X(x', x) + B_Y(y', y) - B_K(u', u) \\ &\geq \frac{1}{2\tau} \|x' - x\|_X^2 + \frac{1}{2\sigma} \|y' - y\|_Y^2 - B_K(u', u). \end{aligned} \quad (22)$$

Using (21), therefore $B^0(u', u) \geq \frac{\tau^{-1} - L_{DK}}{2} \|x' - x\|_X^2 + \frac{\sigma^{-1} - L_{DK}}{2} \|y' - y\|_Y^2$. Thus B^0 is ε -elliptic when $\tau^{-1}, \sigma^{-1} \geq L_{DK} + \varepsilon$. This gives the claim. \square

We now provide several examples of ellipticity. In practise, to guarantee ellipticity, we would choose $\tau, \sigma > 0$ to satisfy the stated conditions.

Example 4 Suppose $K(x, y) = E(x)$ with DE L_{DE} -Lipschitz in $\Omega = X \times Y$. Then $L_{DK} = L_{DE}$, so we recover the standard-for-gradient-descent step length bound $1 \geq \tau L_{DE}$ for B^0 to be semi-elliptic in Ω (elliptic if the inequality is strict).

Example 5 If $K(x, y) = \langle Ax|y \rangle$ for $A \in \mathbb{L}(X; Y^*)$, then B^0 is elliptic under the standard-for-PDPS [17] step length condition

$$1 > \tau\sigma \|A\|^2.$$

Indeed,

$$\langle DK(u + t(u' - u)) - DK(u)|u' - u \rangle = 2t \langle A(x - x')|y - y' \rangle.$$

Therefore, taking any $w > 1$, we easily improve (21) to

$$\begin{aligned} B_K(u', u) &\leq \|A\| \|x' - x\|_X \|y' - y\|_Y \\ &\leq \frac{w\|A\|}{2} \|x' - x\|_X^2 + \frac{w^{-1}\|A\|}{2} \|y' - y\|_Y^2 \quad (u, u' \in X \times Y). \end{aligned} \quad (23)$$

By (22), B^0 is therefore ε -elliptic if $\tau^{-1} \geq w\|A\| + \varepsilon$ and $\sigma^{-1} \geq w^{-1}\|A\| + \varepsilon$. Taking $w = \sigma\|A\|/(1 - \sigma\varepsilon)$ this holds if $1 \geq \tau\sigma\|A\|^2/(1 - \sigma\varepsilon) + \tau\varepsilon$. Since $\varepsilon > 0$ was arbitrary, the claimed step length condition follows.

Example 6 Suppose $K(x, y) = \langle A(x)|y \rangle$ with A and DA Lipschitz with the respective factors $L_A, L_{DA} \geq 0$. Then B^0 is elliptic within $\Omega = X \times B(0, \rho_y)$ if

$$1 > \tau\sigma L_A^2 + \tau \frac{L_{DA}\rho_y}{2}.$$

Indeed, for any $w > 1$, using the mean value equality as in the proof of Lemma 2, we deduce

$$\begin{aligned} B_K(u', u) &= \langle A(x') - A(x)|y' \rangle - \langle DA(x)(x' - x)|y \rangle \\ &= \langle A(x') - A(x)|y' - y \rangle + \langle A(x') - A(x) - DA(x)(x' - x)|y \rangle \\ &\leq L_A \|x' - x\|_X \|y' - y\|_Y + \frac{L_{DA}\|y'\|}{2} \|x' - x\|_X^2 \\ &\leq \frac{wL_A + L_{DA}\|y'\|}{2} \|x' - x\|_X^2 + \frac{w^{-1}L_A}{2} \|y' - y\|_Y^2. \end{aligned} \quad (24)$$

If $\rho_y > 0$ is such that $\|y\| \leq \rho_y$, taking $w = \sigma L_A / (1 - \sigma\varepsilon)$, similarly to [Example 5](#) we deduce the claimed bound.

We can combine the examples above:

Example 7 As in [Example 2](#), take $K(x, (y_1, y_2)) = \langle A_1(x)|y_1 \rangle + \langle A_2x|y_2 \rangle$ with $A_1 \in C^1(X; Y_1^*)$ and $A_2 \in \mathbb{L}(X; Y_2^*)$. Then B^0 is elliptic within $\Omega = X \times B(0, \rho_y)$ if

$$1 > \tau\sigma(L_{A_1}^2 + \|A_2\|^2) + \tau \frac{L_{DA_1}\rho_{y_1}}{2}.$$

Indeed, we bound B_K by summing [\(23\)](#) for A_1 and [\(24\)](#) for A_2 . This yields for any $w_1, w_2 > 0$ the estimate

$$\begin{aligned} B_K(u', u) \leq & \frac{w_1 L_{A_1} + L_{DA_1} \|y_1\|}{2} \|x - x'\|_X^2 + \frac{w_1^{-1} L_{A_1}}{2} \|y'_1 - y_1\|_Y^2 \\ & + \frac{w_2 \|A_2\|}{2} \|x' - x\|_X^2 + \frac{w_2^{-1} \|A_2\|}{2} \|y'_2 - y_2\|_{Y_2}^2. \end{aligned} \quad (25)$$

Taking $w_1 = \sigma L_{A_1} / (1 - \sigma\varepsilon)$ and $w_2 = \sigma \|A_2\| / (1 - \sigma\varepsilon)$, and using [\(22\)](#), we deduce the claimed ellipticity for small enough $\varepsilon > 0$.

Remark 3 In [Examples 6](#) and [7](#) we needed a bound on the dual variable y . In the latter, as an improvement, this was only needed on the subspace Y_1 of non-bilinearity. An ad-hoc solution is to introduce the bound into the problem. In the Hilbert case, [\[23, 24\]](#) secure such bounds by taking the primal step length τ small enough and arguing as in [Theorem 1](#) individually on the primal and dual iterates.

4.3 Ellipticity for block-adapted methods

We now study ellipticity for block-adapted methods. The goal is to obtain faster convergence by adapting the blockwise step length parameters to the problem structure (connections between blocks) and the local (blockwise) properties of the problem.

> Standing assumption

In this subsection, we assume F , G_* , J_X and J_Y to have the form of [Section 3.3](#). In particular, X and Y are (products of) Hilbert spaces, and

$$B^0(u', u) = \sum_{j=1}^m \frac{1}{2\tau_j} \|x'_j - x_j\|_{X_j}^2 + \sum_{\ell=1}^n \frac{1}{2\sigma_\ell} \|y'_\ell - y_\ell\|_{Y_\ell}^2 - B_K(u', u). \quad (26)$$

We start by refining the two-block [Example 7](#) to be adapted to the blocks:

Example 8 Let $K(x, (y_1, y_2)) = \langle A_1(x)|y_1 \rangle + \langle A_2x|y_2 \rangle$ with $A_1 \in C^1(X; Y_1^*)$ and $A_2 \in \mathbb{L}(X; Y_2^*)$ as in [Examples 2](#) and [7](#). Write $\tau = \tau_1$. Using [\(25\)](#) in [\(26\)](#) for $m = 1$ and $n = 2$ with [\(25\)](#), we see B^0 to be ε -elliptic within $\Omega = X \times B(0, \rho_{y_1}) \times Y_2$ if $\tau^{-1} \geq w_1 L_{A_1} + L_{DA_1} \rho_{y_1} + w_2 \|A_2\| + \varepsilon$ and $\sigma_1^{-1} \geq w_1^{-1} L_{A_1}$ as well as $\sigma_2^{-1} \geq w_2^{-1} \|A_2\| + \varepsilon$. Taking $w_1 = \sigma_1 L_{A_1} / (1 - \sigma_1 \varepsilon)$ and $w_2 = \sigma_2 \|A_2\| / (1 - \sigma_2 \varepsilon)$, B^0 is therefore elliptic (some $\varepsilon > 0$) within Ω if $1 > \tau(\sigma_1 L_{A_1}^2 + \sigma_2 \|A_2\|^2) + \tau \frac{L_{DA_1} \rho_{y_1}}{2}$.

Example 9 In [Example 8](#), if both $A_1 \in \mathbb{L}(X; Y_1^*)$ and $A_2 \in \mathbb{L}(X; Y_2^*)$, then B^0 is elliptic within $\Omega = X \times Y_1 \times Y_2$ if $1 > \tau(\sigma_1 \|A_1\|^2 + \sigma_2 \|A_2\|^2)$.

Example 10 Suppose we can write $K(x, y) = \sum_{j=1}^m \sum_{\ell=1}^n K_{j\ell}(x_j, y_\ell)$ with each $K_{j\ell}$ Lipschitz-continuously differentiable with the factor $L_{j\ell}$. Following [Lemma 2](#),

$$B_K(u', u) \leq \sum_{j=1}^m \sum_{\ell=1}^n \frac{L_{j\ell}}{2} (\|x'_j - x_j\|^2 + \|y'_\ell + y_\ell\|^2). \quad (27)$$

Consequently, using [\(26\)](#), we see that B^0 is ε -elliptic if $1 \geq \tau_j(\sum_{\ell=1}^n L_{j\ell} + \varepsilon)$ and $1 \geq \sigma_\ell(\sum_{j=1}^m L_{j\ell} + \varepsilon)$ for all $j = 1, \dots, m$ and $\ell = 1, \dots, n$.

Example 11 If $K(x, y) = \sum_{j=1}^m \sum_{\ell=1}^n \langle A_{j\ell} x_j | y_\ell \rangle$ for some $A_{j\ell} \in \mathbb{L}(X_j; Y_\ell^*)$, then following [Example 5](#), for arbitrary $w_{j\ell} > 0$,

$$\begin{aligned} B_K(u', u) &\leq \sum_{j=1}^m \sum_{\ell=1}^n \|A_{j\ell}\| \|x'_j - x_j\| \|y'_\ell - y_\ell\| \\ &\leq \sum_{j=1}^m \sum_{\ell=1}^n \left(\frac{w_{j\ell} \|A_{j\ell}\|}{2} \|x'_j - x_j\|^2 + \frac{w_{j\ell}^{-1} \|A_{j\ell}\|}{2} \|y'_\ell - y_\ell\|^2 \right). \end{aligned}$$

Using [\(26\)](#), B^0 is thus ε -elliptic if $1 \geq \tau_j(\varepsilon + \sum_{\ell=1}^n w_{j\ell} \|A_{j\ell}\|)$ and $1 \geq \sigma_\ell(\varepsilon + \sum_{j=1}^m w_{j\ell}^{-1} \|A_{j\ell}\|)$ for all $j = 1, \dots, m$ and $\ell = 1, \dots, n$. We can use the factors $w_{j\ell}$ to adapt the algorithm to the different blocks for potentially better convergence.

4.4 Non-smooth second-order conditions

We now study conditions for [\(C\)](#) to hold with $\mathcal{G}(\cdot, \bar{u}) \geq 0$. We start by writing out the condition for the PDBS.

Lemma 3 *Let $\bar{u} = (\bar{x}, \bar{y}) \in X \times Y$ and suppose for some $\mathcal{G}(u, \bar{u}) \in \mathbb{R}$ and a neighbourhood $\Omega_{\bar{u}} \subset X \times Y$ that for all $u = (x, y) \in \Omega_{\bar{u}}$, $x^* \in \partial F(x)$, and $y^* \in \partial G_*(y)$,*

$$\langle x^* + D_x K(x, y) | x - \bar{x} \rangle + \langle y^* - D_y K(x, y) | y - \bar{y} \rangle \geq \mathcal{G}(u, \bar{u}). \quad (\text{C}^2)$$

Let $\{u^{k+1}\}_{k \in \mathbb{N}}$ be generated by the PDBS (16) for some $u^0 \in X \times Y$, and suppose $\{u^k\}_{k \in \mathbb{N}} \subset \Omega_{\bar{u}}$. Then with $B = B^0$ the fundamental condition (C) and the quantitative Δ -Féjer monotonicity (F) hold for all $k \in \mathbb{N}$, and the descent inequality (D) holds for all $N \geq 1$.

Proof Theorem 1 proves (F) and (D) if we show (C²). For H in (12), we have

$$h^{k+1} = \begin{pmatrix} x_{k+1}^* + D_x K(x^{k+1}, y^{k+1}) \\ y_{k+1}^* - D_y K(x^{k+1}, y^{k+1}) \end{pmatrix} \in H(u^{k+1}) \quad \text{with} \quad \begin{cases} x_{k+1}^* \in \partial F(x^{k+1}), \\ y_{k+1}^* \in \partial G_*(y^{k+1}). \end{cases}$$

Thus (C) expands as (C²) for $u = u^{k+1}$ and $(x^*, y^*) = (x_{k+1}^*, y_{k+1}^*)$. \square

In Section 4.7 on the convergence of gap functionals, we will consider general \bar{u} in (C²). For the moment, we however fix a root $\bar{u} = \hat{u} \in H^{-1}(0)$. Then

$$0 = \begin{pmatrix} \hat{x}^* + D_x K(\hat{x}, \hat{y}) \\ \hat{y}^* - D_y K(\hat{x}, \hat{y}) \end{pmatrix} \in H(\hat{u}) \quad \text{with} \quad \begin{cases} \hat{x}^* \in \partial F(\hat{x}), \\ \hat{y}^* \in \partial G_*(\hat{y}). \end{cases} \quad (28)$$

Since we assume F and G_* to be convex, their subdifferentials are monotone. When K is not convex-concave, and to obtain strong convergence of iterates even when it is, we will need some strong monotonicity of the subdifferentials, but only at a solution. Specifically, for $\gamma > 0$, we say that $T : X \rightrightarrows X^*$ is γ -strongly monotone at \hat{x} for $\hat{x}^* \in T(\hat{x})$ if

$$\langle x^* - \hat{x}^* | x - \hat{x} \rangle \geq \gamma \|x - \hat{x}\|_X^2 \quad (x \in X, x^* \in T(x)). \quad (29)$$

If $\gamma = 0$, we drop the word “strong”. For $T = \partial F$, (29) follows from the γ -strong subdifferentiability of F .

➤ Standing assumption

Throughout the rest of this subsection, we assume (28) to hold and that ∂F is (γ_F -strongly) monotone at \hat{x} for \hat{x}^* , and ∂G_* is (γ_{G_*} -strongly) monotone at \hat{y} for \hat{y}^* .

Lemma 4 *The nonsmooth second-order growth condition (C²) holds provided*

$$\gamma_F \|x - \hat{x}\|^2 + \gamma_{G_*} \|y - \hat{y}\|^2 \geq B_K(\hat{u}, u) + B_K(u, \hat{u}) + \mathcal{G}(u, \hat{u}) \quad (u \in \Omega_{\bar{u}}), \quad (30)$$

equivalently

$$\gamma_F \|x - \hat{x}\|^2 + \gamma_{G_*} \|y - \hat{y}\|^2 \geq a_K(\hat{u}, u) + a_K(u, \hat{u}) + \mathcal{G}(u, \hat{u}) \quad (u \in \Omega_{\bar{u}}) \quad (30')$$

for

$$a_K(u, \bar{u}) := K(x, y) - K(\bar{x}, \bar{y}) + \langle D_x K(x, y) | \bar{x} - x \rangle + \langle D_y K(\bar{x}, \bar{y}) | \bar{y} - y \rangle. \quad (31)$$

Note that (30) involves the *symmetrised Bregman divergence* $B_K^S(u, u') := B_K(u, u') + B_K(u', u)$ generated by K .

Proof Inserting the zero of (28) in (C²), we rewrite the latter as

$$\begin{aligned} \langle x^* - \hat{x}^* | x - \hat{x} \rangle + \langle y^* - \hat{y}^* | y - \hat{y} \rangle &\geq \langle D_x K(x, y) - D_x K(\hat{x}, \hat{y}) | \hat{x} - x \rangle \\ &+ \langle D_y K(x, y) - D_y K(\hat{x}, \hat{y}) | y - \hat{y} \rangle + \mathcal{G}(u^{k+1}, \hat{u}). \end{aligned}$$

Using the assumed strong monotonicities, and the definitions of B_K and a_K , this is immediately seen to hold when (30) or (30') does. \square

Example 12 If K is convex-concave, the next Lemma 5 and Lemma 4 prove (C²) for

$$\mathcal{G}(u, \hat{u}) = \gamma_F \|x - \hat{x}\|^2 + \gamma_{G_*} \|y - \hat{y}\|^2 \geq 0 \quad \text{within } \Omega_{\hat{u}} = X \times Y.$$

This is in particular true for $K(x, y) = \langle Ax | y \rangle + E(x)$ with $A \in \mathbb{L}(X; Y^*)$ and $E \in C^1(X)$ convex.

Lemma 5 Suppose $K : X \times Y \rightarrow \mathbb{R}$ is Gâteaux-differentiable and convex-concave. Then $a_K(u, \bar{u}) \leq 0$ and $B_K^S(u, \bar{u}) \leq 0$ for all $u, \bar{u} \in X \times Y$.

Proof The convexity of $K(\cdot, y)$ and the concavity of $K(\bar{x}, \cdot)$ show

$$\begin{aligned} K(x, y) - K(\bar{x}, y) + \langle D_x K(x, y) | \bar{x} - x \rangle &\leq 0 \quad \text{and} \\ K(\bar{x}, y) - K(\bar{x}, \bar{y}) + \langle D_y K(\bar{x}, \bar{y}) | \bar{y} - y \rangle &\leq 0. \end{aligned}$$

Summing these two estimates proves $a_K(u, \bar{u}) \leq 0$, consequently $B_K^S(u, \bar{u}) = a_K(u, \bar{u}) + a_K(\bar{u}, u) \leq 0$. \square

Example 13 Suppose K has L_{DK} -Lipschitz derivative within $\Omega \subset X \times Y$. If $\hat{u} \in \Omega$, then by Lemma 2, $B_K(u, \hat{u}), B_K(\hat{u}, u) \leq \frac{L_{DK}}{2} \|u - \hat{u}\|_{X \times Y}^2$ for $u \in \Omega$. Thus (C²) holds by Lemma 4 with $\Omega_{\hat{u}} = \Omega$ and

$$\mathcal{G}(u, \hat{u}) = (\gamma_F - L_{DK}) \|x - \hat{x}\|^2 + (\gamma_{G_*} - L_{DK}) \|y - \hat{y}\|^2.$$

This is non-negative if $\gamma_F, \gamma_{G_*} \geq L_{DK}$.

Example 14 Let $K(x, y) = \langle A(x) | y \rangle$ for some $A \in \mathbb{L}(X; Y^*)$ such that DA is Lipschitz with the factor $L_{DA} \geq 0$. For some $\tilde{\gamma}_F, \tilde{\gamma}_{G_*} \geq 0$ and $\rho_y, \hat{\rho}_x, \alpha > 0$, let either

- (a) $\tilde{\gamma}_F \geq \frac{L_{DA}}{2} (\rho_y + \|\hat{y}\|_Y)$, $\tilde{\gamma}_{G_*} \geq 0$, and $\Omega_{\hat{u}} = X \times B(0, \rho_y)$; or
- (b) $\tilde{\gamma}_F > L_{DA} (\|\hat{y}\|_Y + \frac{\alpha}{2})$, $\tilde{\gamma}_{G_*} \geq \frac{L_{DA}}{2\alpha} \hat{\rho}_x^2$, and $\Omega_{\hat{u}} = B(\hat{x}, \hat{\rho}_x) \times Y$.

Then Lemma 4 proves (C²) with

$$\mathcal{G}(u, \hat{u}) = (\gamma_F - \tilde{\gamma}_F) \|x - \hat{x}\|^2 + (\gamma_{G_*} - \tilde{\gamma}_{G_*}) \|y - \hat{y}\|^2.$$

To see this, we need to prove (30'). Now

$$a_K(u, \hat{u}) := \langle A(x) - A(\hat{x}) + DA(x)(\hat{x} - x) | y \rangle \quad (u, \hat{u} \in X \times Y). \quad (32)$$

Arguing with the mean value equality and the Lipschitz assumption as in Lemma 2, we get $a_K(\hat{u}, u) + a_K(u, \hat{u}) \leq \frac{L_{DA}}{2} (\|y\|_Y + \|\hat{y}\|_Y) \|x - \hat{x}\|^2$. Thus (a) implies (30'). By (32), the mean-value equality, and the Lipschitz assumption, also

$$\begin{aligned}
a_K(u, \hat{u}) + a_K(\hat{u}, u) &= \langle [DA(x) - DA(\hat{x})](\hat{x} - x) | \hat{y} \rangle \\
&\quad + \langle A(x) - A(\hat{x}) + DA(x)(\hat{x} - x) | y - \hat{y} \rangle \\
&\leq L_{DA} \|x - \hat{x}\|_X^2 (\|\hat{y}\|_Y + \frac{1}{2} \|y - \hat{y}\|_Y).
\end{aligned}$$

Using Cauchy's inequality and (b) we deduce (30').

Remark 4 In the last two examples, we need to bound some of the iterates, and to initialise close enough to a solution. Showing that the iterates stay in a local neighbourhood is a large part of the work in [23, 24], as discussed in Remark 3.

4.5 Second-order growth conditions for block-adapted methods

We now study second-order growth for problems with a block structure as in Section 3.3:

> Standing assumption

In this subsection, F and G_* are as in Section 3.3, each component subdifferential ∂F_j now (γ_{F_j} -strongly) monotone at \hat{x}_j for \hat{x}_j^* and each ∂G_{ℓ^*} ($\gamma_{G_{\ell^*}}$ -strongly) monotone at \hat{y}_ℓ for \hat{y}_ℓ^* . Here $\hat{x}_j, \hat{x}_j^*, \hat{y}_\ell$ and \hat{y}_ℓ^* are the components of $\hat{x}, \hat{x}^*, \hat{y}$, and \hat{y}^* in the corresponding subspace, assumed to satisfy the critical point condition (28).

As only some of the component functions may have $\gamma_{F_j}, \gamma_{G_{\ell^*}} > 0$, through detailed analysis of the block structure, we hope to obtain (strong) convergence on some subspaces even if the entire primal or dual variables might not converge.

Similarly to Lemma 4 we prove:

Lemma 6 Suppose for some neighbourhood $\Omega_{\hat{u}} \subset X \times Y$ that

$$\Delta_{k+1} := \sum_{j=1}^m \tilde{\gamma}_{F_j} \|x_j - \hat{x}_j\|_{X_j}^2 + \sum_{\ell=1}^n \tilde{\gamma}_{G_{\ell^*}} \|y_\ell - \hat{y}_\ell\|_{Y_\ell}^2 \geq a_K(\hat{u}, u) + a_K(u, \hat{u})$$

for some $\tilde{\gamma}_{F_j}, \gamma_{G_{\ell^*}} \geq 0$ for all $u \in \Omega_{\hat{u}}$. Then (C²) holds with

$$\mathcal{G}(u, \hat{u}) = \sum_{j=1}^m (\gamma_{F_j} - \tilde{\gamma}_{F_j}) \|x_j - \hat{x}_j\|_{X_j}^2 + \sum_{\ell=1}^n (\gamma_{G_{\ell^*}} - \tilde{\gamma}_{G_{\ell^*}}) \|y_\ell - \hat{y}_\ell\|_{Y_\ell}^2. \quad (33)$$

In the convex–concave case, we can transfer all strong monotonicity into \mathcal{G} :

Example 15 If K is convex–concave, then by Lemmas 5 and 6, (C²) holds with $\Omega_{\hat{u}} = X \times Y$ and \mathcal{G} as in (33) for $\tilde{\gamma}_{F_j} = 0$ and $\tilde{\gamma}_{G_{\ell^*}} = 0$. We have $\mathcal{G}(\cdot, \hat{u}) \geq 0$ always.

Example 16 As in Example 10, suppose we can write $K(x, y) = \sum_{j=1}^m \sum_{\ell=1}^n K_{j\ell}(x_j, y_\ell)$ with each $K_{j\ell}$ Lipschitz-continuously differentiable with the factor $L_{j\ell}$ in Ω . Then

using (27) and Lemma 6, we see (C²) to hold with $\Omega_{\hat{u}} = \Omega$ and \mathcal{G} as in (33) with

$$\tilde{\gamma}_{F_j} = \sum_{\ell=1}^n L_{j\ell} \quad (j = 1, \dots, m) \quad \text{and} \quad \tilde{\gamma}_{G_{\ell^*}} = \sum_{j=1}^m L_{j\ell} \quad (\ell = 1, \dots, n).$$

Thus $\mathcal{G}(\cdot, \hat{u}) \geq 0$ if $\gamma_{F_j} \geq \sum_{\ell=1}^n L_{j\ell}$ and $\gamma_{G_{\ell^*}} \geq \sum_{j=1}^m L_{j\ell}$ for all ℓ and j .

The special case of Example 10 with each $K_{j\ell}$ bilinear, corresponding to Example 11 for ellipticity, is covered by Example 15.

We consider in detail the two dual block setup of Examples 2 and 8:

Example 17 As in Example 2, let $K(x, y) = \langle A_1(x)|y_1 \rangle + \langle A_2x|y_2 \rangle$ for $A_1 \in C^1(X; Y_1^*)$ and $A_2 \in \mathbb{L}(X; Y_2^*)$. Then, as in (32),

$$a_K(u, \bar{u}) = \langle A_1(x) - A_1(\bar{x}) + DA_1(x)(\bar{x} - x)|y_1 \rangle,$$

which does not depend on A_2 . For any $\alpha, \rho_y, \hat{\rho}_x > 0$ let either

- (a) $\tilde{\gamma}_F \geq \frac{L_{DA_1}}{2}(\rho_{y_1} + \|\hat{y}_1\|_{Y_1})$, $\tilde{\gamma}_{G_{1^*}} \geq 0$, and $\Omega_{\hat{u}} = X \times B(0, \rho_{y_1})$; or
- (b) $\tilde{\gamma}_F > L_{DA_1}(\|\hat{y}_1\|_{Y_1} + \frac{\alpha}{2})$, $\tilde{\gamma}_{G_{1^*}} \geq \frac{L_{DA_1}}{2\alpha}\hat{\rho}_x^2$, and $\Omega_{\hat{u}} = B(\hat{x}, \hat{\rho}_x) \times Y$.

Arguing as in Example 14 and using Lemma 6, we then see (C²) to hold with \mathcal{G} as in (33) and $\tilde{\gamma}_{G_{2^*}} = 0$. In this case $\mathcal{G}(\cdot, \hat{u})$ is non-negative if $\gamma_F \geq \tilde{\gamma}_F$ and $\gamma_{G_{1^*}} \geq \tilde{\gamma}_{G_{1^*}}$.

4.6 Convergence of iterates

We are now ready to prove the convergence of the iterates. We start with weak convergence and proceed to strong and linear convergence. For weak convergence in infinite dimensions, we need some further technical assumptions. We recall that a set-valued map $T : X \rightrightarrows X^*$ is weak-to-strong (weak*-to-strong) outer semicontinuous if $x_k^* \in T(x^k)$ and $x^k \rightharpoonup x$ ($x^k \xrightarrow{*} x$) and $x_k^* \rightarrow x^*$ imply $x^* \in T(x)$. The non-reflexive case of the next assumption covers spaces of functions of bounded variation [1, Remark 3.12], important for total variation based imaging.

Assumption 1 Each of the spaces X and Y is, individually, either a reflexive Banach space or the dual of separable space. The operator $H : X \times Y \rightrightarrows X^* \times Y^*$ is weak(-*)-to-strong outer semicontinuous, where we mean by “weak(-*)” that we take the weak topology if the space is reflexive and weak-* otherwise, individually on X and Y .

Subdifferentials of lower semicontinuous convex functions are weak(-*)-to-strong outer semicontinuous⁵, so the outer semicontinuity of H depends mainly on K .

⁵ This result seems difficult to find in the literature for Banach spaces, but follows easily from the definition of the subdifferential: If $F(x) \geq F(x^k) + \langle x_k^* | x - x^k \rangle$ and $x_k^* \rightarrow \hat{x}^*$ as well as $x^k \rightharpoonup$ (or $\xrightarrow{*}$) \hat{x} , then, using the fact that $\{\|x^k - \hat{x}\|\}_{k \in \mathbb{N}}$ is bounded, in the limit $F(x) \geq F(\hat{x}) + \langle \hat{x}^* | x - \hat{x} \rangle$.

Example 18 If X and Y are finite-dimensional, [Assumption 1](#) holds if $K \in C^1(X; Y)$.

Example 19 More generally, [Assumption 1](#) holds if $K \in C^1(X \times Y)$ and DK is continuous from the weak(-*) topology to the strong topology.

Example 20 If $K = \langle Ax|y \rangle + E(x)$ for $A \in \mathbb{L}(X; Y^*)$ and $E \in C^1(X)$ convex, then H satisfies [Assumption 1](#). Indeed, it can be shown that H is maximal monotone, hence weak(-*) outer semicontinuous similarly to convex subdifferentials.

> Verification of the conditions

To verify the nonsmooth second-order growth condition (C^2) for each of the following [Theorems 2 to 4](#), we point to [Sections 4.4 and 4.5](#). For the verification of the (semi-)ellipticity of B^0 , we point to [Sections 4.2 and 4.3](#). As special cases of the PDBS [\(16\)](#), the theorems apply to the Hilbert-space PDPS [\(17\)](#) and its block-adaptation [\(18\)](#). Then J_X and J_Y are continuously differentiable and convex.

Theorem 2 (Weak convergence) *Let F and G_* be convex, proper, and lower semicontinuous; $K \in C^1(X \times Y)$; and both $J_X \in C^1(X)$ and $J_Y \in C^1(Y)$ convex. Suppose [Assumption 1](#) holds and for some $\hat{u} \in H^{-1}(0)$ that*

- (i) (C^2) holds with $\mathcal{G}(\cdot, \hat{u}) \geq 0$ within $\Omega_{\hat{u}} \subset X \times Y$; and
- (ii) B^0 is elliptic within $\Omega \ni \hat{u}$.

Let $\{u^{k+1}\}_{k \in \mathbb{N}}$ be generated by the PDBS [\(16\)](#) for any initial u^0 , and suppose $\{u^k\}_{k \in \mathbb{N}} \subset \Omega \cap \Omega_{\hat{u}}$. Then there exists at least one cluster point of $\{u^k\}_{k \in \mathbb{N}}$, and all weak(-*) cluster points belong to $H^{-1}(0)$.

Proof [Lemma 3](#) establishes (D) for $B = B^0$ and all $N \geq 1$. With $\varepsilon > 0$ the factor of ellipticity of B^0 , it follows

$$\frac{\varepsilon}{2} \|u^N - \hat{u}\|_{X \times Y}^2 + \frac{\varepsilon}{2} \sum_{k=0}^{N-1} \|u^{k+1} - u^k\|_{X \times Y}^2 \leq B^0(\hat{u}, u^0) \quad (N \geq 1).$$

Clearly $\|u^{k+1} - u^k\| \rightarrow 0$ while $\{\|u^k - \hat{u}\|\}_{k \in \mathbb{N}}$ is bounded. Using the Eberlein–Šmuljan theorem in a reflexive X or Y , and the Banach–Alaoglu theorem otherwise (X or Y the dual of a separable space), we may therefore find a subsequence of $\{u^k\}_{k \in \mathbb{N}}$ converging weakly(-*) to some \bar{x} . Since $J^0 \in C^1(X \times Y)$, we deduce $D_1 B^0(u^{k+1}, u^k) \rightarrow 0$. Consequently [\(15\)](#) implies that $0 \in \limsup_{k \rightarrow \infty} H(u^{k+1})$, where we write “lim sup” for the Painlevé–Kuratowski outer limit of a sequence of sets in the strong topology. Since H is weak(-*)-to-strong outer semicontinuous by [Assumption 1](#), it follows that $0 \in H(\hat{u})$. Therefore, there exists at least one cluster point of $\{u^k\}_{k \in \mathbb{N}}$ belonging to $H^{-1}(0)$. Repeating the argument on any weak(-*) convergent subsequence, we deduce that all cluster points belong to $H^{-1}(0)$. \square

Remark 5 For a unique weak limit we may in Hilbert spaces use the quantitative Féjer monotonicity (F) with Opial's lemma [53, 13]. For bilinear K the result is relatively immediate, as B^0 is a squared matrix-weighted norm; see [63]. Otherwise a variable-metric Opial's lemma [23] and additional work based on the Brezis–Crandall–Pazy lemma [12, Corollary 20.59 (iii)] is required; see [23] for $K(x, y) = \langle A(x)|y \rangle$, and [24] for general K .

Theorem 3 (Strong convergence) *Let F and G_* be convex, proper, and lower semicontinuous; $K \in C^1(X \times Y)$; and both $J_X \in C(X)$ and $J_Y \in C(Y)$ convex and Gâteaux-differentiable. Suppose for some $\hat{u} \in H^{-1}(0)$ that*

- (i) (C^2) holds with $\mathcal{G}(\cdot, \hat{u}) \geq 0$ within $\Omega_{\hat{u}} \subset X \times Y$; and
- (ii) B^0 is semi-elliptic within $\Omega \ni \hat{u}$.

Let $\{u^{k+1}\}_{k \in \mathbb{N}}$ be generated by the PDBS (16) for any initial u^0 , and suppose $\{u^k\}_{k \in \mathbb{N}} \subset \Omega \cap \Omega_{\hat{u}}$. Then $\mathcal{G}(u^{k+1}, \hat{u}) \rightarrow 0$ as $N \rightarrow \infty$.

In particular, if $\mathcal{G}(u, \hat{u}) \geq \|P(u - \hat{u})\|_Z^2$ for some $P \in \mathbb{L}(X; Z)$, then $Px^N \rightarrow P\hat{x}$ strongly in Z and the ergodic sequence $\bar{x}_P^N := \frac{1}{N} \sum_{k=0}^{N-1} Px^{k+1} \rightarrow P\hat{x}$ at rate $O(1/N)$.

Proof Lemma 3 establishes (D). By the semi-ellipticity of B^0 then $\sum_{k=0}^{N-1} \mathcal{G}(u^{k+1}, \hat{u}) \leq B^0(\hat{u}, u^0)$, ($N \in \mathbb{N}$). Since $\mathcal{G}(u^{k+1}, \hat{u}) \geq 0$, this shows that $\mathcal{G}(u^N, \hat{u}) \rightarrow 0$. The strong convergence of the primal variable for quadratically minorised \mathcal{G} is then immediate whereas following by Jensen's inequality gives the ergodic convergence claim. \square

Example 21 In Section 4.4, we can take $Pu = \sqrt{\gamma_F - \tilde{\gamma}_F}x$ if $\gamma_F > \tilde{\gamma}_F$ or $Pu = \sqrt{\gamma_{G_*} - \tilde{\gamma}_{G_*}}y$ if $\gamma_{G_*} > \tilde{\gamma}_{G_*}$. The examples of Section 4.5 for $x = (x_1, \dots, x_m)$, $y = (y_1, \dots, y_n)$ may allow $Pu = \sqrt{\gamma_{F_j} - \tilde{\gamma}_{F_j}}x_j$ or $Pu = \sqrt{\gamma_{G_{\ell_*}} - \tilde{\gamma}_{G_{\ell_*}}}y_{\ell}$.

Remark 6 Under similar conditions as Theorem 3, it is possible to obtain $O(1/N^2)$ convergence rates; see [17, 63] for the convex-concave case and [23, 24] in general.

Theorem 4 (Linear convergence) *Let F and G_* be convex, proper, and lower semicontinuous; $K \in C^1(X \times Y)$; and both $J_X \in C(X)$ and $J_Y \in C(Y)$ convex and Gâteaux-differentiable. Suppose for some $\gamma > 0$ and $\hat{u} \in H^{-1}(0)$ that*

- (i) (C^2) holds with $\mathcal{G}(u, \hat{u}) \geq \gamma B^0(\hat{u}, u)$ within $\Omega_{\hat{u}} \subset X \times Y$; and
- (ii) B^0 is elliptic within $\Omega \supset \hat{u}$.

Let $\{u^{k+1}\}_{k \in \mathbb{N}}$ be generated by the PDBS (16) for any initial u^0 , and suppose $\{u^k\}_{k \in \mathbb{N}} \subset \Omega \cap \Omega_{\hat{u}}$. Then $B^0(\hat{u}, u^N) \rightarrow 0$ and $u^N \rightarrow \hat{u}$ at a linear rate.

In particular, if $\mathcal{G}(u, \hat{u}) \geq \gamma \|u - \hat{u}\|^2$, ($k \in \mathbb{N}$), for some $\gamma > 0$, and J^0 is Lipschitz-continuously differentiable, then $u^N \rightarrow \hat{u}$ at a linear rate.

Proof Lemma 3 establishes the quantitative Δ -Féjer monotonicity (F). Using (i), this yields $(1 + \gamma)B^0(\hat{u}, u^{k+1}) \leq B^0(\hat{u}, u^k)$. By the semi-ellipticity of B^0 , the claimed linear convergence of $B^0(\hat{u}, u^N) \rightarrow 0$ follows. Since B^0 is assumed elliptic, also $u^N \rightarrow \hat{u}$ linearly. If J^0 is Lipschitz-continuously differentiable, then, similarly to Lemma 2, $B^0(\hat{u}, u^{k+1}) \leq L_{DJ} \|u^{k+1} - \hat{u}\|^2$ for some $L_{DJ} > 0$. Thus $\mathcal{G}(u^{k+1}, \hat{u}) \geq \gamma L_{DJ}^{-1} B^0(\hat{u}, u^{k+1})$, so the main claim establishes the particular claim. \square

Example 22 J^0 is Lipschitz-continuously differentiable if X and Y are Hilbert spaces with $J_X = \tau^{-1}N_X$ and $J_Y = \sigma^{-1}N_Y$, and K Lipschitz-continuously differentiable.

4.7 Convergence of gaps in the convex-concave setting

We finish this section by studying the convergence of gap functionals in the convex-concave setting.

Lemma 7 *Suppose F and G_* are convex, proper, and lower semicontinuous, and $K \in C^1(X \times Y)$ is convex-concave on $\text{dom } F \times \text{dom } G_*$. Then (\mathcal{C}^2) holds for all $\bar{u} \in X \times Y$ with $\Omega_{\bar{u}} = X \times Y$ and $\mathcal{G} = \mathcal{G}^{\mathcal{L}}$ the Lagrangian gap*

$$\begin{aligned} \mathcal{G}^{\mathcal{L}}(u, \bar{u}) &:= \mathcal{L}(x, \bar{y}) - \mathcal{L}(\bar{x}, y) \\ &= [F(x) + K(x, \bar{y}) - G_*(\bar{y})] - [F(\bar{x}) + K(\hat{x}, y) - G_*(y)]. \end{aligned}$$

This functional is non-negative if $\bar{u} \in H^{-1}(0)$.

Moreover, if $\sum_{k=0}^{N-1} \mathcal{G}^{\mathcal{L}}(u^{k+1}, \bar{u}) \leq M(\bar{u})$ for some $M(\bar{u}) \geq 0$, for all $\bar{u} \in X \times Y$ and all $N \in \mathbb{N}$, and we define the ergodic sequence $\bar{u}^N := \frac{1}{N} \sum_{k=0}^{N-1} u^{k+1}$, then

- (i) $0 \leq \frac{1}{N} \sum_{k=0}^{N-1} \mathcal{G}^{\mathcal{L}}(u^{k+1}, \hat{u}) \rightarrow 0$ at the rate $O(1/N)$ for $\hat{u} \in H^{-1}(0)$.
- (ii) $0 \leq \mathcal{G}^{\mathcal{L}}(\bar{u}^N, \hat{u}) \rightarrow 0$ at the rate $O(1/N)$ for $\hat{u} \in H^{-1}(0)$.
- (iii) If $M \in C(X \times Y)$ and $\Omega \subset X \times Y$ is bounded with $\Omega \cap H^{-1}(0) \neq \emptyset$, then $0 \leq \mathcal{G}_{\Omega}(\bar{u}^N) \rightarrow 0$ at the rate $O(1/N)$ for the partial gap $\mathcal{G}_{\Omega}(u) := \sup_{\bar{u} \in \Omega} \mathcal{G}^{\mathcal{L}}(u, \bar{u})$.

The convergence results in Lemma 7 are *ergodic* because they apply to sequences of running averages. To understand the partial gap, we recall that with $K(x, y) = \langle Ax|y \rangle$ bilinear Fenchel–Rockafellar’s theorem show that the *duality gap* $\mathcal{G}^D(u) := [F(x) + G_*(Ax)] + [F_*(-A^*y) + G_*(y)] \geq 0$ and is zero if and only if $u \in H^{-1}(0)$. The duality gap can be written $\mathcal{G}^D(u) = \mathcal{G}_{X \times Y}(u)$.

Proof By the convex-concavity of K and the definition of the subdifferential,

$$\begin{aligned} &\langle D_x K(x, y)|x - \bar{x} \rangle - \langle D_y K(x, y)|y - \bar{y} \rangle \\ &\geq [K(x, y) - K(\bar{x}, y)] - [K(x, y) - K(x, \bar{y})] = K(x, \bar{y}) - K(\bar{x}, y). \end{aligned}$$

for all $(x, y) \in X \times Y$. Also using $x^* \in \partial F(x^{k+1})$ and $y^* \in \partial G(y^{k+1})$ with the definition of the convex subdifferential, we see that $\mathcal{G} = \mathcal{G}^{\mathcal{L}}$ satisfies (\mathcal{C}^2) . The non-negativity of $\mathcal{G}(\cdot, \hat{u})$ follows by similar reasoning, first using that

$$K(x, \hat{y}) - K(\hat{x}, y) \geq \langle D_x K(\hat{x}, \hat{y})|x - \hat{x} \rangle - \langle D_y K(\hat{x}, \hat{y})|y - \hat{y} \rangle \quad (34)$$

for all $(x, y) \in X \times Y$, and following by the definition of the subdifferential applied to $-D_x K(\hat{x}, \hat{y}) \in \partial F(\hat{x})$ and $D_y K(\hat{x}, \hat{y}) \in \partial G_*(\hat{y})$.

For (i)–(iii), we first observe that the semi-ellipticity of B^0 and (\mathcal{C}^2) imply $\sum_{k=0}^{N-1} \mathcal{G}^{\mathcal{L}}(u^{k+1}, \bar{u}) \leq M(\bar{u})$. Dividing by N and using that $\mathcal{G}^{\mathcal{L}}(u^{k+1}, \hat{u}) \geq 0$ for $\bar{u} \in H^{-1}(0)$, we obtain (i). Jensen’s inequality then gives $\mathcal{G}^{\mathcal{L}}(\bar{u}^{k+1}, \bar{u}) \leq M(\bar{u})/N$, hence (ii) for $\bar{u} \in H^{-1}(0)$. Finally, taking the supremum over $\bar{u} \in \Omega$ gives (iii) because M is bounded on bounded sets. \square

In the following theorem, we may in particular take $K(x, y) = \langle Ax|y \rangle$ bilinear, or $K(x, y) = \langle Ax|y \rangle + E(x)$ with E convex. [Lemma 2](#) and [Examples 4](#) and [5](#) provide step length conditions that ensure the semi-ellipticity required of B^0 in [Theorem 5](#).

Theorem 5 (Gap convergence) *Let $F : X \rightarrow \bar{\mathbb{R}}$ and $G_* : Y \rightarrow \bar{\mathbb{R}}$ be convex, proper, and lower semicontinuous. Also let $K \in C^1(X \times Y)$ be convex-concave within $\text{dom } F \times \text{dom } G_*$. Finally, let $J_X \in C^1(X)$ and $J_Y \in C^1(Y)$ convex. If B^0 is semi-elliptic, then the iterates $\{u^{k+1}\}_{k \in \mathbb{N}}$ generated by the PDBS (16) for any initial $u^0 \in X \times Y$ satisfy [Lemma 7 \(i\)–\(iii\)](#).*

Proof By [Lemma 7](#), holds with $\mathcal{G} = \mathcal{G}^{\mathcal{L}}$. Hence by [Lemma 3](#), (D) holds. Since B^0 is semi-elliptic, this implies that that $\sum_{k=0}^{N-1} \mathcal{G}(u^{k+1}, \bar{u}) \leq M(\bar{u}) := B^0(\bar{u}, u^0)$ for all $N \in \mathbb{N}$. Since J_X, J_Y , and K are continuously differentiable, $M \in C^1(X \times Y)$. The rest follows from the second part of [Lemma 7](#). \square

5 Inertial terms

We now generalise (BP), making the involved Bregman divergences dependent on the iteration k and earlier iterates:

$$0 \in H(u^{k+1}) + D_1 B_{k+1}(u^{k+1}, u^k) + D_1 B_{k+1}^-(u^k, u^{k-1}), \quad (\text{IPP})$$

for $B_{k+1} := B_{J_{k+1}}$ and $B_{k+1}^- := B_{J_{k+1}^-}$ generated by $J_{k+1}, J_{k+1}^- : U \rightarrow \bar{\mathbb{R}}$. We take $u^{-1} := u^0$ for this to be meaningful for $k = 0$. Our main reason for introducing the dependence on u^{k-1} is improve (16) and (17) to be explicit in K when K is not affine in y : otherwise the dual step of those methods is in general not practical to compute unlike the affine case of [Remark 1](#). Along the way we also construct a more conventional inertial method.

5.1 A generalisation of the fundamental theorem

We realign indices to get a simple fundamental condition to verify on each iteration:

Theorem 6 *On a Banach space U , let $H : U \rightrightarrows U^*$, and let $J_k, J_k^- : U \rightarrow \bar{\mathbb{R}}$ be Gâteaux-differentiable with the corresponding Bregman divergences $B_k := B_{J_k}$ and $B_k^- := B_{J_k^-}$ for all $k = 1, \dots, N$. Suppose (IPP) is solvable for $\{u^{k+1}\}_{k \in \mathbb{N}}$ given an initial iterate $u^0 \in U$. If for all $k = 0, \dots, N-1$, for some $\bar{u} \in U$ and $\mathcal{G}(u^{k+1}, \bar{u}) \in \mathbb{R}$, for all $h^{k+1} \in H(u^{k+1})$ the modified fundamental condition*

$$\langle h^{k+1} | u^{k+1} - \bar{u} \rangle \geq [(B_{k+2} + B_{k+3}^-) - (B_{k+1} + B_{k+2}^-)](u, u^{k+1}) + \mathcal{G}(u^{k+1}, \bar{u}) \quad (\text{IC})$$

holds, and B_{k+1}^- satisfies the general Cauchy inequality

$$\langle D_1 B_{k+1}^-(u^k, u) | u^k - u' \rangle \leq B'_{k+1}(u^k, u) + B''_{k+1}(u', u^k) \quad (u, u' \in X) \quad (35)$$

for some $B'_{k+1}, B''_{k+1} : U \times U \rightarrow \mathbb{R}$, then we have the modified descent inequality

$$\begin{aligned} [B_{N+1} + B_{N+2}^- - B''_{N+1}](\bar{u}, u^N) + \sum_{k=0}^{N-1} [B_{k+1} + B_{k+2}^- - B''_{k+1} - B'_{k+2}](u^{k+1}, u^k) \\ + \sum_{k=0}^{N-1} \mathcal{G}(u^{k+1}, \bar{u}) \leq [B_1 + B_2^-](\bar{u}, u^0). \quad (\text{ID}) \end{aligned}$$

Proof We can write (IPP) as

$$0 = h^{k+1} + D_1 B_{k+1}(u^{k+1}, u^k) + D_1 B_{k+1}^-(u^k, u^{k-1}) \quad \text{for some } h^{k+1} \in H(u^{k+1}). \quad (36)$$

Testing (IPP) by applying $\langle \cdot | u^{k+1} - \bar{u} \rangle$ we obtain

$$0 = \langle h^{k+1} + D_1 B_{k+1}(u^{k+1}, u^k) + D_1 B_{k+1}^-(u^k, u^{k-1}) | u^{k+1} - \bar{u} \rangle.$$

Summing over $k = 0, \dots, N-1$ and using $u^{-1} = u^0$ to eliminate $B_1^-(u^0, u^{-1}) = 0$, we rearrange

$$0 = S_N + \sum_{k=0}^{N-1} \langle h^{k+1} + D_1 [B_{k+1} + B_{k+2}^-](u^{k+1}, u^k) | u^{k+1} - \bar{u} \rangle \quad (37)$$

for

$$S_N := \langle D_1 B_{N+1}^-(u^N, u^{N-1}) | \bar{u} - u^N \rangle + \sum_{k=0}^{N-1} \langle D_1 B_{k+1}^-(u^k, u^{k-1}) | u^{k+1} - u^k \rangle.$$

Abbreviating $\bar{B}_{k+1} := B_{k+1} + B_{k+2}^-$ and using (IC) and the three-point identity (8) in (37) we obtain

$$0 \geq S_N + \sum_{k=0}^{N-1} \left(\bar{B}_{k+2}(\bar{u}, u^{k+1}) - \bar{B}_{k+1}(\bar{u}, u^k) + \bar{B}_{k+1}(u^{k+1}, u^k) + \mathcal{G}(u^{k+1}, \bar{u}) \right).$$

Using the generalised Cauchy inequality (35) and, again, that $u^{-1} = u^0$, we get

$$\begin{aligned} S_N &\geq -B'_{N+1}(u^N, u^{N-1}) - B''_{N+1}(\bar{u}, u^N) - \sum_{k=0}^{N-1} \left(B'_{k+1}(u^k, u^{k-1}) + B''_{k+1}(u^{k+1}, u^k) \right) \\ &= -B''_{N+1}(\bar{u}, u^N) - \sum_{k=0}^{N-1} [B'_{k+1} + B'_{k+2}](u^{k+1}, u^k). \end{aligned}$$

These two inequalities yield (ID). \square

5.2 Inertia (almost) as usually understood

We take $J_{k+1} = J^0$ and $J_{k+1}^- = -\lambda_k J^0$ for some $\lambda_k \in \mathbb{R}$. We then expand (IPP) as

Inertial PDBS

Iteratively over $k \in \mathbb{N}$, solve for x^{k+1} and y^{k+1} :

$$\begin{aligned} (1 + \lambda_k)[DJ_X(x^k) - D_x K(x^k, y^k)] - \lambda_k [DJ_X(x^{k-1}) - D_x K(x^{k-1}, y^{k-1})] \\ \in DJ_X(x^{k+1}) + \partial F(x^{k+1}), \\ (1 + \lambda_k)[DJ_Y(y^k) - D_y K(x^k, y^k)] - \lambda_k [DJ_Y(y^{k-1}) - D_y K(x^{k-1}, y^{k-1})] \\ \in DJ_Y(y^{k+1}) + \partial G_*(y^{k+1}) - 2D_y K(x^{k+1}, y^{k+1}) \end{aligned} \quad (38)$$

If X and Y are Hilbert spaces with $J_X = \tau^{-1}N_X$ and $J_Y = \sigma^{-1}N_Y$ the standard generating functions divided by some step length parameters $\tau, \sigma > 0$, and $K(x, y) = \langle Ax|y \rangle$ for $A \in \mathbb{L}(X; Y)$, (38) reduces to the inertial method of [18]:

Inertial PDPS for bilinear K

With initial $\tilde{x}^0 = x^0$ and $\tilde{y}^0 = y^0$, iterate over $k \in \mathbb{N}$:

$$\begin{aligned} x^{k+1} &:= \text{prox}_{\tau F}(\tilde{x}^k - \tau A^* \tilde{y}^k), \\ y^{k+1} &:= \text{prox}_{\sigma G_*}(\tilde{y}^k + \sigma A(2x^{k+1} - \tilde{x}^k)), \\ \tilde{x}^{k+1} &:= (1 + \lambda_{k+1})x^{k+1} - \lambda_{k+1}x^k, \\ \tilde{y}^{k+1} &:= (1 + \lambda_{k+1})y^{k+1} - \lambda_{k+1}y^k. \end{aligned} \quad (39)$$

More generally, however, (38) does not directly apply inertia to the iterates. It applies inertia to K .

The general Cauchy inequality (35) automatically holds by the three-point identity (8) with $J_{k+1}'' = J_{k+1}' = J_{k+1}^-$ if $B_{k+1}^- \geq 0$, which is to say that J_{k+1}^- is convex. This is the case if $\lambda_k \leq 0$. For usual inertia we, however, want $\lambda_k > 0$. We will therefore use Lemma 1, requiring:

Assumption 2 For some $\beta > 0$, in a domain $\Omega \subset X \times Y$,

$$|\langle D_1 B^0(u^k, u) | u^k - u \rangle| \leq B^0(u^k, u) + \beta B^0(u', u^k) \quad (u, u', u^k \in \Omega). \quad (40)$$

Moreover, the parameters $\{\lambda_k\}_{k \in \mathbb{N}}$ are non-increasing and for some $\varepsilon > 0$,

$$0 \leq \lambda_{k+1} \leq \frac{1 - \varepsilon - \lambda_k \beta}{2} \quad (k \in \mathbb{N}). \quad (41)$$

Example 23 Suppose the generating function J^0 is γ -strongly subdifferentiable (i.e., B^0 is γ -elliptic, see Sections 4.2 and 4.3) within $\Omega \subset X \times Y$ and satisfies the subdifferential smoothness property (10) with the factor $L > 0$. Then by Lemma 1, (40) holds with $\beta = L\gamma^{-1}$ in some domain $\Omega \subset X \times Y$.

As a particular case, let X and Y be Hilbert spaces with the standard generating functions $J_X = \tau^{-1}N_X$, $J_Y = \sigma^{-1}N_Y$. Also let DK be L_{DK} -Lipschitz within Ω . Then J^0 is Lipschitz with factor $L = \max\{\sigma^{-1}, \tau^{-1}\} + L_{DK}$. Consequently the required subdifferential smoothness property (10) holds with the same factor L ; see [5, Theorem 18.15] or [63, Appendix C].

We computed L_{DK} for some specific K in Section 4.2.

Example 24 If $K(x, y) = \langle Ax | y \rangle$ with $A \in \mathbb{L}(X; Y^*)$, and if $J_X = \tau^{-1}N_X$, $J_Y = \sigma^{-1}N_Y$, in Hilbert spaces X and Y , then $B^0(u', u) = \frac{1}{2\tau}\|x - x'\|^2 + \frac{1}{2\sigma}\|y - y'\|^2 + \langle A(x - x') | y - y' \rangle$. By standard Cauchy inequality, (40) holds for $\beta = 1$ in $\Omega = X \times Y$. Consequently the next example recovers the upper bound for λ in [18]:

Example 25 The bound (41) holds for some $\varepsilon > 0$ if $\lambda_k \equiv \lambda$ for $0 \leq \lambda < 1/(2 + \beta)$.

Lemma 8 *Suppose Assumption 2 holds and that (C²) holds within $\Omega_{\bar{u}}$ for some $\bar{u} \in \Omega$ and $\mathcal{G}(u, \bar{u})$. Given $u^0 \in \Omega$, suppose the iterates generated by the inertial PDBS (38) satisfy $\{u^k\}_{k=0}^N \subset \Omega_{\bar{u}} \cap \Omega$. Then*

$$\varepsilon B^0(\bar{u}, u^N) + \varepsilon \sum_{k=0}^{N-1} B^0(u^{k+1}, u^k) + \sum_{k=0}^{N-1} \mathcal{G}(u^{k+1}, \bar{u}) \leq (1 - \lambda_1) B^0(\bar{u}, u^0). \quad (42)$$

Proof Since $B_{k+1} = B^0$ and $B_{k+1}^- = -\lambda_k B^0$ for all $k \in \mathbb{N}$,

$$(B_{k+2} + B_{k+3}^-) - (B_{k+1} + B_{k+2}^-) = (\lambda_{k+1} - \lambda_{k+2}) B^0.$$

Since λ_k is decreasing and B^0 is semi-elliptic within $\Omega \supset \{u^k, \bar{u}\}$, we deduce that $(\lambda_{k+1} - \lambda_{k+2}) B^0(\bar{u}, u^k) \geq 0$. Consequently (IC) holds if (C) does. By the proof of Lemma 3, (IC) then holds if (C²) does. Using (40), (35) holds with $B'_{k+1} = \lambda_k B_0$ and $B''_{k+1} = \lambda_k \beta B_0$. Referring to Theorem 6, we now obtain (ID). We expand

$$\begin{aligned} [B_{N+1} + B_{N+2}^- - B''_{N+1}](\bar{u}, u^N) &= (1 - \lambda_{k+1} - \lambda_k \beta) B^0(\bar{u}, u^N) \quad \text{and} \\ [B_{k+1} + B_{k+2}^- - B''_{k+1} - B'_{k+2}](u^{k+1}, u^k) &= (1 - \lambda_{k+1} - \lambda_k \beta - \lambda_{k+1}) B^0(u^{k+1}, u^k). \end{aligned}$$

Since $\bar{u}, u^k \in \Omega$ for all $k = 0, \dots, N$, using the ellipticity of B^0 within Ω as well as (41) we now estimate the first from below by $\varepsilon B^0(\bar{u}, u^N)$ and the second by $\varepsilon B^0(u^{k+1}, u^k)$. Thus (ID) produces (42). \square

We may now proceed as in Sections 4.6 and 4.7 to prove convergence. For the verification of Assumption 2 we can use Examples 23 to 25.

Theorem 7 (Convergence, inertial method) *Theorems 2, 3 and 5 apply to the iterates $\{u^{k+1}\}_{k \in \mathbb{N}}$ generated by the inertial PDBS (38) if we replace the assumptions of (semi-)ellipticity of B^0 with Assumption 2.*

Proof We replace Lemma 3 and (D) by Lemma 8 and (42) in the proofs of Theorems 2, 3 and 5. Observe that Assumption 2 implies that B^0 is (semi-)elliptic. \square

Remark 7 The inertial PDPS is improved in [65] to yield *non-ergodic* convergence of the Lagrangian gap. To do the “inertial unrolling” that leads to such estimates, one, however, needs to correct for the anti-symmetry introduced by K into H .

Remark 8 Since Theorem 6 does not provide the quantitative Δ -Féjer monotonicity used in Theorem 4, we cannot prove linear convergence using our present simplified “testing” approach lacking the “testing parameters” of [63].

5.3 Improvements to the basic method without dual affinity

We now have the tools to improve the basic PDBS (16) to enjoy prox-simple steps for general K not affine in y . Compared to (14) we amend $J_{k+1} = J^0$ by taking

$$\begin{aligned} J_{k+1}(x, y) &:= J_X(x) + J_Y(y) - K(x, y) + 2K(x^{k+1}, y) \\ &= J^0(x, y) + 2K(x^{k+1}, y). \end{aligned} \quad (43)$$

This would be enough for K to be explicit in the algorithm, however, proofs of convergence would practically require G_* to be strongly convex even in the convex-concave case. To fix this, we introduce the inertial term generated by

$$J_{k+1}^-(u) := [J^0 - J_k](u) = -2K(x^k, y). \quad (44)$$

As always, we write B_{k+1} , B^0 , and B_{k+1}^- for the Bregman divergences generated by J_{k+1} , J^0 , and J_{k+1}^- .

Since

$$D_1[B_{k+1} - B^0](u^k, u^{k-1}) + D_1 B_{k+1}^-(u^k, u^{k-1}) = (0, \tilde{y}_{k+1}^*)$$

for

$$\tilde{y}_{k+1}^* = 2[D_y K(x^{k+1}, y^{k+1}) - D_y K(x^{k+1}, y^k) - D_y K(x^k, y^k) + D_y K(x^k, y^{k-1})],$$

the algorithm (IPP) expands similarly to (16) as the

Modified PDBS

Iteratively over $k \in \mathbb{N}$, solve for x^{k+1} and y^{k+1} :

$$\begin{aligned} DJ_X(x^k) - D_x K(x^k, y^k) &\in DJ_X(x^{k+1}) + \partial F(x^{k+1}) \quad \text{and} \\ DJ_Y(y^k) + [2D_y K(x^{k+1}, y^k) + D_y K(x^k, y^k) - 2D_y K(x^k, y^{k-1})] \\ &\in DJ_Y(y^{k+1}) + \partial G_*(y^{k+1}). \end{aligned} \quad (45)$$

The method reduces to the basic PDBS (16) when K is affine in y . In Hilbert spaces X and Y with $J_X = \tau^{-1}N_X$ and $J_Y = \sigma^{-1}N_Y$, we can rearrange (45) as

Modified PDPS

Iterate over $k \in \mathbb{N}$:

$$\begin{aligned} x^{k+1} &:= \text{prox}_{\tau F}(x^k - \tau \nabla_x K(x^k, y^k)), \\ y^{k+1} &:= \text{prox}_{\sigma G_*}(y^k + \sigma[2\nabla_y K(x^{k+1}, y^k) + \nabla_y K(x^k, y^k) - 2\nabla_y K(x^k, y^{k-1})]). \end{aligned} \quad (46)$$

Remark 9 The modified PDPS (46) is slightly more complicated than the method in [24], which would update

$$y^{k+1} := \text{prox}_{\sigma G_*}(y^k + \sigma \nabla_y K(2x^{k+1} - x^k, y^k)).$$

Likewise, (45) is different from the algorithm presented in [35] for convex-concave K . It would, for the standard generating functions, update⁶

$$y^{k+1} := \text{prox}_{\sigma G_*}(y^k + \sigma[2\nabla_y K(x^{k+1}, y^k) - \nabla_y K(x^k, y^{k-1})]).$$

We could produce this method by taking $J_{k+1}^-(u) = -K(x^k, y)$. However, the convergence proofs would require some additional steps.

The main difference to the overall analysis of Section 4 is in bounding from below the Bregman divergences in (ID). We now have

$$B_{N+1} + B_{N+2}^- - B_{N+1}'' = B^0 - B_{N+1}'' \quad \text{and} \quad (47a)$$

$$B_{k+1} + B_{k+2}^- - B_{k+1}'' - B_{k+2}' = B^0 - B_{k+1}'' - B_{k+2}'. \quad (47b)$$

If $D_y K(x^k, \cdot)$ is $L_{DK,y}$ -Lipschitz,

$$\begin{aligned} \langle D_1 B_{k+1}^-(u^k, u) | u^k - u' \rangle &= 2 \langle D_y K(x^k, y^k) - D_y K(x^k, y) | y^k - y' \rangle \\ &\leq \sqrt{L_{DK,y}} \|y - y^k\|^2 + \sqrt{L_{DK,y}} \|y' - y^k\|^2 \\ &=: B_{k+1}'(u^k, u) + B_{k+1}''(u', u^k). \end{aligned} \quad (48)$$

Therefore, for the modified descent inequality (ID) to be meaningful, we require:

Assumption 3 We assume that $\|D_y K(x, y) - D_y K(x, y')\| \leq L_{DK,y} \|y - y'\|$ when $(x, y), (x, y') \in \Omega$ for some domain $\Omega \subset X \times Y$. Moreover, for some $\varepsilon \geq 0$ we have

$$B^0(u, u') \geq \frac{\varepsilon}{2} \|u - u'\|_{X \times Y}^2 + 2\sqrt{L_{DK,y}} \|y - y'\|_Y^2 \quad (u, u' \in \Omega). \quad (49)$$

We say that the present assumption holds *strongly* if $\varepsilon > 0$.

⁶ Note that [35] uses the historical ordering of the primal and dual updates from [17], prior to the proof-simplifying discovery of the proximal point formulation in [36]. Hence our y^k is their y^{k+1} .

Example 26 If K is affine in y , $L_{DK,y} = 0$. Therefore, [Assumption 3](#) reduces to the (semi-)ellipticity of B^0 , which can be verified as in [Sections 4.2 and 4.3](#).

Example 27 Generally, it is easy to see that if one of the results of [Section 4.2](#) holds with $\tilde{\sigma} = 1/(\sigma^{-1} - 4\sqrt{L_{DK,y}}) > 0$ in place of σ , then (49) holds. In particular, if K has L_{DK} -Lipschitz derivative within Ω , then [Lemma 2](#) gives the condition $1 \geq L_{DK} \max\{\tau, \sigma/(1 - 4\sigma\sqrt{L_{DK,y}})\}$ and $1 > 4\sigma\sqrt{L_{DK,y}}$ for (49) to hold with $\varepsilon = 0$. The assumption holds strongly if the first inequality is strict.

Similarly to [Lemma 8](#), we now have the following replacement for [Lemma 3](#):

Lemma 9 *Suppose [Assumption 3](#) holds and (C^2) holds within $\Omega_{\bar{u}}$ for some $\bar{u} \in X \times Y$ and $\mathcal{G}(u, \bar{u})$. Given $u^0 \in X \times Y$, suppose the iterates generated by the modified PDBS (45) satisfy $\{u^k\}_{k=0}^N \subset \Omega_{\bar{u}}$. Then*

$$\varepsilon B^0(\bar{u}, u^N) + \varepsilon \sum_{k=0}^{N-1} B^0(u^{k+1}, u^k) + \sum_{k=0}^{N-1} \mathcal{G}(u^{k+1}, \bar{u}) \leq [B_1 + B_2^-](\bar{u}, u^0). \quad (50)$$

Proof Inserting (43) and (44), (IC) reduces to (C), which follows from (C^2) as in [Lemma 3](#). We verify (35) via (48) and [Assumption 3](#). Thus [Theorem 6](#) proves (ID). Inserting (47) and (49) with B'_{k+1} and B''_{k+1} from (48) into (ID) proves (50). \square

We may now proceed as in [Sections 4.6 and 4.7](#) to prove convergence. For the verification of [Assumption 3](#) we can use [Examples 26 and 27](#).

Theorem 8 (Convergence, modified method) *Theorems 2, 3 and 5 apply to the iterates $\{u^{k+1}\}_{k \in \mathbb{N}}$ generated by the modified PDBS (45) if we replace the assumptions of semi-ellipticity (resp. ellipticity) of B^0 with [Assumption 3](#) holding (strongly).*

Proof We replace [Lemma 3](#) and (D) by [Lemma 9](#) and (50) in [Theorems 2, 3 and 5](#). Observe that (strong) [Assumption 3](#) implies the (semi-)ellipticity of B^0 . \square

Now we have a locally convergent method (46) with easily implementable steps to tackle problems such as Potts segmentation (4) [24].

6 Further directions

We close by briefly reviewing some things not covered, other possible extensions, and alternative algorithms.

6.1 Acceleration

To avoid technical detail, we did not cover $O(1/N^2)$ acceleration. The fundamental ingredients of proof are, however, exactly the same as we have used: sufficient

second-order growth and ellipticity of the Bregman divergences B_k^0 , which are now iteration-dependent. Additionally, a portion of the second-order growth must be used to make the metrics B_k^0 grow as $k \rightarrow \infty$. For bilinear K in Hilbert spaces, such an argument can be found in [63]; for $K(x, y) = \langle A(x)|y \rangle$ in [23]; and for general K in [24]. As mentioned in [Remarks 1](#) and [9](#), the algorithms in the latter two differ slightly from the ones presented here.

6.2 Stochastic methods

It is possible to refine the block-adapted (18) and its accelerated version into stochastic methods. The idea is to take on each step subsets of primal-blocks $S(i) \subset \{1, \dots, m\}$ and dual blocks $V(i+1) \subset \{1, \dots, n\}$ and to only update the corresponding x_j^{k+1} and y_ℓ^{k+1} . Full discussion of such technical algorithms are outside the scope of our present overview. We refer to [64] for an approach covering block-adapted acceleration and both primal- and dual randomisation in the case of bilinear K , but see also [16] for a more basic version. For more general K affine in y , see [48].

6.3 Alternative Bregman divergences

We have used Bregman divergences as a proof tool, in the end opting for the standard quadratic generating functions on Hilbert spaces. Nevertheless, our theory works for arbitrary Bregman divergences. The practical question is whether F and G_* remain prox-simple with respect to such a divergence. This can be the case for the “entropic distance” generated on $L^1(\Omega; [0, \infty))$ by

$$J(x) := \begin{cases} \int_{\Omega} x(t) \ln x(t) dt, & x \geq 0 \text{ a.e. on } \Omega, \\ \infty, & \text{otherwise} \end{cases}$$

See, for example, [14] for a Landweber method (gradient descent on regularised least squares) based on such a distance.

6.4 Alternative approaches

The derivative $D_1 B^0$ in (15) can be seen as a preconditioner, replacing $\tau(u - u')$ in the proximal point method (13). Our choice of B^0 is not the only option.

Consider the problem

$$\min_{x \in X} F(x) + E(x). \tag{51}$$

Provided E is differentiable and F *prox-simple*, i.e., the proximal map of F has a closed-form expression, (1) can be solved by forward-backward splitting methods as first introduced in [43]. In a Hilbert space X , this can be written

$$x^{k+1} := \text{prox}_{\tau F}(x^k - \tau \nabla E(x^k)). \quad (52)$$

Variants based on Bregman divergences were introduced in [50] under the name “*mirror prox*” or “*mirror descent*”; see also the review [19]. The method and convergence proofs for it can be derived from our primal-dual approach. Indeed, if we take $G_* \equiv \delta_{\{0\}}$ as the **indicator function** of zero, and $K(x, y) = E(x)$ for some $E \in C^1(X)$, then (S) is equivalent to (51). Now the dual step step of (17) is $y^{k+1} := 0$, and the primal step is (52).

Forward-backward splitting is especially popular under the name *iterative soft-thresholding* (ISTA) in the context of *sparse reconstruction* (i.e., regularisation of linear inverse problems with ℓ^1 penalties), see, e.g., [15, 27, 8]. However, forward-backward splitting has limited applicability in imaging and inverse problems due to the joint prox-simplicity and smoothness requirements. Sometimes these can be circumvented by considering so-called dual problems [7].

Let then E be Gâteaux-differentiable and $F = G \circ A$ for a nonsmooth function F and a linear operator A in (51), i.e., consider the problem

$$\min_{x \in X} E(x) + G(Ax),$$

Forward-backward splitting is impractical as $G \circ A$ is in general not prox-simple. Assuming G to have the pre-conjugate G_* , we can write this problem as an instance of (S) with $F = 0$ and $K(x, y) = E(x) + \langle Ax | y \rangle$. Therefore the methods we have presented are applicable. However, in this instance, also $J^0(u) := \frac{1}{2} \|u\|_{X \times Y}^2 + \frac{1}{2} \|A^*y\|_{X^*}^2$ would produce an algorithm with realisable steps. In analogy to the PDPS, it might be called the *primal dual explicit spitting* (PDES). The method was introduced in [45] for $E(z) = \frac{1}{2} \|b - z\|^2$ as the “generalised iterative soft-thresholding” (GIST), but has also been called the *primal-dual fixed point method* (PDFP, [20]) and the *proximal alternating predictor corrector* (PAPC, [30]).

The classical *Augmented Lagrangian* method solves the saddle point problem

$$\min_x \max_y F(x) + \frac{\tau}{2} \|E(x)\|^2 + \langle E(x) | y \rangle, \quad (53)$$

alternatingly for x and y . The *alternating directions method of multipliers* (ADMM) of [33, 3] takes $E(x) = Ax_1 + Bx_2 - c$ and $F(x) = F_1(x_1) + F_2(x_2)$ for $x = (x_1, x_2)$, and alternates between solving (53) for x_1 , x_2 , and y , using the most recent iterate for the other variables. The method cannot be expressed in our Bregman divergence framework, as the preconditioner $D_1 B_{k+1}(\cdot, x^k)$ would need to be non-symmetric. The steps of the method are potentially expensive, each itself being an optimisation problem. Hence the *preconditioned ADMM* of [69], which is equivalent to the PDPS and the classical *Douglas-Rachford splitting* (DRS, [29]) applied to appropriate problems [17, 25]. The preconditioned ADMM was extended to nonlinear E in [10].

Based on derivations avoiding the Lipschitz gradient assumption (cocoercivity) in forward-backward splitting, [47] moves the over-relaxation step $\bar{x}^{k+1} := 2x^{k+1} - x^k$ of the PDPS outside the proximal operators. This amounts to taking $J_{k+1}^- = \lambda_k K$ in Section 5.2 instead of $J_{k+1}^-(x, y) = \lambda_k J^0 = \lambda_k [\tau^{-1} J_X(x) + \sigma^{-1} J_Y(y) - K(x, y)]$, so is “partial inertia”; compare the “corrected inertia” of [65].

An *over-relaxed* variant of the same idea maybe found in [11]. We have not discussed over-relaxation of entire algorithms. To briefly relate it to the basic inertia of (39), the latter “rebases” the algorithm at the inertial iterate \tilde{u}^k constructed from u^k and u^{k-1} , whereas over-relaxation would construct \tilde{u}^k from u^k and \tilde{u}^{k-1} . The derivation in [11] is based on applying Douglas–Rachford splitting on a lifted problem. The basic over-relaxation of the PDPS is known as the Condat–Vũ method [26, 68].

6.5 Functions on manifolds and Hadamard spaces

The PDPS has been extended in [9] to functions on Riemannian manifolds; the problem $\min_{x \in \mathcal{M}} F(x) + G(Ex)$, where $E : \mathcal{M} \rightarrow \mathcal{N}$ with \mathcal{M} and \mathcal{N} Riemannian manifolds. In general, between manifolds, there are no linear maps, so E is nonlinear. Indeed, besides introducing a theory of conjugacy for functions on manifolds, the algorithm presented in [9] is based on the NL-PDPS of [62, 23].

Convergence could only be proved on *Hadamard manifolds*, which are special: a type of three-point inequality holds [28, Lemma 12.3.1]. Indeed, in even more general *Hadamard spaces* with the metric d , for any three points x^{k+1}, x^k, \bar{x} , we have [4, Corollary 1.2.5]

$$\frac{1}{2}d(x^k, x^{k+1})^2 + \frac{1}{2}d(x^{k+1}, \bar{x})^2 - \frac{1}{2}d(x^k, \bar{x})^2 \leq d(x^k, x^{k+1})d(\bar{x}, x^{k+1}). \quad (54)$$

Therefore, given a function f on such a space, to derive a simple proximal point algorithm, having constructed the iterate x^k we might try to find x^{k+1} such that

$$f(x^{k+1}) + d(x^k, x^{k+1}) \leq f(x^k).$$

Multiplying this inequality by $d(\bar{x}, x^{k+1})$ and using the three-point inequality (54),

$$\frac{1}{2}d(x^k, x^{k+1})^2 + \frac{1}{2}d(x^{k+1}, \bar{x})^2 + [f(x^{k+1}) - f(x^k)]d(\bar{x}, x^{k+1}) \leq \frac{1}{2}d(x^k, \bar{x})^2.$$

If the space is bounded, $d(\bar{x}, x^{k+1}) \leq C$, so since $f(x^k) \geq f(x^{k+1})$, we may telescope and proceed as before to obtain convergence.

The Hadamard assumption is restrictive: if a Banach space is Hadamard, it is Hilbert, while a Riemannian manifold is Hadamard if it is simply connected with a non-positive sectional curvature [4, section 1.2].

Acknowledgements Academy of Finland grants 314701 and 320022.

Glossary

The extended reals We define $\overline{\mathbb{R}} := [-\infty, \infty]$.
 A convex function A function $F : X \rightarrow \overline{\mathbb{R}}$ is convex if for all $x, x' \in X$ and $\lambda \in (0, 1)$, we have

$$F(\lambda x + (1 - \lambda)x') \leq \lambda F(x) + (1 - \lambda)F(x').$$

A concave function A function $F : X \rightarrow \overline{\mathbb{R}}$ is concave if $-f$ is convex.
 A convex-concave function A function $K : X \times Y \rightarrow \overline{\mathbb{R}}$ is convex-concave if $K(\cdot, y)$ is convex for all $y \in Y$, and $K(x, \cdot)$ is concave for all $x \in X$.

The dual space We write X^* for the dual space of a topological vector (Banach, Hilbert) space X .

Set-valued map We write $A : X \rightrightarrows Y$ if A is a set-valued map between the spaces X and Y .

Derivative We write $DF : X \rightarrow X^*$ for the derivative of a Gâteaux-differentiable function $F : X \rightarrow \overline{\mathbb{R}}$.

Convex subdifferential This is the map $\partial F : X \rightrightarrows X^*$ for a convex $F : X \rightarrow \overline{\mathbb{R}}$. By definition $x^* \in \partial F(x)$ at $x \in X$ if and only if

$$F(x') - F(x) \geq \langle x^* | x' - x \rangle \quad (x' \in X).$$

Fenchel conjugate This is the function $f^* : X^* \rightarrow \overline{\mathbb{R}}$ defined for $F : X \rightarrow \overline{\mathbb{R}}$ by

$$f^*(x^*) := \sup_{x \in X} \langle x^* | x \rangle - F(x) \quad (x^* \in X^*).$$

Fenchel preconjgate If $X = (X_*)^*$ is the dual space of some space X_* , and $F : X \rightarrow \overline{\mathbb{R}}$, then $f_* : X_* \rightarrow \overline{\mathbb{R}}$ is the preconjgate of f if $f = (f_*)^*$.

Proximal map For a function $F : X \rightarrow \overline{\mathbb{R}}$, this can be defined as

$$\text{prox}_F(x) := \arg \min_{\tilde{x} \in X} \left(F(\tilde{x}) + \frac{1}{2} \|\tilde{x} - x\|_X^2 \right).$$

Distributional derivative It arises from integration by parts: If $u : \mathbb{R}^n \supset \Omega \rightarrow \mathbb{R}$ is differentiable and $\varphi \in C_c^\infty(\Omega; \mathbb{R}^n)$, then

$$\int_{\Omega} \langle \nabla u, \varphi \rangle dx = - \int_{\Omega} u \operatorname{div} \varphi dx.$$

If now u is not differentiable, we *define* the distribution $D \in C_c^\infty(\Omega; \mathbb{R}^n)^*$ by

$$Du(\varphi) := - \int_{\Omega} u \operatorname{div} \varphi \, dx.$$

If Du is bounded (as a linear operator) it can be presented as a vector Radon measure [32], the space denoted $\mathcal{M}(\Omega; \mathbb{R}^n)$.

Indicator function

For a set A , we define

$$\delta_A(x) := \begin{cases} 0, & x \in A, \\ \infty, & x \notin A. \end{cases}$$

References

1. L. Ambrosio, N. Fusco, and D. Pallara, *Functions of Bounded Variation and Free Discontinuity Problems*, Oxford University Press, 2000.
2. S. R. Arridge, J. P. Kaipio, V. Kolehmainen, and T. Tarvainen, Optical Imaging, in *Handbook of Mathematical Methods in Imaging*, O. Scherzer (ed.), Springer, New York, NY, 2011, 735–780, doi:10.1007/978-0-387-92920-0_17.
3. K. J. Arrow, L. Hurwicz, and H. Uzawa, *Studies in Linear and Non-Linear Programming*, Stanford University Press, 1958.
4. M. Bačák, *Convex Analysis and Optimization in Hadamard Spaces*, Nonlinear Analysis and Applications, De Gruyter, 2014.
5. H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, CMS Books in Mathematics, Springer, 2 edition, 2017, doi:10.1007/978-3-319-48311-5.
6. A. Beck, *First-Order Methods in Optimization*, SIAM, 2017, doi:10.1137/1.9781611974997.
7. A. Beck and M. Teboulle, Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems, *IEEE Transactions on Image Processing* 18 (2009), 2419–2434, doi:10.1109/tip.2009.2028250.
8. 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.
9. R. Begmann, R. Herzog, D. Tenbrink, and J. Vidal-Núñez, Fenchel duality for convex optimization and a primal dual algorithm on Riemannian manifolds, 2019, arXiv:1908.02022.
10. M. Benning, F. Knoll, C. B. Schönlieb, and T. Valkonen, Preconditioned ADMM with nonlinear operator constraint, in *System Modeling and Optimization: 27th IFIP TC 7 Conference, CSMO 2015, Sophia Antipolis, France, June 29–July 3, 2015, Revised Selected Papers*, Springer, 2016, 117–126, doi:10.1007/978-3-319-55795-3_10, arXiv:1511.00425.
11. K. Bredies and H. Sun, Preconditioned Douglas–Rachford splitting methods for convex-concave saddle-point problems, *SIAM Journal on Numerical Analysis* 53 (2015), 421–444, doi:10.1137/140965028.
12. H. Brezis, M. G. Crandall, and A. Pazy, Perturbations of nonlinear maximal monotone sets in Banach space, *Communications on Pure and Applied Mathematics* 23 (1970), 123–144, doi:10.1002/cpa.3160230107.
13. F. E. Browder, Convergence theorems for sequences of nonlinear operators in Banach spaces, *Mathematische Zeitschrift* 100 (1967), 201–225, doi:10.1007/bf01109805.

14. M. Burger, E. Resmerita, and M. Benning, An entropic Landweber method for linear ill-posed problems, 2019, [arXiv:1906.10032](https://arxiv.org/abs/1906.10032).
15. A. Chambolle, R. A. DeVore, N. y. Lee, and B. J. Lucier, Nonlinear wavelet image processing: variational problems, compression, and noise removal through wavelet shrinkage, *IEEE Transactions on Image Processing* 7 (1998), 319–335, [doi:10.1109/83.661182](https://doi.org/10.1109/83.661182).
16. A. Chambolle, M. Ehrhardt, P. Richtárik, and C. Schönlieb, Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications, *SIAM Journal on Optimization* 28 (2018), 2783–2808, [doi:10.1137/17m1134834](https://doi.org/10.1137/17m1134834).
17. 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).
18. 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](https://doi.org/10.1007/s10107-015-0957-3).
19. A. Chambolle and T. Pock, An introduction to continuous optimization for imaging, *Acta Numerica* 25 (2016), 161–319, [doi:10.1017/s096249291600009x](https://doi.org/10.1017/s096249291600009x).
20. P. Chen, J. Huang, and X. Zhang, A primal-dual fixed point algorithm for convex separable minimization with applications to image restoration, *Inverse Problems* 29 (2013), 025011, [doi:10.1088/0266-5611/29/2/025011](https://doi.org/10.1088/0266-5611/29/2/025011).
21. G. Chierchia, E. Chouzenoux, P.L. Combettes, and J.C. Pesquet, The Proximity Operator Repository, 2019, <http://proximity-operator.net>. Online resource.
22. F. Clarke, *Optimization and Nonsmooth Analysis*, Society for Industrial and Applied Mathematics, 1990, [doi:10.1137/1.9781611971309](https://doi.org/10.1137/1.9781611971309).
23. C. Clason, S. Mazurenko, and T. Valkonen, Acceleration and global convergence of a first-order primal-dual method for nonconvex problems, *SIAM Journal on Optimization* 29 (2019), 933–963, [doi:10.1137/18m1170194](https://doi.org/10.1137/18m1170194), [arXiv:1802.03347](https://arxiv.org/abs/1802.03347).
24. C. Clason, S. Mazurenko, and T. Valkonen, Primal-dual proximal splitting and generalized conjugation in nonsmooth nonconvex optimization, *Applied Mathematics and Optimization* (2020), [doi:10.1007/s00245-020-09676-1](https://doi.org/10.1007/s00245-020-09676-1), [arXiv:1901.02746](https://arxiv.org/abs/1901.02746).
25. C. Clason and T. Valkonen, Introduction to Nonsmooth Analysis and Optimization, 2020, [arXiv:2001.00216](https://arxiv.org/abs/2001.00216). Work in progress.
26. L. Condat, A Primal–Dual Splitting Method for Convex Optimization Involving Lipschitzian, Proximable and Linear Composite Terms, *Journal of Optimization Theory and Applications* 158 (2013), 460–479, [doi:10.1007/s10957-012-0245-9](https://doi.org/10.1007/s10957-012-0245-9).
27. 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).
28. M. P. do Carmo, *Riemannian Geometry*, Mathematics: Theory & Applications, Birkhäuser, 2013.
29. 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), 421–439, [doi:10.2307/1993056](https://doi.org/10.2307/1993056).
30. Y. Drori, S. Sabach, and M. Teboulle, A simple algorithm for a class of nonsmooth convex–concave saddle-point problems, *Operations Research Letters* 43 (2015), 209–214, [doi:10.1016/j.orl.2015.02.001](https://doi.org/10.1016/j.orl.2015.02.001).
31. I. Ekeland and R. Temam, *Convex analysis and variational problems*, SIAM, 1999.
32. H. Federer, *Geometric Measure Theory*, Springer, 1969.
33. D. Gabay, Applications of the Method of Multipliers to Variational Inequalities, in *Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems*, M. Fortin and R. Glowinski (eds.), volume 15 of Studies in Mathematics and its Applications, North-Holland, 1983, 299–331.
34. S. Geman and D. Geman, Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 6 (1984), 721–741, [doi:10.1109/tpami.1984.4767596](https://doi.org/10.1109/tpami.1984.4767596).
35. E. Y. Hamedani and N. S. Aybat, A primal-dual algorithm for general convex-concave saddle point problems, 2018, [arXiv:1803.01401](https://arxiv.org/abs/1803.01401).

36. 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.
37. J. B. Hiriart-Urruty and C. Lemaréchal, *Fundamentals of Convex Analysis*, Grundlehren Text Editions, Springer, 2004.
38. T. Hohage and C. Homann, A Generalization of the Chambolle-Pock Algorithm to Banach Spaces with Applications to Inverse Problems, 2014, arXiv:1412.0126.
39. A. Hunt, Weighing without touching: applying electrical capacitance tomography to mass flowrate measurement in multiphase flows, *Measurement and Control* 47 (2014), 19–25, doi:10.1177/0020294013517445.
40. J. Jauhainen, P. Kuusela, A. Seppänen, and T. Valkonen, Relaxed Gauss–Newton methods with applications to electrical impedance tomography, *SIAM Journal on Imaging Sciences* (2020), arXiv:2002.08044. in press.
41. P. Kingsley, Introduction to diffusion tensor imaging mathematics: Parts I-III, *Concepts in Magnetic Resonance Part A* 28 (2006), 101–179, doi:10.1002/cmr.a.20048.
42. P. Kuchment and L. Kunyansky, Mathematics of Photoacoustic and Thermoacoustic Tomography, in *Handbook of Mathematical Methods in Imaging*, O. Scherzer (ed.), Springer, New York, NY, 2011, 817–865, doi:10.1007/978-0-387-92920-0_19.
43. P. Lions and B. Mercier, Splitting algorithms for the sum of two nonlinear operators, *SIAM Journal on Numerical Analysis* 16 (1979), 964–979, doi:10.1137/0716071.
44. A. Lipponen, A. Seppänen, and J. P. Kaipio, Nonstationary approximation error approach to imaging of three-dimensional pipe flow: experimental evaluation, *Measurement Science and Technology* 22 (2011), 104013, doi:10.1088/0957-0233/22/10/104013.
45. 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.
46. M. Lustig, D. Donoho, and J. M. Pauly, Sparse MRI: The application of compressed sensing for rapid MR imaging, *Magnetic Resonance in Medicine* 58 (2007), 1182–1195, doi:10.1002/mrm.21391.
47. Y. Malitsky and M. K. Tam, A forward-backward splitting method for monotone inclusions without cocoercivity, 2018, arXiv:1808.04162.
48. S. Mazurenko, J. Jauhainen, and T. Valkonen, Primal-dual block-proximal splitting for a class of non-convex problems, *Electronic Transactions on Numerical Analysis* (2020), arXiv:1911.06284. accepted.
49. G. J. Minty, On the Maximal domain of a “monotone” function, *The Michigan Mathematical Journal* 8 (1961), 135–137.
50. A. S. Nemirovski and D. Yudin, Problem Complexity and Method Efficiency in Optimization (translated from Russian), *Wiley Interscience Series in Discrete Mathematics* (1983).
51. D. Nishimura, *Principles of Magnetic Resonance Imaging*, Stanford University, 1996.
52. J. M. Ollinger and J. A. Fessler, Positron-emission tomography, *IEEE Signal Processing Magazine* 14 (1997), 43–55, doi:10.1109/79.560323.
53. Z. Opial, Weak convergence of the sequence of successive approximations for nonexpansive mappings, *Bulletin of the American Mathematical Society* 73 (1967), 591–597, doi:10.1090/s0002-9904-1967-11761-0.
54. 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*, IEEE, 2011, 1762–1769, doi:10.1109/iccv.2011.6126441.
55. T. Pock, D. Cremers, H. Bischof, and A. Chambolle, An algorithm for minimizing the Mumford-Shah functional, in *12th IEEE Conference on Computer Vision*, IEEE, 2009, 1133–1140, doi:10.1109/iccv.2009.5459348.
56. R. T. Rockafellar, *Convex Analysis*, Princeton University Press, 1972.
57. R. T. Rockafellar, Monotone operators and the proximal point algorithm, *SIAM Journal on Optimization* 14 (1976), 877–898, doi:10.1137/0314056.
58. L. Rudin, S. Osher, and E. Fatemi, Nonlinear Total Variation based noise removal algorithms, *Physica D* 60 (1992), 259–268.

59. J. Shen and T. F. Chan, Mathematical Models for Local Nontexture Inpaintings, *SIAM Journal on Applied Mathematics* 62 (2002), 1019–1043, doi:10.1137/s0036139900368844.
60. D. Trucu, D. B. Ingham, and D. Lesnic, An inverse coefficient identification problem for the bio-heat equation, *Inverse Problems in Science and Engineering* 17 (2009), 65–83, doi:10.1080/17415970802082880.
61. G. Uhlmann, Electrical impedance tomography and Calderón’s problem, *Inverse Problems* 25 (2009), 123011, doi:10.1088/0266-5611/25/12/123011.
62. T. Valkonen, A primal-dual hybrid gradient method for non-linear operators with applications to MRI, *Inverse Problems* 30 (2014), 055012, doi:10.1088/0266-5611/30/5/055012, arXiv:1309.5032.
63. T. Valkonen, Testing and non-linear preconditioning of the proximal point method, *Applied Mathematics and Optimization* (2018), doi:10.1007/s00245-018-9541-6, arXiv:1703.05705.
64. T. Valkonen, Block-proximal methods with spatially adapted acceleration, *Electronic Transactions on Numerical Analysis* 51 (2019), 15–49, doi:10.1553/etna_vol51s15, arXiv:1609.07373.
65. T. Valkonen, Inertial, corrected, primal-dual proximal splitting, *SIAM Journal on Optimization* 30 (2020), 1391–1420, doi:10.1137/18m1182851, arXiv:1804.08736.
66. T. Valkonen and T. Pock, Acceleration of the PDHGM on partially strongly convex functions, *Journal of Mathematical Imaging and Vision* 59 (2017), 394–414, doi:10.1007/s10851-016-0692-2, arXiv:1511.06566.
67. C. R. Vogel and M. E. Oman, Fast, robust total variation-based reconstruction of noisy, blurred images, *IEEE Transactions on Image Processing* 7 (1998), 813–824, doi:10.1109/83.679423.
68. B. C. Vũ, A splitting algorithm for dual monotone inclusions involving cocoercive operators, *Advances in Computational Mathematics* 38 (2013), 667–681, doi:10.1007/s10444-011-9254-8.
69. X. Zhang, M. Burger, and S. Osher, A unified primal-dual algorithm framework based on Bregman iteration, *Journal of Scientific Computing* 46 (2011), 20–46, doi:10.1007/s10915-010-9408-8.