

# ACCELERATION AND GLOBAL CONVERGENCE OF A FIRST-ORDER PRIMAL–DUAL METHOD FOR NONCONVEX PROBLEMS

Christian Clason\*      Stanislav Mazurenko†      Tuomo Valkonen‡

2018-02-09 (revised 2018-08-07)

**Abstract** The primal–dual hybrid gradient method (PDHGM, also known as the Chambolle–Pock method) has proved very successful for convex optimization problems involving linear operators arising in image processing and inverse problems. In this paper, we analyze an extension to nonconvex problems that arise if the operator is nonlinear. Based on the idea of testing, we derive new step length parameter conditions for the convergence in infinite-dimensional Hilbert spaces and provide acceleration rules for suitably (locally and/or partially) monotone problems. Importantly, we prove linear convergence rates as well as global convergence in certain cases. We demonstrate the efficacy of these step length rules for PDE-constrained optimization problems.

## 1 INTRODUCTION

Many optimization problems can be represented as minimizing a sum of two terms of the form

$$(P) \quad \min_x G(x) + F(K(x))$$

for some (extended) real-valued functionals  $F$  and  $G$  and a (possibly nonlinear) operator  $K$ . For instance, in inverse problems,  $G$  will typically be a fidelity term, measuring fit to data, and  $F \circ K$  is a regularization term introduced to avoid ill-posedness and promote desired features in the solution. In imaging problems in particular, quite often total variation type regularization is used, in which case  $K$  is composed of differential operators [1, 5, 8]. In optimal control,  $K$  frequently denotes the solution operator to partial or ordinary differential equations as a function of the control input. In this case  $G$  and  $F$  stand for control- and state-dependent contributions to the cost function, respectively. The function  $F$  might also account for state constraints [11].

Since the above applications usually involve high and possibly infinite-dimensional spaces, first-order numerical methods can provide the best trade-off between precision and computation time. This, however, depends on the exact formulation of the problem and the specific algorithm

---

\*Faculty of Mathematics, University Duisburg-Essen, 45117 Essen, Germany ([christian.clason@uni-due.de](mailto:christian.clason@uni-due.de))

†Department of Mathematical Sciences, University of Liverpool, United Kingdom ([stan.mazurenko@gmail.com](mailto:stan.mazurenko@gmail.com))

‡ModeMat, Escuela Politécnica Nacional, Quito, Ecuador; *previously* Department of Mathematical Sciences, University of Liverpool, United Kingdom ([tuomo.valkonen@iki.fi](mailto:tuomo.valkonen@iki.fi))

used. Nonsmooth first-order methods roughly divide into two classes: ones based on explicit subgradients, and ones based on proximal maps as introduced in [19]. The former can exhibit very slow convergence, while taking a step in the latter is often tantamount to solving the original problem. As both  $G$  and  $F$  are often convex, introducing a dual variable  $y$  and the convex conjugate  $F^*$  of  $F$ , we can rewrite (P) as

$$(S) \quad \min_x \max_y G(x) + \langle K(x), y \rangle - F^*(y).$$

Now, if we can decouple the primal and dual variables, and, instead of the proximal map of  $x \mapsto G(x) + F(K(x))$ , individually and efficiently compute the proximal maps  $(I + \tau \partial G)^{-1}$  and  $(I + \sigma \partial F^*)^{-1}$ , methods based on proximal steps can be highly efficient. Based on this idea, for linear  $K$  a decoupling algorithm – now commonly known as the Chambolle–Pock method – was suggested in [7, 17]. In [7, 9] the authors proved the  $O(1/N)$  convergence of an ergodic duality gap to zero and provided an  $O(1/N^2)$  acceleration scheme when either the primal or dual objective is strongly convex. In [12], the method was classified as the Primal–Dual Hybrid Gradient method, Modified (PDHGM).

However, frequently in applications,  $K$  is not linear, making (P) nonconvex. This situation is the focus in the present work. Our starting point is the extension of the PDHGM to nonlinear  $K$  suggested in [11, 22], where the authors proved local weak convergence without a rate under a metric regularity assumption. The method, called the NL-PDHGM (for “nonlinear PDHGM”), and its ADMM variants have successfully been applied to problems in magnetic resonance imaging and PDE-constrained optimization [4, 11, 22]. We state it in [Algorithm 1.1](#), also incorporating references to the step length rules of the present work.

**Algorithm 1.1 (NL-PDHGM).** Pick a starting point  $(x^0, y^0)$ . Select step length parameters  $\tau_i, \sigma_i, \omega_i > 0$  according to a suitable rule from one of [Theorems 4.1, 4.3](#) and [4.4](#). The iterate

$$\begin{aligned} x^{i+1} &:= (I + \tau_i \partial G)^{-1}(x^i - \tau_i [\nabla K(x^i)]^* y^i), \\ \bar{x}^{i+1} &:= x^{i+1} + \omega_i (x^{i+1} - x^i), \\ y^{i+1} &:= (I + \sigma_{i+1} \partial F^*)^{-1}(y^i + \sigma_{i+1} K(\bar{x}^{i+1})). \end{aligned}$$

Besides nonconvex ADMM [4, 26] (which is a closely related algorithm), first-order alternatives to the NL-PDHGM include iPiano [15], iPalm [18], and an extension of the PDHGM to semiconvex functions [14]. The former two are inertial variants of forward–backward splitting, with iPalm further splitting the proximal step into two sub-blocks. We stress that none of these can be applied directly to (P) if  $F$  is nonsmooth *and*  $K$  is nonlinear, which is the focus of this work. Another advantage of the approach based on the saddle point formulation (S) which moves all nonconvexity to  $K$  is the following. Consider

$$(1.1) \quad \min_x \frac{1}{2} \|T(x) - z\|^2 + F_0(K_0 x),$$

where  $K$  is linear and  $T$  nonlinear. Such problems arise, e.g., from total variation regularized nonlinear inverse problems, in which case  $K_0 = \nabla$  and  $T$  is a nonlinear forward operator [22]. As

the function  $F_0$  is typically nonsmooth (e.g.,  $F_0 = \|\cdot\|$  for total variation regularization), to apply a simple forward–backward scheme to this problem one would have to compute the proximal map of  $F_0 \circ K_0$ , which is seldom feasible. On the other hand, even if  $T$  were linear, solving the dual problem instead as in [3] will not work either unless  $T$  is unitary. However, we can rewrite (1.1) in the form (S) with  $y = (y_1, y_2)$ ,  $G \equiv 0$ ,  $K(x) := (K_0x, T(x) - z)$ , and  $F^*(y) := F_0^*(y_1) + \frac{1}{2}\|y_2\|^2$ . Now we only need to be able to compute  $K$ ,  $\nabla K$ , and the proximal map of  $F_0^*$ , all of which are typically easy. Observe also how  $F^*$  is strongly convex on the subspace corresponding to the nonlinear part of  $K$ . This will be useful for estimating convergence rates.

In [11], based on small modifications to our original analysis in [22], we showed that the acceleration scheme from [7] for strongly convex problems can also be used with Algorithm 1.1 and nonlinear  $K$  *provided we stop the acceleration at some iteration*. Hence, no convergence rates could be obtained. *In the present paper, based on a completely new and simplified analysis, we provide such rates and show that the acceleration does not have to be stopped*. To the best of our knowledge, this is the first work to prove convergence rates for a primal–dual method for nonsmooth saddle point problems with nonlinear operators. Our new analysis of the NL-PDHGM is based on the “testing” framework introduced in [25] for preconditioned proximal point methods. In particular, we relax the metric regularity required in [22] to mere monotonicity at a solution together with a three-point growth condition on  $K$  around this solution. Both are essentially “nonsmooth” formulations of standard second-order growth conditions. We prove weak convergence to a critical point as well as  $O(1/N^2)$  convergence (which is even global in some situations) with an acceleration rule if  $\partial G$  or  $[\nabla K(x)]^*y$  is strongly monotone at a primal critical point  $\hat{x}$ . If  $\partial F^*$  is also strongly monotone at a dual critical point  $\hat{y}$ , we present step length rules that lead to linear convergence. We emphasize that *all the time we allow  $K$  to be nonlinear, and through this the problem (P) to be globally nonconvex*. In addition, our local monotonicity assumptions are comparable nonsmooth counterparts to standard  $C^2$  and positive Hessian assumptions in smooth nonconvex optimization.

This work is organized as follows. We summarize the “testing” framework introduced in [25] for preconditioned proximal point methods in Section 2. We state our main results in Section 3. Since block-coordinate methods have been receiving more and more attention lately – including in the primal–dual algorithm designed in [23] based on the same testing framework – the main technical derivations of Section 3.2 are implemented in a generalized operator form. Once we have obtained these generic estimates, we devote Section 4 to scalar step length parameters and formulate our main convergence results. These amount to basically standard step length rules for the PDHGM combined with bounds on the initial step lengths. Finally, in Section 5, we illustrate our theoretical results with numerical evidence. We study parameter identification with  $L^1$  fitting and optimal control with state constraints, where the nonlinear operator  $K$  involves the mapping from a potential term in an elliptic partial differential equation to the corresponding solution.

## 2 PROBLEM FORMULATION

Throughout this paper, we write  $\mathbb{L}(X; Y)$  for the space of bounded linear operators between Hilbert spaces  $X$  and  $Y$ . We write  $I$  for the identity operator,  $\langle x, x' \rangle$  for the inner product,

and  $\mathbb{B}(x, r)$  for the closed unit ball of the radius  $r$  at  $x$  in the corresponding space. We set  $\langle x, x' \rangle_T := \langle Tx, x' \rangle$  and  $\|x\|_T := \sqrt{\langle x, x \rangle_T}$ . For  $T, S \in \mathbb{L}(X; Y)$ , the inequality  $T \geq S$  means  $T - S$  is positive semidefinite. Finally,  $\llbracket x_1, x_2 \rrbracket^\alpha := (1 - \alpha)x_1 + \alpha x_2$ ; in particular,  $\bar{x}^{i+1} := \llbracket x^{i+1}, x^i \rrbracket^{-\omega_i}$  in [Algorithm 1.1](#).

We generally assume  $G : X \rightarrow \overline{\mathbb{R}}$  and  $F^* \rightarrow \overline{\mathbb{R}}$  to be convex, proper, and lower semicontinuous, so that their subgradients  $\partial G$  and  $\partial F^*$  are well-defined maximally monotone operators [[2](#), [Theorem 20.25](#)]. Under a constraint qualification, e.g., when  $K$  is  $C^1$  and either the null space of  $[\nabla K(x)]^*$  is trivial or  $\text{dom } F = X$  [[20](#), [Example 10.8](#)], the critical point conditions for (P) and (S) can be written as  $0 \in H(\hat{u})$  for the set-valued operator  $H : X \times Y \rightrightarrows X \times Y$ ,

$$(2.1) \quad H(u) := \begin{pmatrix} \partial G(x) + [\nabla K(x)]^* y \\ \partial F^*(y) - K(x) \end{pmatrix},$$

and  $u = (x, y) \in X \times Y$ . Throughout the paper,  $\hat{u} := (\hat{x}, \hat{y})$  always denotes an arbitrary root  $H$ , which can equivalently be characterized as  $\hat{u} \in H^{-1}(0)$ .

To formulate [Algorithm 1.1](#) in terms suitable for the testing framework of [[25](#)], we define the step length and testing operator

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

respectively, where  $T_i, \Phi_i \in \mathbb{L}(X; X)$  and  $\Sigma_{i+1}, \Psi_{i+1} \in \mathbb{L}(Y; Y)$  are the primal step length and testing operators as well as their dual counterparts.

We also define the nonlinear preconditioner  $M_{i+1}$  and the partial linearization  $\tilde{H}_{i+1}$  of  $H$  by

$$(2.2) \quad M_{i+1} := \begin{pmatrix} I & -T_i[\nabla K(x^i)]^* \\ -\omega_i \Sigma_{i+1} \nabla K(x^i) & I \end{pmatrix}, \quad \text{and}$$

$$(2.3) \quad \tilde{H}_{i+1}(u) := \begin{pmatrix} \partial G(x) + [\nabla K(x^i)]^* y \\ \partial F^*(y) - K(\llbracket x, x^i \rrbracket^{-\omega_i}) - \nabla K(x^i)(x - \llbracket x, x^i \rrbracket^{-\omega_i}) \end{pmatrix}.$$

Note that  $\tilde{H}_{i+1}(u)$  simplifies to  $H(u)$  for linear  $K$ . Now [Algorithm 1.1](#) (which coincides with the “exact” NL-PDHGM of [[22](#)]) can be written as

$$(PP) \quad 0 \in W_{i+1} \tilde{H}_{i+1}(u^{i+1}) + M_{i+1}(u^{i+1} - u^i).$$

(For the “linearized” NL-PDHGM of [[22](#)], we would replace  $\llbracket x, x^i \rrbracket^{-\omega}$  in (2.3) by  $x^i$ .) Following [[25](#)], the step length operator  $W_{i+1}$  in (PP) acts on  $\tilde{H}_{i+1}$  rather than on the step  $u^{i+1} - u^i$  so as to eventually allow zero-length steps on sub-blocks of variables as employed in [[23](#)]. The testing operator  $Z_{i+1}$  does not yet appear in (PP) as it does not feature in the algorithm. We will shortly see that when we apply it to (PP), the product  $Z_{i+1}M_{i+1}$  will form a metric (in the differential-geometric sense) that encodes convergence rates.

Finally, we will also make use of the (possibly empty) subspace  $Y_{\text{NL}}$  of  $Y$  in which  $K$  acts linearly, i.e.,

$$Y_{\text{L}} := \{y \in Y \mid \text{the mapping } x \mapsto \langle y, K(x) \rangle \text{ is linear}\} \quad \text{and} \quad Y_{\text{NL}} := Y_{\text{L}}^\perp.$$

(For examples of such subspaces, we refer to the introduction or, in particular, to [22].) Furthermore,  $P_{\text{NL}}$  will denote the orthogonal projection to  $Y_{\text{NL}}$ . We also write  $\mathbb{B}_{\text{NL}}(\widehat{y}, r) := \{y \in Y \mid \|y - \widehat{y}\|_{P_{\text{NL}}} \leq r\}$  for a closed cylinder in  $Y$  of the radius  $r$  with axis orthogonal to  $Y_{\text{NL}}$ .

Our goal in the rest of the paper is to analyze the convergence of (PP) for the choices (2.1)–(2.3). We will base this analysis on the following abstract “meta-theorem”, which formalizes common steps in convergence proofs of optimization methods. Its purpose is to reduce the proof of convergence to showing that the “iteration gaps”  $\Delta_{i+1}$  – which encode differences in function values and whose specific form depend on the details of the algorithm – are non-positive. The proof of the meta-theorem itself is relatively trivial, being based on telescoping and Pythagoras’ (three-point) formula.

**Theorem 2.1** ([25, Theorem 2.1]). *Suppose (PP) is solvable, and denote the iterates by  $\{u^i\}_{i \in \mathbb{N}}$ . If  $Z_{i+1}M_{i+1}$  is self-adjoint, and for some  $\Delta_{i+1} \in \mathbb{R}$  we have*

$$(CI) \quad \langle \widetilde{H}_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}} \geq \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+2}M_{i+2} - Z_{i+1}M_{i+1}}^2 - \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 - \Delta_{i+1}$$

for all  $i \leq N - 1$  and some  $\widehat{u} \in U$ , then

$$(DI) \quad \frac{1}{2} \|u^N - \widehat{u}\|_{Z_{N+1}M_{N+1}}^2 \leq \frac{1}{2} \|u^0 - \widehat{u}\|_{Z_1M_1}^2 + \sum_{i=0}^{N-1} \Delta_{i+1}.$$

Note that the theorem always holds for *some* choice of the  $\Delta_{i+1} \in \mathbb{R}$ . Our goal will be to choose the step length and testing operators  $T_i, \Sigma_{i+1}, \Phi_i$  and  $\Sigma_{i+1}$  as well as the over-relaxation parameter  $\omega_i$  such that  $\Delta_{i+1} \leq 0$  and – in order to obtain rates –  $Z_{i+1}M_{i+1}$  grows fast as  $i \rightarrow \infty$ . For example, if  $\Delta_{i+1} \leq 0$  and  $Z_{N+1}M_{N+1} \geq \mu_N I$  with  $\mu_N \rightarrow \infty$ , then clearly  $\|u^N - \widehat{u}\|^2 \rightarrow 0$  at the rate  $O(1/\mu_N)$ . In other contexts,  $\Delta_{i+1}$  can be used to encode duality gaps [25] or a penalty on convergence rates due to inexact, stochastic, updates of the local metric  $Z_{i+1}M_{i+1}$  [23].

To motivate the following, consider the “generalized descent inequality” (CI) in the simple case  $\widetilde{H}_{i+1} = H$ . If we now had *at*  $\widehat{u}$  for  $\widehat{w} := 0 \in H(\widehat{u})$  the “operator-relative strong monotonicity”

$$\langle H(u^{i+1}) - \widehat{w}, u^{i+1} - \widehat{u} \rangle_{Z_{i+1}W_{i+1}} \geq \|u^{i+1} - \widehat{u}\|_{Z_{i+1}\Gamma_{i+1}}^2$$

for some suitable operator  $\Gamma_{i+1}$ , then the local metrics should ideally be updated as  $Z_{i+1}M_{i+2} = Z_{i+1}(M_{i+1} + 2\Gamma_{i+1})$ . Part of our work in the following sections is to find such a  $\Gamma_{i+1}$  while maintaining self-adjointness and obtaining fast growth of the metrics. However, our specific choices of  $\widetilde{H}_{i+1}$  and  $M_{i+1}$  switch parts of  $H$  to take the gradient step  $-\nabla K(x^i)^* y^i$  in the primal update and an over-relaxed step in the dual update. We will approximately undo these changes using the term  $-\frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$  in (CI). This component of (CI) can also be related to the “three-point hypomonotonicity”  $\langle \nabla G(x^i) - \nabla G(\widehat{x}), x^{i+1} - \widehat{x} \rangle \geq -\frac{L}{4} \|x^{i+1} - x^i\|^2$  that holds for convex  $G$  with an  $L$ -Lipschitz gradient [25].

Before proceeding with deriving convergence rates using this approach, we show that we can still obtain weak convergence even if  $Z_{N+1}M_{N+1}$  does not grow quickly.

**Proposition 2.2 (weak convergence).** *Suppose the iterates of (PP) satisfy (CI) for some  $\widehat{u} \in H^{-1}(0)$  with  $Z_{i+1}M_{i+1}$  self-adjoint and  $\Delta_{i+1} \leq -\frac{\delta}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$  for some  $\delta > 0$ . Assume that*

(i)  $\varepsilon I \leq Z_{i+1}M_{i+1}$  for some  $\varepsilon > 0$ ;

(ii) for some nonsingular  $W \in \mathbb{L}(U; U)$ ,

$$Z_{i+1}M_{i+1}(u^{i+1} - u^i) \rightarrow 0, \quad u^{i_k} \rightarrow \bar{u} \implies 0 \in WH(\bar{u});$$

(iii) there exists a constant  $C$  such that  $\|Z_i M_i\| \leq C^2$  for all  $i$ , and for any subsequence  $u^{i_k} \rightarrow u$  there exists  $A_\infty \in \mathbb{L}(U; U)$  such that  $Z_{i_k+1}M_{i_k+1}u \rightarrow A_\infty u$  strongly in  $U$  for all  $u \in U$ .

Then  $u^i \rightarrow \bar{u}$  weakly in  $U$  for some  $\bar{u} \in H^{-1}(0)$ .

*Proof.* This is an improvement of [25, Proposition 2.5] that permits nonconstant  $Z_{i+1}M_{i+1}$  and a nonconvex solution set. The proof is based on the corresponding improvement of Opial's lemma (Lemma A.2) together with Theorem 2.1. Using  $\Delta_{i+1} \leq -\frac{\delta}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$ , (DI) applied with  $N = 1$  and  $u^i$  in place of  $u^0$  shows that  $i \mapsto \|u^i - \hat{u}\|_{Z_{i+1}M_{i+1}}^2$  is nonincreasing. This verifies Lemma A.2 (i). Further use of (DI) shows that  $\sum_{i=0}^{\infty} \frac{\delta}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 < \infty$ . Thus  $Z_{i+1}M_{i+1}(u^{i+1} - u^i) \rightarrow 0$ . By (PP) and (ii), any weak limit point  $\bar{u}$  of the  $\{u^i\}_{i \in \mathbb{N}}$  therefore satisfies  $\bar{u} \in H^{-1}(0)$ . This verifies Lemma A.2 (ii) with  $\hat{X} = H^{-1}(0)$ . The remaining assumptions of Lemma A.2 are verified by conditions (i) and (iii), which yields the claim.  $\square$

### 3 ABSTRACT ANALYSIS OF THE NL-PDHGM

We will apply Theorem 2.1 to Algorithm 1.1, for which we have to verify (CI). This inequality always holds for some  $\Delta_{i+1}$ , but for obvious reasons we aim for  $\Delta_{i+1} \leq 0$ . To obtain fast convergence rates, our second goal is to make the metric  $Z_{i+1}M_{i+1}$  grow as quickly as possible; the rate of this growth is constrained through (CI) by the term  $\frac{1}{2}\|u^{i+1} - \hat{u}\|_{Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2}}^2$ . In this section, we therefore reduce (CI) into a few simple conditions on the step length and testing operators. After stating our fundamental assumptions in Section 3.1, we first derive in Section 3.2 explicit (albeit somewhat technical) bounds on the step length operators to ensure (CI). These require that the iterates  $\{u^i\}_{i \in \mathbb{N}}$  stay in a neighborhood of the critical point  $\hat{u}$ . Therefore, in Section 3.3, we provide sufficient conditions for this requirement to hold in the form of additional step length bounds. We will use these conditions in Section 4, where we will derive the actual convergence rates for scalar step lengths.

#### 3.1 FUNDAMENTAL ASSUMPTIONS

In what follows, we will need  $K$  to be locally Lipschitz differentiable.

**Assumption 3.1 (locally Lipschitz  $\nabla K$ ).** The operator  $K : X \rightarrow Y$  is Fréchet differentiable, and for some  $L \geq 0$  and a neighborhood  $\mathcal{X}_K$  of  $\hat{x}$ ,

$$(3.1) \quad \|\nabla K(x) - \nabla K(x')\| \leq L\|x - x'\| \quad (x, x' \in \mathcal{X}_K).$$

**Remark 3.1.** Using Assumption 3.1 and the mean value equality

$$K(x') = K(x) + \nabla K(x)(x' - x) + \int_0^1 (\nabla K(x + s(x' - x)) - \nabla K(x))(x' - x) ds,$$

we obtain for any  $x, x' \in \mathcal{X}_K$  and  $y \in \text{dom } F^*$  the useful inequality

$$(3.2) \quad \langle K(x') - K(x) - \nabla K(x)(x' - x), y \rangle \leq (L/2) \|x - x'\|^2 \|y\|_{P_{\text{NL}}},$$

where the norm in the dual space consists of only the  $Y_{\text{NL}}$  component because by definition, the function  $x \mapsto \langle K(x), y \rangle$  is linear in  $x$  for  $y \in Y_{\text{L}}$ . Consequently, for such  $y$ , the left-hand side of (3.2) is zero.

We also require a form of “local operator-relative strong monotonicity” of the saddle-point mapping  $H$ . Let  $U$  be a Hilbert space, and  $\Gamma \in \mathbb{L}(U; U)$ ,  $\Gamma \geq 0$ . We say that the set-valued map  $H : U \rightrightarrows U$  is  $\Gamma$ -strongly monotone at  $\hat{u}$  for  $\hat{w} \in H(\hat{u})$  if there exists a neighborhood  $\mathcal{U} \ni \hat{u}$  such that

$$(3.3) \quad \langle w - \hat{w}, u - \hat{u} \rangle \geq \|u - \hat{u}\|_{\Gamma}^2, \quad (u \in \mathcal{U}, w \in H(u)).$$

If  $\Gamma = 0$ , we say that  $H$  is monotone at  $\hat{u}$  for  $\hat{w}$ .

In particular, we will assume this monotonicity in terms of  $\partial G$  and  $\partial F^*$ . The idea is that  $G$  and  $F^*$  can have different level of strong convexity on sub-blocks of the variables  $x$  and  $y$ ; we will in particular use this approach to assume strong convexity from  $F^*$  on the subspace  $Y_{\text{NL}}$  only. We will first need the following assumption in [Lemma 3.6](#).

**Assumption 3.2 (monotone  $\partial G$  and  $\partial F^*$ ).** The set-valued map  $\partial G$  is ( $\Gamma_G$ -strongly) monotone at  $\hat{x}$  for  $-\nabla K(\hat{x})^* \hat{y}$  in the neighborhood  $\mathcal{X}_G$  of  $\hat{x}$ , and the set-valued map  $\partial F^*$  is ( $\Gamma_{F^*}$ -strongly) monotone at  $\hat{y}$  for  $K(\hat{x})$  in the neighborhood  $\mathcal{Y}_{F^*}$  of  $\hat{y}$ .

Of course, in view of the assumed convexity of  $G$  and  $F^*$ , [Assumption 3.2](#) is always satisfied with  $\Gamma_G = \Gamma_{F^*} = 0$ .

Our next three-point assumption on  $K$  is central to our analysis. It combines a second-order growth condition with a smoothness estimate, and the operator  $\tilde{\Gamma}_G$  that we now introduce will later be employed as an acceleration factor.

**Assumption 3.3 (three-point condition on  $K$ ).** For given  $\Gamma_G, \tilde{\Gamma}_G \in \mathbb{L}(X; X)$ , neighborhood  $\mathcal{X}_K$  of  $\hat{x}$ , and some  $\Lambda \in \mathbb{L}(X; X)$ ,  $\theta \geq 0$ , and  $p \in [1, 2]$  we have

$$(3.4) \quad \langle [\nabla K(x') - \nabla K(\hat{x})]^* \hat{y}, x - \hat{x} \rangle + \|x - \hat{x}\|_{\Gamma_G - \tilde{\Gamma}_G}^2 \\ \geq \theta \|K(\hat{x}) - K(x) - \nabla K(x)(\hat{x} - x)\|^p - \frac{1}{2} \|x - x'\|_{\Lambda}^2, \quad (x, x' \in \mathcal{X}_K).$$

We typically have that  $0 \leq \tilde{\Gamma}_G \leq \Gamma_G$ . For linear  $K$ , [Assumption 3.3](#) trivially holds for any  $\tilde{\Gamma}_G \leq \Gamma_G$ ,  $\Lambda = 0$ , and  $\theta \geq 0$ . To motivate the assumption in nontrivial cases, consider the following example.

**Example 3.1.** Let  $F^* = \delta_{\{1\}}$  and take  $\tilde{\Gamma}_G = \Gamma_G$  as well as  $K(x) = J(x)$  for some  $J \in C^2(X)$ , which corresponds to the problem  $\min_{x \in X} G(x) + J(x)$  where  $J$  is smooth and possibly nonconvex but  $G$  can be nonsmooth. In this case, both the over-relaxation step and dual update of [Algorithm 1.1](#) are superfluous, and the entire algorithm reduces to conventional forward-backward splitting. If  $x$  and  $x'$  are suitably close to  $\hat{x}$ , Taylor expansion shows that



(3.4) can be expressed as

$$(3.5) \quad \langle x' - \widehat{x}, x - \widehat{x} \rangle_{\nabla^2 J(\widehat{x}')} \geq \theta \|x - \widehat{x}\|_{\nabla^2 J(\widehat{x})}^{2p} - \frac{1}{2} \|x - x'\|_{\Lambda}^2$$

for some  $\widehat{x} = \widehat{x}(x, \widehat{x})$ ,  $\widehat{x}' = \widehat{x}'(x', \widehat{x}) \in X$ . If  $\nabla^2 J(\widehat{x}')$  is positive definite, i.e.  $\nabla^2 J(\widehat{x}') \geq \varepsilon I$  for some  $\varepsilon > 0$ , then writing

$$(3.6) \quad \begin{aligned} \langle x' - \widehat{x}, x - \widehat{x} \rangle_{\nabla^2 J(\widehat{x}')} &= \|x - \widehat{x}\|_{\nabla^2 J(\widehat{x}')}^2 + \langle x' - x, x - \widehat{x} \rangle_{\nabla^2 J(\widehat{x}')} \\ &\geq \|x - \widehat{x}\|_{\nabla^2 J(\widehat{x}')}^2 - (1 - \alpha) \|x - \widehat{x}\|_{\nabla^2 J(\widehat{x}')}^2 \\ &\quad - \frac{1}{4(1 - \alpha)} \|x' - x\|_{\nabla^2 J(\widehat{x}')}^2 \\ &\geq \alpha \varepsilon \|x - \widehat{x}\|^2 - \frac{L}{4(1 - \alpha)} \|x' - x\|^2, \end{aligned}$$

we see that (3.5) holds in some neighborhood  $\mathcal{X}_K$  of  $\widehat{x}$ , for any  $p \in [1, 2]$ ,  $\theta > 0$  small enough, and  $\Lambda > 0$  large enough. The positivity of  $\nabla^2 J(\widehat{x}')$  is guaranteed by the positivity of  $\nabla^2 J(\widehat{x})$  for  $x'$  close to  $\widehat{x}$ . Alternatively, recalling the full expression (3.4), we can use the strong monotonicity of  $\partial G$  at  $\widehat{x}$ . Overall, we therefore require  $\Gamma_G + \nabla^2 J(\widehat{x})$  to be positive, which is a standard condition in nonconvex optimization.

If  $\text{dom } F^*$  is not a singleton, we can apply the reasoning of [Example 3.1](#) to  $J(x) := K(x)^* \widehat{y}$ . The positivity of  $\Gamma_G + \nabla^2 (K(\cdot)^* \widehat{y})(\widehat{x})$  then amounts to a second-order optimality condition on the solution  $\widehat{x}$  to the problem  $\min_x G(x) + \langle K(x), \widehat{y} \rangle$ . Indeed, we can verify [Assumption 3.3](#) simply based on the monotonicity of  $\partial G + \nabla K(\cdot)^* \widehat{y}$  at  $\widehat{x}$ .

**Proposition 3.2.** *Suppose [Assumption 3.1](#) (locally Lipschitz  $\nabla K$ ) and [Assumption 3.2](#) (monotone  $\partial G$  and  $\partial F^*$ ) hold and for some  $\gamma_x > 0$ ,*

$$(3.7) \quad \|x - \widehat{x}\|_{\Gamma_G - \widetilde{\Gamma}_G}^2 + \langle (\nabla K(x) - \nabla K(\widehat{x}))(x - \widehat{x}), \widehat{y} \rangle \geq \gamma_x \|x - \widehat{x}\|^2 \quad (\forall x \in \mathcal{X}_K).$$

*Then [Assumption 3.3](#) holds with  $p = 1$ ,  $\theta = 2(\gamma_x - \xi)L^{-1}$ , and  $\Lambda = L^2 \|P_{\text{NL}} \widehat{y}\|^2 (2\xi)^{-1} I$  for any  $\xi \in (0, \gamma_x]$ .*

*Proof.* An application of Cauchy's inequality, [Assumption 3.1](#), and (3.7) yields for any  $\xi > 0$  the estimate

$$\begin{aligned} \langle [\nabla K(x') - \nabla K(\widehat{x})]^* \widehat{y}, x - \widehat{x} \rangle + \|x - \widehat{x}\|_{\Gamma_G - \widetilde{\Gamma}_G}^2 &= \langle [\nabla K(x) - \nabla K(\widehat{x})]^* \widehat{y}, x - \widehat{x} \rangle + \|x - \widehat{x}\|_{\Gamma_G - \widetilde{\Gamma}_G}^2 \\ &\quad + \langle (\nabla K(x') - \nabla K(x))(x - \widehat{x}), \widehat{y} \rangle \\ &\geq (\gamma_x - \xi) \|x - \widehat{x}\|^2 - L^2 \|P_{\text{NL}} \widehat{y}\|^2 (4\xi)^{-1} \|x' - x\|^2. \end{aligned}$$

At the same time, using (3.1) and the reasoning of (3.2),  $\|K(\widehat{x}) - K(x) - \nabla K(x)(\widehat{x} - x)\| \leq (L/2) \|x - \widehat{x}\|^2$ . So [Assumption 3.3](#) holds if we take  $p = 1$ ,  $\theta \leq 2(\gamma_x - \xi)/L$ , and  $\Lambda = L^2 \|P_{\text{NL}} \widehat{y}\|^2 (2\xi)^{-1} I$ .  $\square$

**Remark 3.3.** *Observe that if  $\Gamma_G - \widetilde{\Gamma}_G \geq \varepsilon I$  for some  $\varepsilon > L \|P_{\text{NL}} \widehat{y}\|$ , then [Assumption 3.1](#) (locally Lipschitz  $\nabla K$ ) guarantees (3.7) for  $\gamma_x = \varepsilon - L \|P_{\text{NL}} \widehat{y}\|$ . This requires  $\|P_{\text{NL}} \widehat{y}\|$  to be small, which was*



a central assumption in [22] that we intend to avoid in the present work. Also note that if  $\langle K(\cdot), \widehat{y} \rangle$  is convex, then (3.7) holds for  $\gamma_x = \varepsilon$ , so estimating  $\gamma_x$  based on the Lipschitz continuity of  $\nabla K$  alone provides a too conservative estimate.

More generally, while based on our discussion above the satisfaction of (3.3) seems reasonable to expect, its verification can demand some effort. To demonstrate that the condition can be satisfied, we verify this in [Appendix B](#) for a simple example of reconstructing the phase and amplitude of a complex number from a noisy measurement.

Combining [Assumptions 3.1 to 3.3](#), we assume throughout the rest of the paper that for some  $\rho_y \geq 0$ , the corresponding neighborhood

$$(3.8) \quad \mathcal{U}(\rho_y) := (\mathcal{X}_G \cap \mathcal{X}_K) \times (\mathbb{B}_{\text{NL}}(\widehat{y}, \rho_y) \cap \mathcal{Y}_{F^*})$$

of  $\widehat{u}$  is nonempty.

### 3.2 GENERAL ESTIMATES

We verify the conditions of [Theorem 2.1](#) in several steps. First, we ensure that the operator  $Z_{i+1}M_{i+1}$  giving rise to the local metric is self-adjoint. Then we show that  $Z_{i+2}M_{i+2}$  and the update  $Z_{i+1}(M_{i+1} + \Xi_{i+1})$  performed by the algorithm yield identical norms, where  $\Xi_{i+1}$  represents some off-diagonal components from the algorithm as well as any strong monotonicity. Finally, we estimate  $\widetilde{H}_{i+1}(u)$  in order to verify [\(CI\)](#).

We require for some  $\kappa \in [0, 1)$ ,  $\eta_i > 0$ ,  $\widetilde{\Gamma}_G \in \mathbb{L}(X; X)$ , and  $\widetilde{\Gamma}_{F^*} \in \mathbb{L}(Y; Y)$  the following relationships:

$$(3.9a) \quad \omega_i := \eta_i / \eta_{i+1},$$

$$\Psi_i \Sigma_i = \eta_i I,$$

$$(3.9b) \quad \Phi_i T_i = \eta_i I,$$

$$(1 - \kappa) \Psi_{i+1} \geq \eta_i^2 \nabla K(x^i) \Phi_i^{-1} [\nabla K(x^i)]^*,$$

$$(3.9c) \quad \Phi_i = \Phi_i^* \geq 0,$$

$$\Psi_{i+1} = \Psi_{i+1}^* \geq 0,$$

$$(3.9d) \quad \Phi_{i+1} = \Phi_i (1 + 2T_i \widetilde{\Gamma}_G),$$

$$\Psi_{i+2} = \Psi_{i+1} (1 + 2\Sigma_{i+1} \widetilde{\Gamma}_{F^*}).$$

In [Section 4](#), we will verify these relationships for specific scalar step length rules in [Algorithm 1.1](#).

**Lemma 3.4.** *Fix  $i \in \mathbb{N}$  and suppose (3.9) holds. Then  $Z_{i+1}M_{i+1}$  is self-adjoint and satisfies*

$$Z_{i+1}M_{i+1} \geq \begin{pmatrix} \delta \Phi_i & 0 \\ 0 & (\kappa - \delta)(1 - \delta)^{-1} \Psi_{i+1} \end{pmatrix} \quad \text{for any } \delta \in [0, \kappa].$$

*Proof.* From (2.2) and (3.9), we have  $\Phi_i T_i = \eta_i I$  and  $\Psi_{i+1} \Sigma_{i+1} \omega_i = \eta_i I$ . Hence

$$(3.10) \quad Z_{i+1}M_{i+1} = \begin{pmatrix} \Phi_i & -\eta_i [\nabla K(x^i)]^* \\ -\eta_i \nabla K(x^i) & \Psi_{i+1} \end{pmatrix},$$

and therefore  $Z_{i+1}M_{i+1}$  is self-adjoint. Cauchy's inequality furthermore implies that

$$(3.11) \quad Z_{i+1}M_{i+1} \geq \begin{pmatrix} \delta \Phi_i & 0 \\ 0 & \Psi_{i+1} - \frac{\eta_i^2}{1 - \delta} \nabla K(x^i) \Phi_i^{-1} [\nabla K(x^i)]^* \end{pmatrix}.$$

Now (3.9) ensures the remaining part of the statement.  $\square$

Our next step is to simplify  $Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2}$  in (CI) while keeping the option to accelerate the method when some of the blocks of  $H$  exhibit strong monotonicity.

**Lemma 3.5.** *Fix  $i \in \mathbb{N}$ , and suppose (3.9) holds. Then  $\frac{1}{2} \|\cdot\|_{Z_{i+1}(M_{i+1} + \Xi_{i+1}) - Z_{i+2}M_{i+2}}^2 = 0$  for*

$$(3.12) \quad \Xi_{i+1} := \begin{pmatrix} 2T_i \tilde{\Gamma}_G & 2T_i [\nabla K(x^i)]^* \\ -2\Sigma_{i+1} \nabla K(x^{i+1}) & 2\Sigma_{i+1} \tilde{\Gamma}_{F^*} \end{pmatrix}.$$

*Proof.* Using (3.9) and (3.10) can write

$$Z_{i+1}(M_{i+1} + \Xi_{i+1}) - Z_{i+2}M_{i+2} = \begin{pmatrix} 0 & [\eta_{i+1} \nabla K(x^{i+1}) + \eta_i \nabla K(x^i)]^* \\ -\eta_{i+1} \nabla K(x^{i+1}) - \eta_i \nabla K(x^i) & 0 \end{pmatrix}.$$

Inserting this into the definition of the weighted norm yields the claim.  $\square$

The next somewhat technical lemma estimates the linearizations of  $\tilde{H}_{i+1}$  that are needed to make the abstract algorithm (PP) computable for nonlinear  $K$ .

**Lemma 3.6.** *Suppose Assumption 3.1 (locally Lipschitz  $\nabla K$ ), Assumption 3.2 (monotone  $\partial G$  and  $\partial F^*$ ), and (3.9) hold. For a fixed  $i \in \mathbb{N}$ , let  $\bar{x}^{i+1} \in X_K$  and let  $\rho_y \geq 0$  be such that  $u^i, u^{i+1} \in \mathcal{U}(\rho_y)$ . Also suppose Assumption 3.3 (three-point condition on  $K$ ) holds with  $\theta \geq \rho_y^{2-p} p^{-p} \omega_i^{-1} \zeta^{1-p}$  for some  $\zeta > 0$  and  $p \in [1, 2]$ . Then*

$$\begin{aligned} \langle \tilde{H}_{i+1}(u^{i+1}), u^{i+1} - \hat{u} \rangle_{Z_{i+1}M_{i+1}} - \frac{1}{2} \|u^{i+1} - \hat{u}\|_{Z_{i+1}\Xi_{i+1}}^2 \\ \geq \|y^{i+1} - \hat{y}\|_{\eta_{i+1}[\Gamma_{F^*} - \tilde{\Gamma}_{F^*} - (p-1)\zeta P_{NL}]}^2 - \frac{1}{2} \|x^{i+1} - x^i\|_{\eta_i[\Lambda + L(2+\omega_i)\rho_y I]}^2. \end{aligned}$$

*Proof.* From (2.3), (3.9), and (3.12), we have

$$(3.13) \quad \begin{aligned} D &:= \langle \tilde{H}_{i+1}(u^{i+1}), u^{i+1} - \hat{u} \rangle_{Z_{i+1}M_{i+1}} - \frac{1}{2} \|u^{i+1} - \hat{u}\|_{Z_{i+1}\Xi_{i+1}}^2 \\ &= \langle H(u^{i+1}), u^{i+1} - \hat{u} \rangle_{Z_{i+1}W_{i+1}} \\ &\quad + \eta_i \langle [\nabla K(x^i) - \nabla K(x^{i+1})](x^{i+1} - \hat{x}), y^{i+1} \rangle \\ &\quad + \eta_{i+1} \langle K(x^{i+1}) - K(\bar{x}^{i+1}) - \nabla K(x^i)(x^{i+1} - \bar{x}^{i+1}), y^{i+1} - \hat{y} \rangle \\ &\quad + \langle (\eta_{i+1} \nabla K(x^{i+1}) - \eta_i \nabla K(x^i))(x^{i+1} - \hat{x}), y^{i+1} - \hat{y} \rangle \\ &\quad - \eta_i \|x^{i+1} - \hat{x}\|_{\Gamma_G}^2 - \eta_{i+1} \|y^{i+1} - \hat{y}\|_{\Gamma_{F^*}}^2. \end{aligned}$$

Since  $0 \in H(\hat{u})$ , we have  $z_G := -[\nabla K(\hat{x})]^* \hat{y} \in \partial G(\hat{x})$  and  $z_{F^*} := K(\hat{x}) \in \partial F^*(\hat{y})$ . Using (3.9), we can therefore expand

$$\begin{aligned} \langle H(u^{i+1}), u^{i+1} - \hat{u} \rangle_{Z_{i+1}W_{i+1}} &= \eta_i \langle \partial G(x^{i+1}) - z_G, x^{i+1} - \hat{x} \rangle + \eta_{i+1} \langle \partial F^*(y^{i+1}) - z_{F^*}, y^{i+1} - \hat{y} \rangle \\ &\quad + \eta_i \langle [\nabla K(x^{i+1})]^* y^{i+1} - [\nabla K(\hat{x})]^* \hat{y}, x^{i+1} - \hat{x} \rangle \\ &\quad + \eta_{i+1} \langle K(\hat{x}) - K(x^{i+1}), y^{i+1} - \hat{y} \rangle. \end{aligned}$$

Using the local (strong) monotonicity of  $G$  and  $F^*$  (Assumption 3.2) and rearranging terms, we obtain

$$(3.14) \quad \begin{aligned} \langle H(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}W_{i+1}} &\geq \eta_i \|x^{i+1} - \widehat{x}\|_{\Gamma_G}^2 + \eta_{i+1} \|y^{i+1} - \widehat{y}\|_{\Gamma_{F^*}}^2 \\ &\quad + \eta_i \langle \nabla K(x^{i+1})(x^{i+1} - \widehat{x}), y^{i+1} \rangle \\ &\quad - \eta_i \langle \nabla K(\widehat{x})(x^{i+1} - \widehat{x}), \widehat{y} \rangle \\ &\quad + \eta_{i+1} \langle K(\widehat{x}) - K(x^{i+1}), y^{i+1} - \widehat{y} \rangle. \end{aligned}$$

Now we plug the estimate (3.14) into (3.13) and rearrange to arrive at

$$\begin{aligned} D &\geq \eta_i \|x^{i+1} - \widehat{x}\|_{\Gamma_G - \widetilde{\Gamma}_G}^2 + \eta_{i+1} \|y^{i+1} - \widehat{y}\|_{\Gamma_{F^*} - \widetilde{\Gamma}_{F^*}}^2 \\ &\quad - \eta_i \langle \nabla K(\widehat{x})(x^{i+1} - \widehat{x}), \widehat{y} \rangle + \eta_i \langle \nabla K(x^i)(x^{i+1} - \widehat{x}), y^{i+1} \rangle \\ &\quad + \eta_{i+1} \langle K(\widehat{x}) - K(\bar{x}^{i+1}) - \nabla K(x^i)(x^{i+1} - \bar{x}^{i+1}), y^{i+1} - \widehat{y} \rangle \\ &\quad + \langle (\eta_{i+1} \nabla K(x^{i+1}) - \eta_i \nabla K(x^i))(x^{i+1} - \widehat{x}), y^{i+1} - \widehat{y} \rangle \\ &= \eta_i \|x^{i+1} - \widehat{x}\|_{\Gamma_G - \widetilde{\Gamma}_G}^2 + \eta_{i+1} \|y^{i+1} - \widehat{y}\|_{\Gamma_{F^*} - \widetilde{\Gamma}_{F^*}}^2 \\ &\quad + \eta_i \langle [\nabla K(x^i) - \nabla K(\widehat{x})](x^{i+1} - \widehat{x}), \widehat{y} \rangle \\ &\quad + \eta_{i+1} \langle K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1}), y^{i+1} - \widehat{y} \rangle \\ &\quad + \eta_{i+1} \langle K(x^{i+1}) - K(\bar{x}^{i+1}) + \nabla K(x^{i+1})(\bar{x}^{i+1} - x^{i+1}), y^{i+1} - \widehat{y} \rangle \\ &\quad + \eta_{i+1} \langle (\nabla K(x^i) - \nabla K(x^{i+1}))(\bar{x}^{i+1} - x^{i+1}), y^{i+1} - \widehat{y} \rangle. \end{aligned}$$

Applying Assumption 3.1, (3.2), and  $\bar{x}^{i+1} - x^{i+1} = \omega_i(x^{i+1} - x^i)$  to the last two terms, we obtain

$$\langle K(x^{i+1}) - K(\bar{x}^{i+1}) + \nabla K(x^{i+1})(\bar{x}^{i+1} - x^{i+1}), y^{i+1} - \widehat{y} \rangle \geq -\frac{L\omega_i^2}{2} \|x^{i+1} - x^i\|^2 \|y^{i+1} - \widehat{y}\|_{P_{NL}}$$

and

$$\langle (\nabla K(x^i) - \nabla K(x^{i+1}))(\bar{x}^{i+1} - x^{i+1}), y^{i+1} - \widehat{y} \rangle \geq -L\omega_i \|x^{i+1} - x^i\|^2 \|y^{i+1} - \widehat{y}\|_{P_{NL}}.$$

These estimates together with (3.9) and  $u^{i+1} \in \mathcal{U}(\rho_y)$  now imply that

$$(3.15) \quad D \geq \eta_i D_{i+1}^K + \eta_{i+1} \|y^{i+1} - \widehat{y}\|_{\Gamma_{F^*} - \widetilde{\Gamma}_{F^*}}^2$$

for

$$\begin{aligned} D_{i+1}^K &:= \langle [\nabla K(x^i) - \nabla K(\widehat{x})](x^{i+1} - \widehat{x}), \widehat{y} \rangle + \|x^{i+1} - \widehat{x}\|_{\Gamma_G - \widetilde{\Gamma}_G}^2 - L(1 + \omega_i/2)\rho_y \|x^{i+1} - x^i\|^2 \\ &\quad - \|y^{i+1} - \widehat{y}\|_{P_{NL}} \|K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1})\|/\omega_i. \end{aligned}$$

Finally, we use Assumption 3.3 to estimate

$$(3.16) \quad \begin{aligned} D_{i+1}^K &\geq \theta \|K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1})\|^p - \frac{1}{2} \|x^{i+1} - x^i\|_{\Lambda + L(2 + \omega_i)\rho_y I}^2 \\ &\quad - \|y^{i+1} - \widehat{y}\|_{P_{NL}} \|K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1})\|/\omega_i. \end{aligned}$$

We now use the following Young's inequality for any positive  $a, b, p$  and  $q$  such that  $q^{-1} + p^{-1} = 1$ :

$$ab = \left(ab^{\frac{2-p}{p}}\right) b^{\frac{2-p-1}{p}} \leq \frac{1}{p} \left(ab^{\frac{2-p}{p}}\right)^p + \frac{1}{q} b^{2\frac{p-1}{p}q} = \frac{1}{p} a^p b^{2-p} + \left(1 - \frac{1}{p}\right) b^2.$$

With this inequality applied to the last term of (3.16) for

$$a = (\zeta p)^{-1/2} \|K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1})\|, \quad b = (\zeta p)^{1/2} \|y^{i+1} - \widehat{y}\|_{P_{\text{NL}}},$$

and any  $\zeta > 0$ , we arrive at the estimate

$$\begin{aligned} D_{i+1}^K &\geq \theta \|K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1})\|^p - \frac{1}{2} \|x^{i+1} - x^i\|_{\Lambda+L(2+\omega_i)\rho_y I}^2 \\ &\quad - \frac{\|y^{i+1} - \widehat{y}\|_{P_{\text{NL}}}^{2-p}}{p^p \omega_i \zeta^{p-1}} \|K(\widehat{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\widehat{x} - x^{i+1})\|^p - \frac{p-1}{\omega_i} \zeta \|y^{i+1} - \widehat{y}\|_{P_{\text{NL}}}^2. \end{aligned}$$

Now observe that  $\theta - \|y^{i+1} - \widehat{y}\|_{P_{\text{NL}}}^{2-p} (p^p \omega_i \zeta^{p-1})^{-1} \geq \theta - \rho_y^{2-p} (p^p \omega_i \zeta^{p-1})^{-1} \geq 0$ . Therefore

$$(3.17) \quad D_{i+1}^K \geq -\frac{1}{2} \|x^{i+1} - x^i\|_{\Lambda+L(2+\omega_i)\rho_y I}^2 - \frac{p-1}{\omega_i} \zeta \|y^{i+1} - \widehat{y}\|_{P_{\text{NL}}}^2.$$

Combining this with (3.15) we finally obtain

$$D \geq \|y^{i+1} - \widehat{y}\|_{\eta_{i+1}[\Gamma_{F^*} - \widetilde{\Gamma}_{F^*} - (p-1)\zeta P_{\text{NL}}]}^2 - \frac{1}{2} \|x^{i+1} - x^i\|_{\eta_i[\Lambda+L(2+\omega_i)\rho_y I]}^2,$$

which was our claim.  $\square$

We now have all the necessary tools in hand to formulate the main estimate.

**Theorem 3.7.** *Fix  $i \in \mathbb{N}$ , and suppose (3.9) and Assumption 3.1 (locally Lipschitz  $\nabla K$ ), Assumption 3.2 (monotone  $\partial G$  and  $\partial F^*$ ), and Assumption 3.3 (three-point condition on  $K$ ) hold. Also suppose  $\bar{x}^{i+1} \in \mathcal{X}_K$  and that  $u^i, u^{i+1} \in \mathcal{U}(\rho_y)$  for some  $\rho_y \geq 0$ . Furthermore, for  $0 \leq \delta \leq \kappa < 1$  define*

$$S_{i+1} := \begin{pmatrix} \delta \Phi_i - \eta_i[\Lambda+L(2+\omega_i)\rho_y I] & 0 \\ 0 & \Psi_{i+1} - \frac{\eta_i^2}{1-\kappa} \nabla K(x^i) \Phi_i^{-1} [\nabla K(x^i)]^* \end{pmatrix}.$$

Finally, suppose Assumption 3.3 holds with  $\theta \geq \rho_y^{2-p} p^{-p} \omega_i^{-1} \zeta^{1-p}$  for some  $\zeta > 0$ . Then (CI) is satisfied (for this  $i$ ) if

$$(3.18) \quad \frac{1}{2} \|u^{i+1} - u^i\|_{S_{i+1}}^2 + \|y^{i+1} - \widehat{y}\|_{\eta_{i+1}[\Gamma_{F^*} - \widetilde{\Gamma}_{F^*} - (p-1)\zeta P_{\text{NL}}]}^2 \geq -\Delta_{i+1}.$$

In particular, under the above assumptions, we may take  $\Delta_{i+1} = 0$  in (CI) provided

$$(3.19a) \quad \Phi_i \geq \eta_i \delta^{-1} [\Lambda + L(2 + \omega_i)\rho_y I],$$

$$(3.19b) \quad \Psi_{i+1} \geq \frac{\eta_i^2}{1-\kappa} \nabla K(x^i) \Phi_i^{-1} [\nabla K(x^i)]^*, \quad \text{and}$$

$$(3.19c) \quad \Gamma_{F^*} \geq \widetilde{\Gamma}_{F^*} + (p-1)\zeta P_{\text{NL}}.$$

*Proof.* Using the definition of  $S_{i+1}$  and (3.10), we have that

$$\frac{1}{2} \|u^{i+1} - u^i\|_{S_{i+1}}^2 \leq \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 - \frac{1}{2} \|x^{i+1} - x^i\|_{\eta_i[\Lambda+L(2+\omega_i)\rho_y I]}^2.$$

Since [Assumption 3.3](#) holds with  $\theta \geq \rho_y^{2-p} p^{-p} \omega_i^{-1} \zeta^{1-p}$  for some  $\zeta > 0$ , we can apply [Lemma 3.6](#) to further bound

$$\begin{aligned} & \frac{1}{2} \|u^{i+1} - u^i\|_{S_{i+1}}^2 + \|y^{i+1} - \widehat{y}\|_{\eta_{i+1}[\Gamma_{F^*} - \widetilde{\Gamma}_{F^*} - (p-1)\zeta P_{\text{NL}}]}^2 \\ & \leq \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \langle \widetilde{H}_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}M_{i+1}} - \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}\Xi_{i+1}}^2. \end{aligned}$$

Using [Lemma 3.5](#), we may insert  $\|u^{i+1} - \widehat{u}\|_{Z_{i+1}\Xi_{i+1}}^2 = \|u^{i+1} - \widehat{u}\|_{Z_{i+2}M_{i+2} - Z_{i+1}M_{i+1}}^2$  and use (3.18) to obtain

$$\begin{aligned} -\Delta_{i+1} & \leq \frac{1}{2} \|u^{i+1} - u^i\|_{S_{i+1}}^2 + \|y^{i+1} - \widehat{y}\|_{\eta_{i+1}[\Gamma_{F^*} - \widetilde{\Gamma}_{F^*} - (p-1)\zeta P_{\text{NL}}]}^2 \\ & \leq \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \langle \widetilde{H}_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}M_{i+1}} \\ & \quad - \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+2}M_{i+2} - Z_{i+1}M_{i+1}}^2. \end{aligned}$$

Rearranging the terms, we arrive at

$$\langle \widetilde{H}_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}M_{i+1}} \geq \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+2}M_{i+2} - Z_{i+1}M_{i+1}}^2 - \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 - \Delta_{i+1}.$$

Hence, (CI) is satisfied.

Finally, if in addition the relations (3.19) are satisfied, then the left-hand side of (3.18) is trivially bounded from below by zero.  $\square$

We close this section with some remarks on the conditions (3.19):

- While (3.19a) and (3.19b) are stated in terms of  $\Phi_i$  and  $\Psi_{i+1}$ , they actually provide bounds on the step length operators  $T_i$  and  $\Sigma_i$ : Since  $\eta_i I = \Phi_i T_i = \Psi_i \Sigma_i$  by (3.9),  $\Phi_i$  and  $\Psi_{i+1}$  can be eliminated from (3.19), and we will do so for scalar step lengths in [Section 4](#). Thus, while (3.9) provides valid update rules for the parameters in [Algorithm 1.1](#), (3.19) will provide upper bounds on step lengths under which convergence can be proven.
- If  $K$  is linear, (3.19a) reduces to  $\Phi_i \geq 0$  since  $P_{\text{NL}} = 0$  and hence  $\rho_y = 0$  and  $\Lambda = 0$ . We can thus take  $\kappa = 0$ , so that (3.19b) turns into an operator analogue of the step length bound  $\tau_i \sigma_i \|K\|^2 < 1$  of [7].
- Recall from (3.8) that  $\rho_y$  only bounds the dual variable on the subspace  $Y_{\text{NL}}$ . Therefore, most of the requirements for convergence introduced in [Section 4.2](#) to account for non-linear  $K$  (e.g., upper bounds on the primal step length, initialization of the dual variable close to a critical point, or the strong convexity of  $F^*$  at  $\widehat{y}$ ) will only be required with respect to  $Y_{\text{NL}}$ .
- Comparing (3.19) with the requirements of [22], a crucial difference is that in (3.19c),  $\Gamma_{F^*}$  is allowed to be zero when  $p = 1$  and hence we do not require strong convexity from  $F^*$ ; see also [11]. In fact, for  $p = 1$  the inequality on  $\theta$  in [Theorem 3.7](#) reduces to  $\rho_y \leq \omega_i \theta$ . We therefore only need to ensure that the dual variable is initialized close to a critical point

within the subspace  $Y_{\text{NL}}$ . If  $p \in (1, 2]$ , the factor  $\theta$  imposes a lower bound on the dual factor of strong monotonicity over  $Y_{\text{NL}}$ . Indeed, the minimal  $\zeta$  allowed in [Theorem 3.7](#) is given by  $\zeta = \rho_y^{(2-p)/(p-1)} (p^p \omega_i \theta)^{-1/(p-1)}$ , and [\(3.19c\)](#) requires that the factor of strong convexity of  $F^*$  at  $\widehat{y}$  with respect to the subspace  $Y_{\text{NL}}$  is not smaller than this  $\zeta$ .

- Finally, while [\(3.19\)](#) says nothing about  $\Gamma_G$  or  $\widetilde{\Gamma}_G$ , the discussion in and after [Example 3.1](#) indicates that the solution  $\widehat{x}$  to  $G(\cdot) + \langle K(\cdot), \widehat{y} \rangle$  should satisfy a “nonsmooth” second-order growth condition to compensate for the nonlinearity of  $K$ . Therefore,  $\widetilde{\Gamma}_G$  is implicitly bounded from above by the strong convexity factor of the primal problem in [Assumption 3.3](#).

### 3.3 LOCAL STEP LENGTH BOUNDS

In order to apply [Lemma 3.6](#) and therefore [Theorem 3.7](#), we need one final technical result to ensure that the new iterates  $u^{i+1}$  remain in the local neighborhood  $\mathcal{U}(\rho_y)$  of  $\widehat{u}$ . The following lemma provides the basis from which we further work in [Section 4.3](#) and puts a limit on how far the next iterate can escape from a given neighborhood of  $\widehat{u}$  in terms of bounds on the step lengths.

**Lemma 3.8.** *Fix  $i \in \mathbb{N}$ . Suppose [Assumption 3.1](#) (locally Lipschitz  $\nabla K$ ) holds and  $u^{i+1}$  solves (PP). For simplicity, assume  $\omega_i \leq 1$ . For some  $r_{x,i}, r_y > 0$ , and  $\delta_{x,i}, \delta_y \geq 0$ , let  $\mathbb{B}(\widehat{x}, r_{x,i} + \delta_{x,i}) \subset \mathcal{X}_K$ ,  $x^i \in \mathbb{B}(\widehat{x}, r_{x,i})$ , and  $y^i \in \mathbb{B}(\widehat{y}, r_y)$ . If*

$$(3.20) \quad \|T_i\| \leq \frac{\delta_{x,i}/2}{\|\nabla K(x^i)\| r_y + L \|P_{\text{NL}} \widehat{y}\| r_{x,i}}, \quad \text{and} \quad \|\Sigma_{i+1}\| \leq \frac{2\delta_y (r_{x,i} + \delta_{x,i})^{-1}}{L(r_{x,i} + \delta_{x,i}) + 2\|\nabla K(\widehat{x})\|},$$

then  $x^{i+1}, \bar{x}^{i+1} \in \mathbb{B}(\widehat{x}, r_{x,i} + \delta_{x,i})$  and  $y^{i+1} \in \mathbb{B}(\widehat{y}, r_y + \delta_y)$ .

*Proof.* We want to show that the step length conditions [\(3.20\)](#) imply that

$$\|x^{i+1} - \widehat{x}\| \leq r_{x,i} + \delta_{x,i}, \quad \|\bar{x}^{i+1} - \widehat{x}\| \leq r_{x,i} + \delta_{x,i}, \quad \text{and} \quad \|y^{i+1} - \widehat{y}\| \leq r_y + \delta_y.$$

We do this by applying the testing argument on the primal and dual variables separately. Multiplying (PP) by  $Z_{i+1}^*(u^{i+1} - \widehat{u})$  with  $\Phi_i = I$  and  $\Psi_{i+1} = 0$ , we get

$$0 \in \langle \partial G(x^{i+1}) + [\nabla K(x^i)]^* y^i, x^{i+1} - \widehat{x} \rangle_{T_i} + \langle x^{i+1} - x^i, x^{i+1} - \widehat{x} \rangle.$$

Using the three-point version of Pythagoras’ identity,

$$(3.21) \quad \langle x^{i+1} - x^i, x^{i+1} - \widehat{x} \rangle = \frac{1}{2} \|x^{i+1} - x^i\|^2 - \frac{1}{2} \|x^i - \widehat{x}\|^2 + \frac{1}{2} \|x^{i+1} - \widehat{x}\|^2,$$

we obtain

$$\|x^i - \widehat{x}\|^2 \in 2\langle \partial G(x^{i+1}) + [\nabla K(x^i)]^* y^i, x^{i+1} - \widehat{x} \rangle_{T_i} + \|x^{i+1} - x^i\|^2 + \|x^{i+1} - \widehat{x}\|^2.$$

Using  $0 \in \partial G(\widehat{x}) + [\nabla K(\widehat{x})]^* \widehat{y}$  and the monotonicity of  $\partial G$ , we then arrive at

$$\|x^{i+1} - x^i\|^2 + \|x^{i+1} - \widehat{x}\|^2 + 2\langle [\nabla K(x^i)]^* y^i - [\nabla K(\widehat{x})]^* \widehat{y}, x^{i+1} - \widehat{x} \rangle_{T_i} \leq \|x^i - \widehat{x}\|^2.$$

With  $C_x := \|\nabla K(x^i)^* y^i - \nabla K(\widehat{x})^* \widehat{y}\|_{T_i^2}$ , this implies that

$$(3.22) \quad \|x^{i+1} - x^i\|^2 + \|x^{i+1} - \widehat{x}\|^2 \leq 2C_x \|x^{i+1} - \widehat{x}\| + \|x^i - \widehat{x}\|^2.$$

Rearranging the terms and using  $\|x^{i+1} - \widehat{x}\| \leq \|x^{i+1} - x^i\| + \|x^i - \widehat{x}\|$  yields

$$(\|x^{i+1} - x^i\| - C_x)^2 + \|x^{i+1} - \widehat{x}\|^2 \leq (\|x^i - \widehat{x}\| + C_x)^2,$$

which further leads to

$$(3.23) \quad \|x^{i+1} - \widehat{x}\| \leq \|x^i - \widehat{x}\| + C_x.$$

To estimate the dual variable, we multiply (PP) by  $Z_{i+1}^*(u^{i+1} - \widehat{u})$  with  $\Phi_i = 0, \Psi_{i+1} = I$ , yielding

$$0 \in \langle \partial F^*(y^{i+1}) - K(\bar{x}^{i+1}), y^{i+1} - \widehat{y} \rangle_{\Sigma_{i+1}} + \langle y^{i+1} - y^i, y^{i+1} - \widehat{y} \rangle.$$

Using  $0 \in \partial F^*(\widehat{y}) - K(\widehat{x})$  and following the steps leading to (3.23), we deduce

$$(3.24) \quad \|y^{i+1} - \widehat{y}\| \leq \|y^i - \widehat{y}\| + C_y$$

with  $C_y := \|K(\widehat{x}) - K(\bar{x}^{i+1})\|_{\Sigma_{i+1}^2}$ .

We now proceed to deriving bounds on  $C_x$  and  $C_y$  with the goal of bounding (3.23) and (3.24) from above. Using Assumption 3.1 and arguing as in (3.2), we estimate

$$(3.25) \quad C_x \leq \|T_i\|(\|\nabla K(x^i)\| \|y^i - \widehat{y}\| + L\|P_{\text{NL}}\widehat{y}\| \|x^i - \widehat{x}\|) =: R_x,$$

and, if  $\bar{x}^{i+1} \in \mathcal{X}_K$ ,

$$(3.26) \quad C_y \leq \|\Sigma_{i+1}\|(L\|\bar{x}^{i+1} - \widehat{x}\|/2 + \|\nabla K(\widehat{x})\|)\|\bar{x}^{i+1} - \widehat{x}\| =: R_y.$$

We thus need to verify first that  $\bar{x}^{i+1} \in \mathcal{X}_K$ . By definition,

$$\begin{aligned} \|\bar{x}^{i+1} - \widehat{x}\|^2 &= \|x^{i+1} - \widehat{x} + \omega_i(x^{i+1} - x^i)\|^2 \\ &= \|x^{i+1} - \widehat{x}\|^2 + \omega_i^2 \|x^{i+1} - x^i\|^2 + 2\omega_i \langle x^{i+1} - \widehat{x}, x^{i+1} - x^i \rangle \\ &= (1 + \omega_i) \|x^{i+1} - \widehat{x}\|^2 + \omega_i(1 + \omega_i) \|x^{i+1} - x^i\|^2 - \omega_i \|x^i - \widehat{x}\|^2 \\ &\leq (1 + \omega_i)(\|x^{i+1} - \widehat{x}\|^2 + \|x^{i+1} - x^i\|^2) - \omega_i \|x^i - \widehat{x}\|^2. \end{aligned}$$

Now, the bound (3.20) on  $T_i$  together with the definition of  $R_x$  implies that  $C_x \leq R_x \leq \delta_{x,i}/2$ . Applying (3.22) and (3.23), we obtain from this that

$$\|\bar{x}^{i+1} - \widehat{x}\|^2 \leq 4C_x \|x^{i+1} - \widehat{x}\| + \|x^i - \widehat{x}\|^2 \leq 4C_x(r_{x,i} + C_x) + r_{x,i}^2 \leq (r_{x,i} + \delta_{x,i})^2.$$

In addition, (3.23) shows that  $\|x^{i+1} - \widehat{x}\| \leq r_{x,i} + \delta_{x,i}$  as well. Similarly, the bound (3.20) on  $\Sigma_{i+1}$  implies that  $C_y \leq R_y \leq \delta_y$ , and hence (3.24) and (3.26) lead to  $\|y^{i+1} - \widehat{y}\| \leq r_y + \delta_y$ , completing the proof.  $\square$

Note that only the radii  $r_{x,i}, r_{x,i} + \delta_{x,i}$  of the primal neighborhoods depend on the iteration, and we will later seek to control these based on actual convergence estimates.



**Remark 3.9.** Suppose that  $\mathcal{X}_K = \mathcal{X}_G = X$ . Then we can take  $\delta_{x,i}$  arbitrarily large in order to satisfy the bound on  $T_i$  while still satisfying  $\mathbb{B}(\widehat{x}, r_{x,i} + \delta_{x,i}) \subset \mathcal{X}_K$ . On the other hand, the bound on  $\Sigma_{i+1}$  will go to zero as  $\delta_{x,i} \rightarrow \infty$ , which seems at first very limiting. However, we observe from the proof that this bound is actually not required to satisfy  $x^{i+1}, \bar{x}^{i+1} \in \mathbb{B}(\widehat{x}, r_{x,i} + \delta_{x,i})$ . Furthermore, if  $\text{dom } F^*$  is bounded and we take  $r_y$  large enough that  $\text{dom } F^* \subseteq \mathbb{B}(\widehat{y}, r_y)$ , the property  $y^{i+1} \in \mathbb{B}(\widehat{y}, r_y + \delta_y)$  is automatically satisfied by the iteration and does not require dual step length bounds. Hence if  $\mathcal{X}_K = \mathcal{X}_G = X$  and  $\text{dom } F^*$  is bounded, we can expect global convergence. We will return to this topic in [Remark 4.5](#).

## 4 REFINEMENT TO SCALAR STEP LENGTHS

To derive convergence rates, we now simplify [Theorem 3.7](#) to scalar step lengths. Specifically, we assume for some scalars  $\gamma_G, \gamma_L, \gamma_{NL}, \tau_i, \phi_i, \sigma_i, \psi_i \geq 0$ , and  $\lambda \in \mathbb{R}$  the structure

$$(4.1) \quad \begin{cases} T_i = \tau_i I, & \Phi_i = \phi_i I, & \Gamma_G = \gamma_G I, \\ \Sigma_i = \sigma_i I, & \Psi_i = \psi_i I, & \Gamma_{F^*} = \gamma_L P_L + \gamma_{NL} P_{NL}, \quad \text{and} \quad \Lambda = \lambda I. \end{cases}$$

Consequently, the preconditioning, step length, and testing operators simplify to

$$\begin{aligned} M_{i+1} &:= \begin{pmatrix} I & -\tau_i [\nabla K(x^i)]^* \\ -\omega_i \sigma_{i+1} \nabla K(x^i) & I \end{pmatrix}, \\ W_{i+1} &:= \begin{pmatrix} \tau_i I & 0 \\ 0 & \sigma_{i+1} I \end{pmatrix} \quad \text{and} \quad Z_{i+1} := \begin{pmatrix} \phi_i I & 0 \\ 0 & \psi_{i+1} I \end{pmatrix}. \end{aligned}$$

This reduces (PP) to [Algorithm 1.1](#), which for convex, proper, lower semicontinuous  $G$  and  $F^*$  is always solvable for the iterates  $\{u^i := (x^i, y^i)\}_{i \in \mathbb{N}}$ . Before proceeding to the main results, we state next all our assumptions in scalar form. We then derive in [Section 4.2](#) our main convergence rates results for [Algorithm 1.1](#) under specific update rules depending on monotonicity properties of  $G$  and  $F^*$ . The final [Section 4.3](#) is devoted to giving sufficient conditions for the scalar version of the assumptions of [Lemma 3.8](#) to hold.

### 4.1 GENERAL DERIVATIONS AND ASSUMPTIONS

Under the setup [\(4.1\)](#), the update rules [\(3.9\)](#) and the conditions [\(3.19\)](#) simplify to

$$\begin{aligned} (4.2a) \quad \omega_i &= \eta_i / \eta_{i+1}, & \eta_i &= \psi_i \sigma_i = \phi_i \tau_i, \\ (4.2b) \quad \phi_{i+1} &= \phi_i (1 + 2\tau_i \widetilde{\gamma}_G), & \psi_{i+2} &= \psi_{i+1} (1 + 2\sigma_{i+1} \widetilde{\gamma}_{F^*}), \\ (4.2c) \quad \phi_i &\geq \eta_i \delta^{-1} (\lambda + (\omega_i + 2)L\rho_y), & \psi_{i+1} &\geq \frac{\eta_i^2 \phi_i^{-1}}{1 - \kappa} \|\nabla K(x^i)\|^2, \\ (4.2d) \quad \gamma_L &\geq \widetilde{\gamma}_{F^*}, & \gamma_{NL} &\geq \widetilde{\gamma}_{F^*} + (p - 1)\zeta, \end{aligned}$$

for some  $\eta_i > 0$ ,  $p \in [1, 2]$ ,  $\zeta > 0$ ,  $0 \leq \delta \leq \kappa < 1$ , and  $\widetilde{\gamma}_{F^*}$ , for which we will from now on further assume  $\widetilde{\gamma}_{F^*} \geq 0$ . To formulate the scalar version of [Assumption 3.3](#), we also introduce the corresponding factor  $\widetilde{\gamma}_G \geq 0$ ; see [Assumption 4.1 \(iv\)](#) below.

Let us comment on these relations in turn. The conditions (4.2a) and (4.2b) limit the rate of growth of the testing parameters – and thus the convergence rate – and set basic coupling conditions for the step length parameters. They are virtually unchanged from the standard case of linear  $K$ ; see [25, Example 3.2].

The conditions in (4.2c) are essentially step length bounds: Substituting  $\eta_i = \phi_i \tau_i$  and  $\eta_i^2 = \phi_i \tau_i \psi_i \sigma_i$  in (4.2c), we obtain

$$(4.3) \quad \tau_i \leq \frac{\delta}{\lambda + (\omega_i + 2)L\rho_y}, \quad \text{and} \quad \sigma_i \tau_i \leq \frac{1 - \kappa}{R_K^2},$$

where  $R_K = \sup_{\mathcal{X}} \|\nabla K(x)\|$ . In the latter bound, we also used  $\psi_{i+1} \geq \psi_i$ , which follows from (4.2b) and  $\tilde{\gamma}_{F^*} \geq 0$ . We point out that this condition is simply a variant for nonlinear  $K$  of the standard initialization condition  $\tau \sigma \|K\|^2 < 1$  for the PDHGM for linear  $K$ . From (4.3), we see that we need to initialize and keep the dual iterates at a known finite distance  $\rho_y$  from  $\hat{y}$ . It also individually bounds the primal step length based on further properties of the specific saddle-point problem. The nonconvexity enters here via the factor  $\lambda$  from the three-point condition on  $K$  (Assumption 3.3).

The final condition (4.2d) bounds the acceleration parameters  $\tilde{\gamma}_{F^*}$  based on the actually available strong monotonicity minus any penalties we get from nonlinearity of  $K$  if Assumption 3.3 is satisfied with  $p \in (1, 2]$ . However,  $K$  may contribute to strong monotonicity of the primal problem at  $\hat{x}$ , so it can in some specific problems be possible to choose  $\tilde{\gamma}_G > \gamma_G$ .

Before we will further refine these bounds in the following sections, we collect the scalar refinements of all the structural assumptions of Section 3.

**Assumption 4.1.** Suppose  $G : X \rightarrow \overline{\mathbb{R}}$  and  $F^* \rightarrow \overline{\mathbb{R}}$  are convex, proper, and lower semicontinuous, and  $K \in C^1(X; Y)$ . Furthermore:

- (i) (*locally Lipschitz  $\nabla K$* ) There exists  $L \geq 0$  with  $\|\nabla K(x) - \nabla K(x')\| \leq L\|x - x'\|$  for any  $x, x' \in \mathcal{X}_K$ .
- (ii) (*locally bounded  $\nabla K$* ) There exists  $R_K > 0$  with  $\sup_{x \in \mathcal{X}_K} \|\nabla K(x)\| \leq R_K$ .
- (iii) (*monotone  $\partial G$  and  $\partial F^*$* ) The mapping  $\partial G$  is ( $\gamma_G I$ -strongly) monotone at  $\hat{x}$  for  $-\nabla K(\hat{x})^* \hat{y}$  in  $\mathcal{X}_G$  with  $\gamma_G \geq 0$ ; and the mapping  $\partial F^*$  is ( $\gamma_L P_L + \gamma_{NL} P_{NL}$ -strongly) monotone at  $\hat{y}$  for  $\hat{\xi} + K(\hat{x})$  in  $\mathcal{Y}_{F^*}$  with  $\gamma_L, \gamma_{NL} \geq 0$ .
- (iv) (*three-point condition on  $K$* ) For some  $p \in [1, 2]$ , some  $\lambda, \theta \geq 0$ , and any  $x, x' \in \mathcal{X}_K$ ,

$$\begin{aligned} & \langle [\nabla K(x') - \nabla K(\hat{x})]^* \hat{y}, x - \hat{x} \rangle + (\gamma_G - \tilde{\gamma}_G) \|x - \hat{x}\|^2 \\ & \geq \theta \|K(\hat{x}) - K(x) - \nabla K(x)(\hat{x} - x)\|^p - \frac{\lambda}{2} \|x - x'\|^2. \end{aligned}$$

- (v) (*neighborhood-compatible iterations*) The iterates of Algorithm 1.1 satisfy  $\{u^i\}_{i \in \mathbb{N}} \in \mathcal{U}(\rho_y)$  and  $\{\bar{x}^{i+1}\}_{i \in \mathbb{N}} \in \mathcal{X}_K$  for some  $\rho_y \geq 0$ , where  $\mathcal{U}(\rho_y)$  is given by (3.8).

We again close with remarks on the assumptions.

- Assumptions (i) and (ii) are standard assumptions in nonlinear optimization of smooth functions.
- Assumption (iii) is always satisfied due to the assumed convexity of  $G$  and  $F^*$ ; it only becomes restrictive under the additional requirement that  $\gamma_G$  or  $\gamma_L, \gamma_{NL}$  are positive, which will be needed to derive convergence rates in the next section. However, we stress that we never require the functions to be strongly monotone globally; rather we only require the strong monotonicity *at*  $\hat{x}$  or *at*  $\hat{y}$ . Furthermore – and this is a crucial feature of our operator-based analysis – we can split the strong monotonicity of  $\partial F^*$  on the subspaces  $Y_{NL}$  and  $Y_L$ . We will more depend on the former, which is often automatic as in the example given in the introduction.
- We have already elaborated on (iv) in [Example 3.1](#) and [Proposition 3.2](#). In particular, this condition holds if the saddle-point problem satisfies a standard second-order growth condition at  $\hat{u}$  with respect to the primal variable.
- Finally, only assumption (v) is specific to the actual algorithm (i.e., the choice of step sizes); it requires that the iterates of [Algorithm 1.1](#) remain in the neighborhood in which the first four assumptions are valid. We will prove in [Section 4.3](#) that this can be guaranteed under additional bounds on the step lengths. Moreover, we will demonstrate in [Remark 4.5](#) that [Assumption 4.1 \(v\)](#) is always satisfied for a specific class of problems, for which we therefore obtain global convergence.

## 4.2 CONVERGENCE RESULTS

We now come to the core of our work, where we apply the analysis of the preceding sections to derive convergence results under explicit step lengths rules. We start with a weak convergence result that requires no (partial) strong monotonicity, which however needs to be replaced with additional assumptions on  $K$ .

**Theorem 4.1 (weak convergence).** *Suppose [Assumption 4.1](#) holds for some  $L \geq 0, R_K > 0, \lambda \geq 0$ , and  $\rho_y \geq 0$  and choose the step lengths as  $\tau_i \equiv \tau, \sigma_i \equiv \sigma$ , and  $\omega_i \equiv 1$ . Assume that for some  $\zeta > 0$  and  $p \in [1, 2]$  that*

$$(4.4) \quad \gamma_{NL} \geq (p-1)\zeta \quad \text{and} \quad \theta \geq \rho_y^{2-p} p^{-p} \zeta^{1-p},$$

and for some  $0 < \delta < \kappa < 1$  that

$$(4.5) \quad 0 < \tau < \frac{\delta}{3L\rho_y + \lambda} \quad \text{and} \quad 0 < \sigma\tau < \frac{1-\kappa}{R_K^2}.$$

Furthermore, suppose that

$$(i) \quad x^i \rightarrow \bar{x} \text{ implies that } \nabla K(x^i)x \rightarrow \nabla K(\bar{x})x \text{ for all } x \in X,$$

and either

$$(ia) \quad H(u) \text{ is maximally monotone in } \mathcal{U}(\rho_y);$$

- (iib) the mapping  $(x, y) \mapsto ([\nabla K(x)]^* y, K(x))$  is weak-to-strong continuous in  $\mathcal{U}(\rho_y)$ ; or
- (iic) the mapping  $(x, y) \mapsto ([\nabla K(x)]^* y, K(x))$  is weak-to-weak continuous, and in addition Assumption 4.1 (iii) (monotone  $\partial G$  and  $\partial F^*$ ) and Assumption 4.1 (iv) (three-point condition on  $K$ ) hold at any weak limit  $\bar{u} = (\bar{x}, \bar{y})$  of  $\{u^i\}$  in addition to  $\hat{u}$  with  $\theta \geq (2\rho_y)^{2-p} p^{-p} \zeta^{1-p}$ .

Then the sequence  $\{u^i\}$  generated by (PP) converges weakly to some  $\bar{u} \in H^{-1}(0)$  (possibly different from  $\hat{u}$ ).

*Proof.* We wish to apply Proposition 2.2 and therefore need to verify its assumptions. For the basic assumptions of (CI) and the self-adjointness of  $Z_{i+1}M_{i+1}$ , we will use Theorem 3.7 together with Lemma 3.4. Most of their assumptions are directly verified by Assumption 4.1; it only remains to verify (3.9) and (3.19), which reduce to (4.2) in the scalar case. By taking  $\tilde{\gamma}_G = \tilde{\gamma}_{F^*} = 0$  and any positive constants  $\psi$  and  $\phi$  such that  $\psi\sigma = \phi\tau$ , the relations (4.2a), (4.2b), and (4.2d) hold. Furthermore, since  $\omega_i \equiv 1$ , (4.5) is equivalent to (4.3). This yields (4.2c), completing the verification of (4.2) and thus the conditions of Theorem 2.1 and Proposition 2.2.

Since the inequalities in (4.5) are strict, we can deduce from Theorem 3.7 that (3.18) even holds for  $\Delta_{i+1} \leq -\hat{\delta}\|u^{i+1} - u^i\|^2$  for some  $\hat{\delta} > 0$ . Combining (3.11) and (4.5), we thus obtain that condition (i) of Proposition 2.2 holds. Furthermore, the condition (iii) follows from the assumed constant step lengths and the assumption (i).

It remains to show the condition (ii) of Proposition 2.2. First, if the assumption (iia) holds, the inclusion in condition (ii) follows directly from the fact that maximally monotone operators have sequentially weakly-strongly closed graphs [2, Proposition 20.38].

The two other cases are more difficult to verify. First, we note that for any  $x^{i+1} \rightharpoonup \bar{x}$  and  $y^{i+1} \rightharpoonup \bar{y}$  we have  $W_{i+1} \equiv W$ , and (PP) implies that  $v_{i+1} \in WA(u^{i+1})$  for

$$A(u^{i+1}) := \begin{pmatrix} \partial G(x^{i+1}) - \gamma_G(x^{i+1} - \bar{x}) \\ \partial F^*(y^{i+1}) - \gamma_{\text{NL}} P_{\text{NL}}(y^{i+1} - \bar{y}) \end{pmatrix}$$

and

$$(4.6) \quad v_{i+1} := W \begin{pmatrix} -[\nabla K(x^{i+1})]^* y^{i+1} - \gamma_G(x^{i+1} - \bar{x}) \\ K(x^{i+1}) - \gamma_{\text{NL}} P_{\text{NL}}(y^{i+1} - \bar{y}) \end{pmatrix} - W \begin{pmatrix} [\nabla K(x^i) - \nabla K(x^{i+1})]^* y^{i+1} \\ K(x^{i+1}) - K(\bar{x}^{i+1}) - \nabla K(x^i)(x^{i+1} - \bar{x}^{i+1}) \end{pmatrix} - M_{i+1}(u^{i+1} - u^i).$$

Therefore, we need to show that

$$v_{i+1} \rightharpoonup \bar{v} := W \begin{pmatrix} -[\nabla K(\bar{x})]^* \bar{x} \\ K(\bar{x}) \end{pmatrix} \quad \text{and} \quad \bar{v} \in A(\bar{u}),$$

which by construction is equivalent to the inclusion  $\bar{u} \in H^{-1}(0)$ . Note that due to Assumption 4.1 (iii),  $A$  is maximally monotone since it only involves subgradient mappings of proper, convex, and lower semicontinuous functions. From Theorem 2.1, we obtain (DI), which in turn implies that  $Z_{i+1}M_{i+1}(u^{i+1} - u^i) \rightarrow 0$ . The scalar case of Lemma 3.4 together with the condition  $0 < \delta < \kappa < 1$  and positive  $\psi$  and  $\phi$  then results in  $\|u^{i+1} - u^i\| \rightarrow 0$ , so the last two terms in

(4.6) go to zero. We therefore only have to consider the first term, for which we make a case distinction:

If assumption (iib) holds, we obtain that  $v_{i+1} \rightarrow \bar{v}$ , and the required inclusion  $\bar{v} \in A(\bar{u})$  follows from the fact that the graph of the maximally monotone operator  $A$  is sequentially weakly–strongly closed; see [2, Proposition 16.36].

If assumption (iic) holds, then only  $v_{i+1} \rightarrow \bar{v}$ . In this case, we can apply the Brezis–Crandall–Pazy Lemma [2, Corollary 20.59 (iii)] to obtain the required inclusion under the additional condition that  $\limsup_{i \rightarrow \infty} \langle u_i - \bar{u}, v_i - \bar{v} \rangle \leq 0$ . In our case,  $\limsup_{i \rightarrow \infty} \langle u_i - \bar{u}, v_i - \bar{v} \rangle = \limsup_{i \rightarrow \infty} q_i$  for

$$q_i := \langle [\nabla K(\bar{x})]^* \bar{y} - [\nabla K(x^{i+1})]^* y^{i+1}, x^{i+1} - \bar{x} \rangle + \langle K(x^{i+1}) - K(\bar{x}), y^{i+1} - \bar{y} \rangle - \gamma_{\text{NL}} \|y^{i+1} - \bar{y}\|_{P_{\text{NL}}}^2 - \gamma_G \|x^{i+1} - \bar{x}\|^2.$$

Note that  $\|y^{i+1} - \bar{y}\|_{P_{\text{NL}}} \leq 2\rho_y$  because  $\|y^{i+1} - \widehat{y}\|_{P_{\text{NL}}}, \|\widehat{y} - \bar{y}\|_{P_{\text{NL}}} \leq \rho_y$ . With this, (3.2), and both Assumption 4.1 (iii) and (iv) at  $\bar{u}$ , we similarly to (3.17) estimate

$$\begin{aligned} (4.7) \quad q_i &= \langle K(x^{i+1}) - K(\bar{x}) + \nabla K(x^{i+1})(\bar{x} - x^{i+1}), y^{i+1} - \bar{y} \rangle \\ &\quad - \langle (\nabla K(x^i) - \nabla K(\bar{x}))(\bar{x} - x^i), \bar{y} \rangle + \gamma_G \|x^{i+1} - \bar{x}\|^2 \\ &\quad - \gamma_{\text{NL}} \|y^{i+1} - \bar{y}\|_{P_{\text{NL}}}^2 + \langle (\nabla K(x^i) - \nabla K(x^{i+1}))(\bar{x} - x^i), \bar{y} \rangle \\ &\leq (\|y^{i+1} - \bar{y}\|_{P_{\text{NL}}}^{2-p} p^{-p} \zeta^{1-p} - \theta) \|K(\bar{x}) - K(x^{i+1}) - \nabla K(x^{i+1})(\bar{x} - x^{i+1})\|^p \\ &\quad + ((p-1)\zeta - \gamma_{\text{NL}}) \|y^{i+1} - \bar{y}\|_{P_{\text{NL}}}^2 + O(\|x^{i+1} - \bar{x}\|). \end{aligned}$$

Since

$$\|y^{i+1} - \bar{y}\|_{P_{\text{NL}}}^{2-p} p^{-p} \zeta^{1-p} - \theta \leq (2\rho_y)^{2-p} p^{-p} \zeta^{1-p} - \theta \leq 0,$$

$(p-1)\zeta - \gamma_{\text{NL}} \leq 0$ , and  $Z_{i+1} M_{i+1}(u^{i+1} - u^i) \rightarrow 0$ , we obtain that  $\limsup_{i \rightarrow \infty} q_i \leq 0$ . The Brezis–Crandall–Pazy Lemma thus yields the desired inclusion  $\bar{v} \in A(\bar{u})$ .

Hence in all three cases, the condition (ii) of Proposition 2.2 holds with  $u^i \rightarrow \bar{u} \in H^{-1}(0)$ , which completes the proof.  $\square$

**Remark 4.2.** *It is instructive to consider the two limiting cases  $p = 1$  and  $p = 2$  in the conditions (4.4) of Theorem 4.1:*

- (i) *If  $p = 1$ , the condition on  $\gamma_{\text{NL}}$  is trivially satisfied, while the one on  $\theta$  reduces to  $\rho_y \leq \theta$ . Therefore, we only require the dual variable to be initialized close to  $\widehat{y}$  (and only when projected into the subspace  $Y_{\text{NL}}$ ).*
- (ii) *In contrast, if  $p = 2$ , the condition on  $\theta$  does not involve  $\rho_y$  and hence there will be no dual initialization bound; on the other hand,  $\zeta \geq (4\theta)^{-1}$  will be required. Therefore, we need  $F^*$  to be strongly convex with the factor  $\gamma_{\text{NL}} \geq (4\theta)^{-1}$ , but only at  $\widehat{y}$  within the subspace  $Y_{\text{NL}}$ .*

*The remaining cases  $p \in (1, 2)$  can be seen as an interpolation between these conditions. The same observations hold for the other results of this section.*

We now turn to convergence rates under strong monotonicity assumptions, starting with the case that merely  $G$  is strongly convex. Since we obtain *a fortiori* strong convergence from the rates, we do not require the additional assumptions on  $K$  needed to apply [Proposition 2.2](#); on the other hand, we only obtain convergence of the primal iterates. In the proof of the following result, observe how the step length choice follows directly from having to satisfy [\(4.2b\)](#) and the desire to keep  $\sigma_i \tau_i$  constant to satisfy [\(4.2c\)](#) via the bound [\(4.3\)](#) on the initial choice.

**Theorem 4.3 (convergence rates under acceleration).** *Suppose [Assumption 4.1](#) holds for some  $L \geq 0, R_K > 0, \lambda \geq 0$ , and  $\rho_y \geq 0$  with  $\tilde{\gamma}_G > 0, \gamma_{\text{NL}} \geq (p-1)\zeta$ , and  $\theta \geq \rho_y^{2-p} p^{-p} \sqrt{1+2\tau_0 \tilde{\gamma}_G} \zeta^{1-p}$  for some  $\zeta > 0$  and  $p \in [1, 2]$ . Choose*

$$(4.8) \quad \tau_{i+1} = \tau_i \omega_i, \quad \sigma_{i+1} = \sigma_i / \omega_i, \quad \omega_i = 1 / \sqrt{1 + 2\tau_i \tilde{\gamma}_G}$$

with

$$(4.9) \quad 0 < \tau_0 \leq \frac{\delta}{3L\rho_y + \lambda}, \quad \text{and} \quad 0 < \tau_0 \sigma_0 \leq \frac{1 - \kappa}{R_K^2}.$$

for some  $0 < \delta \leq \kappa < 1$ . Then  $\|x^i - \hat{x}\|^2$  converges to zero at the rate  $O(1/N^2)$ .

*Proof.* The first stage of the proof is similar to [Theorem 4.1](#), where we verify [\(4.2\)](#) to use [Theorem 3.7](#) (but need not apply [Proposition 2.2](#)). Since we do not assume any (partial) strong convexity of  $F^*$ , we have to take  $\tilde{\gamma}_{F^*} = 0$ , and thus [\(4.2d\)](#) is satisfied by assumption. Note that by [\(4.8\)](#), we have  $\sigma_i \tau_i = \sigma_0 \tau_0$  for all  $i \in \mathbb{N}$ ,  $\omega_i < 1$ , and  $\tau_{i+1} < \tau_0$ . Then [\(4.9\)](#) leads to [\(4.3\)](#), which is equivalent to [\(4.2c\)](#). Taking now  $\eta_i := \sigma_i > 0$ ,  $\psi_i \equiv 1$ , and  $\phi_i := \sigma_0 \tau_0 \tau_i^{-2} > 0$ , [\(4.2a\)](#) and [\(4.2b\)](#) follow from [\(4.8\)](#) since  $\eta_i := \sigma_i = \sigma_{i+1} \omega_i = \eta_{i+1} \omega_i$  and

$$\phi_i \tau_i = \frac{\sigma_0 \tau_0}{\tau_i} = \frac{\sigma_i \tau_i}{\tau_i} = \psi_i \sigma_i \quad \text{and} \quad \phi_{i+1} := \frac{\psi \sigma_0 \tau_0}{\tau_{i+1}^2} = \frac{\psi \sigma_0 \tau_0}{\tau_i^2 \omega_i^2} = \phi_i (1 + 2\tau_i \tilde{\gamma}_G).$$

Furthermore, [\(4.8\)](#) also implies that

$$1/\omega_i \leq 1/\omega_0 = \sqrt{1 + 2\tau_0 \tilde{\gamma}_G},$$

and together with our assumption on  $\theta$  we obtain that  $\theta \geq \rho_y^{2-p} p^{-p} \omega_i^{-1} \zeta^{1-p}$ . We can thus apply [Theorems 2.1](#) and [3.7](#) to arrive at [\(DI\)](#) with each  $\Delta_{i+1} \leq 0$ .

We now estimate the convergence rate from [\(DI\)](#) by bounding  $Z_{N+1} M_{N+1}$  from below. Using [Lemma 3.4](#), we obtain that

$$(4.10) \quad \delta \phi_N \|x^N - \hat{x}\|^2 \leq \|u^0 - \hat{u}\|_{Z_1 M_1}^2.$$

But from [\[7, Corollary 1\]](#), we know that  $\tau_N = O(N^{-1})$  as  $N \rightarrow \infty$  and hence  $\phi_N \sim \tau_N^{-2} = O(N^2)$  by our choice of  $\phi_N$ , which yields the desired convergence rate.  $\square$

If both  $\partial G$  and  $\partial F^*$  are strongly monotone, [Algorithm 1.1](#) with constant step lengths leads to convergence of both primal and dual iterates at a linear rate.

**Theorem 4.4 (linear convergence).** *Suppose Assumption 4.1 holds for some  $L \geq 0, R_K > 0, \lambda \geq 0$ , and  $\rho_y \geq 0, \tilde{\gamma}_G > 0, \tilde{\gamma}_{F^*} := \min\{\gamma_L, \gamma_{NL} - (p-1)\zeta\} > 0$ , and  $\theta \geq \rho_y^{2-p} p^{-p} \omega^{-1} \zeta^{1-p}$  for some  $\zeta > 0$  and  $p \in [1, 2]$ . Choose*

$$(4.11) \quad 0 < \tau_i \equiv \tau \leq \min\left\{\frac{\delta}{3L\rho_y + \lambda}, \frac{\sqrt{(1-\kappa)\tilde{\gamma}_{F^*}/\tilde{\gamma}_G}}{R_K}\right\}, \quad \sigma_i \equiv \sigma := \frac{\tilde{\gamma}_G}{\tilde{\gamma}_{F^*}}\tau, \quad \omega_i \equiv \omega := \frac{1}{1+2\tilde{\gamma}_G\tau}$$

for some  $0 \leq \delta \leq \kappa < 1$ . Then  $\|u^i - \hat{u}\|^2$  converges to zero at the rate  $O(\omega^N)$ .

*Proof.* The proof follows that of Theorem 4.3. To verify (4.2), we take  $\psi_0 := 1/\sigma$ , and  $\phi_0 := 1/\tau$ , for which (4.2b) is satisfied due to the second relation of (4.11). By induction, we further obtain from this

$$(4.12) \quad \phi_i \tau = \psi_i \sigma = (1 + 2\tilde{\gamma}_G \tau)^i$$

for all  $i \in \mathbb{N}$ , verifying (4.2a). Inequality (4.2d) holds due to  $\tilde{\gamma}_{F^*} = \min\{\gamma_L, \gamma_{NL} - (p-1)\zeta\} > 0$ . It remains to prove (4.2c), which follow via (4.3) from the bound on  $\tau$  in (4.11). Finally, we apply Lemma 3.4, (4.12) for  $i = N$ , and Theorem 2.1 to conclude that

$$(1 + 2\tilde{\gamma}_G \tau)^N \left( \frac{\delta}{2\tau} \|x^N - \hat{x}\|^2 + \frac{\kappa - \delta}{2\sigma(1 - \delta)} \|y^N - \hat{y}\|^2 \right) \leq \frac{1}{2} \|u^0 - \hat{u}\|_{Z_1 M_1}^2,$$

which yields the desired convergence rate.  $\square$

**Remark 4.5 (global convergence).** *Following Remark 3.9, suppose that Assumption 4.1 (i)–(iv) hold for  $\mathcal{X}_K = \mathcal{X}_G = X$  and that  $\text{dom } F^*$  is bounded. If we then take  $\rho_y$  large enough that  $\text{dom } F^* \subseteq \mathbb{B}_{NL}(\hat{y}, \rho_y)$ , Assumption 4.1 (v) will be satisfied for any choice of starting point  $u^0 = (x^0, y^0) \in X \times Y$ , i.e., we have global convergence. Note, however, that in this case we need  $\nabla K$  to be bounded on the whole space, i.e.,  $R_K = \sup_{x \in X} \|\nabla K(x)\| < \infty$  has to hold.*

### 4.3 NEIGHBORHOOD-COMPATIBLE ITERATIONS

To conclude our analysis, we provide explicit conditions on the initialization of Algorithm 1.1 to ensure that Assumption 4.1 (v) (neighborhood-compatible iterations) holds in cases where the global convergence of Remark 4.5 cannot be guaranteed.

To begin, the following result shows that the rules (4.2) are consistent with the sequence  $\{u^i\}_{i \in \mathbb{N}}$  generated by (PP) remaining in  $\mathcal{U}(\rho_y)$ , provided that  $\tau_0$  is sufficiently small and that the starting point  $u^0 = (x^0, y^0)$  is sufficiently far inside the interior of  $\mathcal{U}(\rho_y)$ .

**Lemma 4.6.** *Let  $0 \leq \delta \leq \kappa < 1$  and  $\rho_y > 0$  be given, and assume (4.2) holds with*

$$(4.13) \quad 1/\sqrt{1 + 2\tau_i \tilde{\gamma}_G} \leq \omega_i \leq \omega_{i+1} \leq 1 \quad (i \in \mathbb{N}).$$

*Assume further that  $\sup_{x \in \mathcal{X}_K} \|\nabla K(x)\| \leq R_K$ . Define*

$$(4.14) \quad r_{\max} := \sqrt{2\delta^{-1}(\|x^0 - \hat{x}\|^2 + \mu^{-1}\|y^0 - \hat{y}\|^2)} \quad \text{with} \quad \mu := \sigma_1 \omega_0 / \tau_0.$$



Assume also that  $\mathbb{B}(\widehat{x}, r_{\max} + \delta_x) \times \mathbb{B}(\widehat{y}, r_y + \delta_y) \subseteq \mathcal{U}(\rho_y)$  for some  $\delta_x, \delta_y > 0$  as well as  $r_y \geq r_{\max} \sqrt{\mu(1-\delta)\delta/(\kappa-\delta)}$ . If the initial primal step length  $\tau_0$  satisfies

$$(4.15) \quad \tau_0 \leq \min \left\{ \frac{\delta_x}{2R_K r_y + 2L\|P_{\text{NL}}\widehat{y}\|r_{\max}}, \frac{2\delta_y \omega_0 (r_{\max} + \delta_x)^{-1}}{(L(r_{\max} + \delta_x) + 2R_K)\mu} \right\},$$

then *Assumption 4.1 (v)* (neighborhood-compatible iterations) holds.

*Proof.* We first set up some basic set inclusions. Without loss of generality, we can assume that  $\psi_1 = 1$ , as we can always rescale the testing variables  $\phi_i$  and  $\psi_i$  by the same constant without violating (4.2). We then define  $r_{x,i} := \|u^0 - \widehat{u}\|_{Z_1 M_1} / \sqrt{\delta} \phi_i$ ,  $\delta_{x,i} := \sqrt{\phi_0 / \phi_i} \delta_x$ , and

$$\mathcal{U}_i := \{(x, y) \in X \times Y \mid \|x - \widehat{x}\|^2 + \frac{\psi_{i+1}}{\phi_i} \frac{\kappa - \delta}{(1 - \delta)\delta} \|y - \widehat{y}\|^2 \leq r_{x,i}^2\}.$$

We then observe from *Lemma 3.4* that

$$(4.16) \quad \{u \in X \times Y \mid \|u - \widehat{u}\|_{Z_{i+1} M_{i+1}} \leq \|u^0 - \widehat{u}\|_{Z_1 M_1}\} \subset \mathcal{U}_i.$$

From (4.2), we also deduce that  $\phi_{i+1} \geq \phi_i$  and hence that  $r_{x,i+1} \leq r_{x,i}$  as well as  $\delta_{x,i} \leq \delta_x$ . Consequently, if  $r_{x,0} \leq r_{\max}$ , then

$$(4.17) \quad \mathbb{B}(\widehat{x}, r_{x,i} + \delta_{x,i}) \times \mathbb{B}(\widehat{y}, r_y + \delta_y) \subseteq \mathbb{B}(\widehat{x}, r_{\max} + \delta_x) \times \mathbb{B}(\widehat{y}, r_y + \delta_y) \subseteq \mathcal{U}(\rho_y),$$

so it will suffice to show that  $u^i \in \mathbb{B}(\widehat{x}, r_{x,i} + \delta_{x,i}) \times \mathbb{B}(\widehat{y}, r_y + \delta_y)$  for each  $i \in \mathbb{N}$  to prove the claim. We do this in two steps. The first step shows that  $r_{x,i} \leq r_{\max}$  and

$$(4.18) \quad \mathcal{U}_i \subseteq \mathbb{B}(\widehat{x}, r_{x,i}) \times \mathbb{B}(\widehat{y}, r_y) \quad (i \in \mathbb{N}).$$

In the second step, we then show that  $u^i \in \mathcal{U}_i$  as well as  $\bar{x}^{i+1} \in \mathcal{X}_K$  for  $i \in \mathbb{N}$ . The two inclusion (4.17) and (4.18) then imply that *Assumption 4.1 (v)* holds.

**Step 1** We first prove (4.18). Since  $\mathcal{U}_i \subseteq \mathbb{B}(\widehat{x}, r_{x,i}) \times Y$ , we only have to show that  $\mathcal{U}_i \subseteq X \times \mathbb{B}(\widehat{y}, r_y)$ . First, note that (4.2) implies that  $\bar{y}_{F^*} \geq 0$  and therefore  $\psi_{i+1} \geq \psi_i \geq \psi_1 = 1$  as well as  $\phi_{i+1} \geq \phi_i \geq \phi_0 = \eta_1 \omega_0 / \tau_0 = \mu$ . We then obtain from the definition of  $r_{x,i}$  that

$$r_{x,i}^2 \delta \phi_i = \|u^0 - \widehat{u}\|_{Z_1 M_1}^2 = \mu \|x^0 - \widehat{x}\|^2 - 2\eta_0 \langle x^0 - \widehat{x}, [\nabla K(x^0)]^*(y^0 - \widehat{y}) \rangle + \|y^0 - \widehat{y}\|^2.$$

Using Cauchy's inequality, the fact that  $\phi_i \geq \mu$ , and the assumption  $\|\nabla K(x^0)\| \leq R_K$ , we arrive at

$$r_{x,i}^2 \leq (2\mu \|x^0 - \widehat{x}\|^2 + (1 + \eta_0^2 \phi_0^{-1} R_K^2) \|y^0 - \widehat{y}\|^2) / (\delta \mu).$$

We obtain from (4.2c) that  $\eta_0^2 \phi_0^{-1} R_K^2 \leq 1 - \kappa \leq 1$  and hence that  $r_{x,i}^2 \leq r_{x,0}^2$ . The assumption on  $r_y$  then yields that

$$(4.19) \quad r_y^2 \geq r_{\max}^2 \phi_0 \frac{(1-\delta)\delta}{\kappa-\delta} \geq \frac{r_{x,0}^2 \phi_0 (1-\delta)\delta}{\psi_{i+1} (\kappa-\delta)} = \frac{r_{x,i}^2 \phi_i (1-\delta)\delta}{\psi_{i+1} (\kappa-\delta)}$$

for all  $i \in \mathbb{N}$ , and (4.18) follows from the definition of  $\mathcal{U}_i$ .

**Step 2** We next show by induction that  $u^i \in \mathcal{U}_i$ ,  $\bar{x}^{i+1} \in \mathcal{X}_K$ , and

$$(4.20) \quad \tau_i \leq \frac{\delta_{x,i}}{2R_K r_y + 2L\|P_{\text{NL}}\widehat{y}\|_{r_{x,i}}}, \quad \text{and} \quad \sigma_{i+1}(r_{x,i} + \delta_{x,i}) \leq \frac{2\delta_y}{L(r_{x,i} + \delta_{x,i}) + 2R_K}$$

hold for all  $i \in \mathbb{N}$ .

Since (4.2) holds, we can apply Lemma 3.4 to  $\|u^0 - \widehat{u}\|_{Z_1 M_1}$  to verify that  $u^0 \in \mathcal{U}_0$ . Moreover, since  $\sigma_1 = \mu\tau_0/\omega_0$ , the bound (4.20) for  $i = 0$  follows from (4.15) and the bound  $r_{x,0} \leq r_{\max}$  from Step 1. This gives the induction basis.

Suppose now that  $u^N \in \mathcal{U}_N$  and that (4.20) holds for  $i = N$ . By (4.18), we have that  $u^N \in \mathbb{B}(\widehat{x}, r_{x,N}) \times \mathbb{B}(\widehat{y}, r_y)$ . Since (4.20) guarantees (3.20), we can apply Lemma 3.8 to obtain

$$u^{N+1} \in \mathbb{B}(\widehat{x}, r_{x,N} + \delta_{x,N}) \times \mathbb{B}(\widehat{y}, r_y + \delta_y) \quad \text{and} \quad \bar{x}^{N+1} \in \mathbb{B}(\widehat{x}, r_{x,N} + \delta_{x,N}).$$

Together with (4.17), we obtain that  $u^{N+1} \in \mathcal{U}(\rho_y)$  and  $\bar{x}^{N+1} \in \mathcal{X}_K$ . Theorems 2.1 and 3.7 now imply that (DI) is satisfied for  $i \leq N$  with  $\Delta_{N+1} \leq 0$ , which together with (DI) and (4.16) yields that  $u^{N+1} \in \mathcal{U}_{N+1}$ . This is the first part of the claim. To show (4.20), we deduce from (4.2) that

$$(4.21) \quad \tau_{N+1} = \frac{\tau_N \phi_N}{\phi_{N+1} \omega_N} = \frac{\tau_N}{\omega_N(1 + 2\tau_N \widetilde{Y}_G)}, \quad r_{x,N+1} = \frac{r_{x,N}}{\sqrt{1 + 2\tau_N \widetilde{Y}_G}},$$

$$(4.22) \quad \sigma_{N+2} = \frac{\sigma_{N+1} \psi_{N+1}}{\psi_{N+2} \omega_{N+1}} = \frac{\sigma_{N+1}}{\omega_{N+1}(1 + 2\sigma_{N+1} \widetilde{Y}_{F^*})}, \quad \delta_{x,N+1} = \frac{\delta_{x,N}}{\sqrt{1 + 2\tau_N \widetilde{Y}_G}}.$$

Hence, using  $\omega_{N+1} \geq \omega_N$ ,  $\omega_N \sqrt{1 + 2\tau_N \widetilde{Y}_G} \geq 1$ , and  $r_{x,N+1} \leq r_{x,N}$ , as well as the inductive assumption (4.20) shows that

$$\begin{aligned} \tau_{N+1} &= \frac{\delta_{x,N+1}}{\delta_{x,N}} \frac{\tau_N}{\omega_N \sqrt{1 + 2\tau_N \widetilde{Y}_G}} \leq \frac{\delta_{x,N+1}}{2R_K r_y + 2L\|P_{\text{NL}}\widehat{y}\|_{r_{x,N+1}}}, \quad \text{and} \\ \sigma_{N+2}(r_{x,N+1} + \delta_{x,N+1}) &\leq \frac{\sigma_{N+1}(r_{x,N} + \delta_{x,N})}{\omega_N \sqrt{1 + 2\tau_N \widetilde{Y}_G}} \leq \frac{2\delta_y}{L(r_{x,N+2} + \delta_{x,N+2}) + 2R_K}. \end{aligned}$$

This completes the induction step and hence the proof.  $\square$

If the step lengths and the over-relaxation parameter  $\omega_i$  are constant, we can remove the lower bound on  $\omega_i$  in Lemma 4.6.

**Lemma 4.7.** *The claims of Lemma 4.6 also hold for  $\tau_i \equiv \tau_0$ ,  $\sigma_i \equiv \sigma_1$ , and any choice of  $\omega_i \equiv \omega \leq 1$ . In particular,  $\omega$  can be chosen according to (4.11).*

*Proof.* The proof proceeds exactly as that of Lemma 4.6, replacing  $r_{x,i}$  by  $r_{x,0}$  and  $\delta_{x,i}$  by  $\delta_x$  everywhere. Since  $r_{x,i} \leq r_{x,0}$  and  $\delta_{x,i} \leq \delta_{x,0}$ , the bound (4.16) holds in this case as well, while (4.18) reduces to the case  $i = 0$  that was shown in Step 1. Observe also that the only place where we needed the lower bound on  $\omega_i$  was in the proof of (4.20) as part of the inductive step of Step 2. As the bound and the step lengths remain unchanged between iterations in this case, this is sufficient to conclude the proof, and the lower bound on  $\omega_i$  is thus not required.  $\square$

Finally, as long as we start close enough to a solution, no additional step length bounds are needed to guarantee Assumption 4.1(v).

**Proposition 4.8.** *Under the assumptions of Theorem 4.1, 4.3, or 4.4, suppose that  $\rho_y > 0$ . Then there exists an  $\varepsilon > 0$  such that for all  $u^0 = (x^0, y^0)$  satisfying*

$$(4.23) \quad \sqrt{2\delta^{-1}(\|x^0 - \widehat{x}\|^2 + \mu^{-1}\|y^0 - \widehat{y}\|^2)} \leq \varepsilon \quad \text{with} \quad \mu := \sigma_1\omega_0/\tau_0,$$

Assumption 4.1(v) (neighborhood-compatible iterations) holds.

*Proof.* For given  $\varepsilon > 0$ , set  $r_y = \varepsilon\sqrt{\mu(1-\delta)\delta}/(\kappa-\delta)$  as well as  $\delta_x = \sqrt{\varepsilon}$  and  $\delta_y = \rho_y - r_y$ . Then for  $\varepsilon > 0$  sufficiently small, both  $\delta_y > 0$  and  $\mathbb{B}(\widehat{x}, r_{\max} + \delta_x) \times \mathbb{B}(\widehat{y}, r_y + \delta_y) \subseteq \mathcal{U}(\rho_y)$  hold. Furthermore, (4.23) yields that  $r_{\max} \leq \varepsilon$  in Lemma 4.6. Since

$$\min \left\{ \frac{\varepsilon^{-1/2}}{2R_K\sqrt{\mu(1-\delta)\delta}/(\kappa-\delta) + 2L\|P_{\text{NL}}\widehat{y}\|}, \frac{2(r_y - \varepsilon)\omega_0\varepsilon^{-1/2}}{(L(\varepsilon + \sqrt{\varepsilon}) + 2R_K)(\sqrt{\varepsilon} + 1)\mu} \right\} \rightarrow \infty$$

for  $\varepsilon \rightarrow 0$ , we can guarantee that (4.15) holds for any given  $\tau_0 > 0$  by further reducing  $\varepsilon > 0$ . The claim then follows from Lemma 4.6.  $\square$

## 5 NUMERICAL EXAMPLES

We now illustrate the convergence and the effects of acceleration for a nontrivial example from PDE-constrained optimization. Following [11], we consider as the nonlinear operator the mapping from a potential coefficient in an elliptic equation to the corresponding solution, i.e., for a Lipschitz domain  $\Omega \subset \mathbb{R}^d$ ,  $d \leq 3$  and  $X = Y = L^2(\Omega)$ , we set  $S : x \mapsto z$  for  $z$  satisfying

$$(5.1) \quad \begin{cases} \Delta z + xz = f & \text{on } \Omega, \\ \partial_\nu z = 0 & \text{on } \partial\Omega. \end{cases}$$

Here  $f \in L^2(\Omega)$  is given; for our examples below we take  $f \equiv 1$ . The operator  $S$  is uniformly bounded for all  $x \geq \varepsilon > 0$  almost everywhere as well as completely continuous and twice Fréchet differentiable with uniformly bounded derivatives. Furthermore, for any  $h \in X$ , the application  $\nabla S(x)^*h$  of the adjoint Fréchet derivative can be computed by solving a similar elliptic equation; see [11, Section 3]. For our numerical examples, we take  $\Omega = (-1, 1)$  and approximate  $S$  by a standard finite element discretization on a uniform mesh with 1000 elements with piecewise constant  $x$  and piecewise linear  $z$ . We use the MATLAB codes accompanying [11] that can be downloaded from [10].

The first example is the  $L^1$  fitting problem

$$(5.2) \quad \min_{x \in L^2(\Omega)} \frac{1}{\alpha} \|S(x) - z^\delta\|_{L^1} + \frac{1}{2} \|x\|_{L^2}^2,$$

for some noisy data  $z^\delta \in L^2(\Omega)$  and a regularization parameter  $\alpha > 0$ ; see [11, Section 3.1] for details. For the purpose of this example, we take  $z^\delta$  as arising from random-valued impulsive

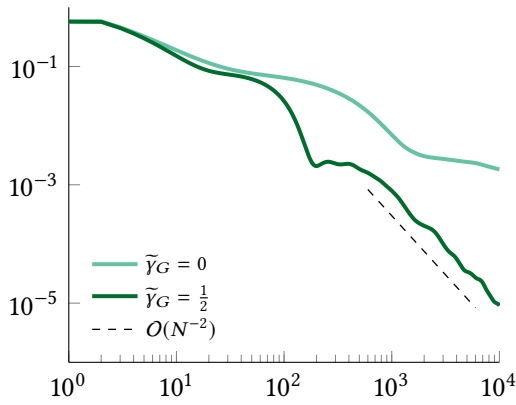


Figure 1:  $L^1$  fitting:  $\|x^N - \hat{x}\|_{L^2}^2$  for different values of  $\tilde{\gamma}_G$

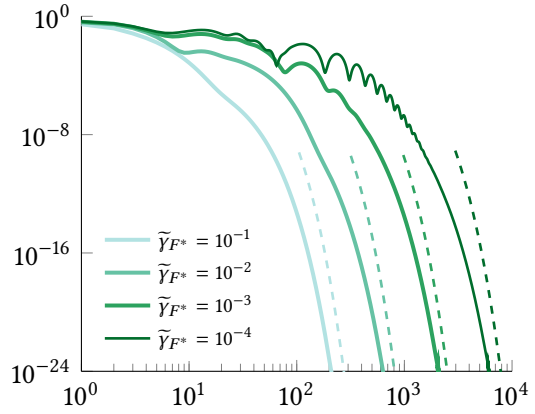


Figure 2:  $L^1$  fitting:  $\|u^N - \hat{u}\|_{L^2 \times L^2}^2$  (solid) and bounds  $(1 + 2\tilde{\gamma}_G\tau)^{-N}$  (dashed) for strongly convex  $F^*$  and different values of  $\tilde{\gamma}_{F^*}$

noise applied to  $z^\dagger = S(x^\dagger)$  for  $x^\dagger(t) = 2 - |t|$  and  $\alpha = 10^{-2}$ . This fits into the framework of problem (P) with  $F(y) = \frac{1}{\alpha}\|y\|_{L^1}$ ,  $G(x) = \frac{1}{2}\|x\|_{L^2}^2$ , and  $K(x) = S(x) - z^\delta$ . (Note that in contrast to [11], we do not introduce a Moreau–Yosida regularization of  $F$  here.) Due to the properties of  $S$ , the gradient of  $K$  is uniformly bounded and Lipschitz continuous; cf. [13]. Hence, following Remark 4.5, we can expect that Assumption 4.1 (v) holds independent of the initialization. Furthermore,  $G$  and  $F^*$  are convex and hence Assumption 4.1 (iii) is satisfied for  $\tilde{\gamma}_G, \tilde{\gamma}_{F^*} \geq 0$ . This leaves Assumption 4.1 (iv), which amounts to quadratic growth condition of (5.2) near the minimizer; cf. Proposition 3.2. Similar assumptions are needed for the convergence of Newton-type methods, see, e.g., [21]. In the context of PDE-constrained optimization, they are generally difficult to prove a priori and have to be assumed. To set the initial step lengths, we estimate as in [11] the Lipschitz constant  $L$  by  $\tilde{L} = \max\{1, \|\nabla S'(u^0)u^0\|/\|u^0\|\} \approx 1$ . We then set  $\tau_0 = (4\tilde{L})^{-1}$  and  $\sigma_0 = (2\tilde{L})^{-1}$ . The starting points are chosen as  $x_0 \equiv 1$  and  $y^0 \equiv 0$  (which are not close to the expected saddle point). Figure 1 shows the convergence behavior  $\|x^N - \hat{x}\|_{L^2}^2$  of the primal iterates for  $N \in \{1, \dots, N_{\max}\}$  for  $N_{\max} = 10^4$ , both without and with acceleration. Since the exact minimizer to (5.2) is unavailable, here we take  $\hat{x} := x^{2N_{\max}}$  as an approximation. As can be seen, the convergence in the first case (corresponding to  $\tilde{\gamma}_G = 0$ ) is at best  $O(N^{-1})$ , while the accelerated algorithm according to Theorem 4.3 with  $\tilde{\gamma}_G = \frac{1}{2} < \gamma_G$  indeed eventually enters a region where the rate is  $O(N^{-2})$ . If we replace  $F$  by its Moreau–Yosida regularization  $F_\gamma$ , i.e., replace  $F^*$  by  $F_\gamma^* := F^* + \frac{\gamma}{2}\|\cdot\|_{L^2}^2$ , Theorem 4.4 is applicable for  $\tilde{\gamma}_{F^*} = \gamma > 0$ . As Figure 2 shows for different choices of  $\gamma$  and constant step sizes  $\tau = \sqrt{\tilde{\gamma}_{F^*}/\tilde{\gamma}_G\tilde{L}^{-1}}$ ,  $\sigma = (\tilde{\gamma}_G/\tilde{\gamma}_{F^*})\tau$ , the corresponding algorithm leads to linear convergence of the full iterates  $\|u^N - \hat{u}\|_{L^2 \times L^2}^2$  with a rate of  $(1 + 2\tilde{\gamma}_G\tau)^{-N}$  (which depends on  $\gamma$  by way of  $\tau$ ).

We also consider the example of optimal control with state constraints mentioned in the Introduction, i.e.,

$$(5.3) \quad \min_{x \in L^2} \frac{1}{2\alpha} \|S(x) - z^d\|_{L^2}^2 + \frac{1}{2} \|x\|_{L^2}^2 \quad \text{s. t.} \quad [S(x)](t) \leq c \quad \text{a. e. in } \Omega,$$

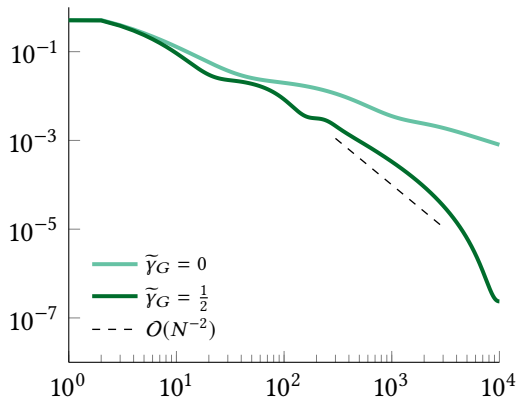


Figure 3: State constraints:  $\|x^N - \hat{x}\|_{L^2}^2$  for different values of  $\tilde{\gamma}_G$

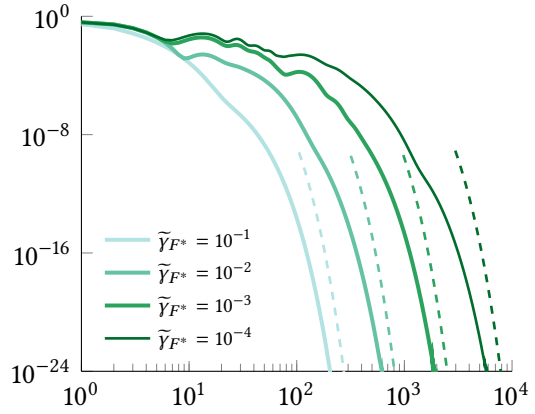


Figure 4: State constraints:  $\|u^N - \hat{u}\|_{L^2 \times L^2}^2$  (solid) and bounds  $(1 + 2\tilde{\gamma}_G\tau)^{-N}$  (dashed) for strongly convex  $F^*$  and different values of  $\tilde{\gamma}_{F^*}$

see [11, Section 3.3] for details. Here we choose  $z^d = S(x^\dagger)$  with  $x^\dagger$  as above,  $\alpha = 10^{-3}$ , and  $c = 0.68$  such that the state constraints are violated for  $z^d$ . Again, this fits into the framework of problem (P) with  $F(y) = \frac{1}{2\alpha}\|y - z^d\|_{L^2}^2 + \delta_{(-\infty, c]}(y)$ ,  $G(x) = \frac{1}{2}\|x\|_{L^2}^2$ , and  $K(x) = S(x)$ . With the same parameter choice as in the last example, we again eventually observe the convergence rate of  $O(N^{-2})$  for the accelerated algorithm (see Figure 3) as well as linear convergence if the state constraints are replaced by a Moreau–Yosida regularization (see Figure 4).

## 6 CONCLUSIONS

We have developed sufficient conditions on primal and dual step lengths that ensure (in some cases global) convergence and higher convergence rates of the NL-PDHGM method for nonsmooth nonconvex optimization. We have proved that usual acceleration rules give local  $O(1/N^2)$  convergence, justifying their use in previously published numerical examples [11]. Furthermore, we have derived novel linear convergence results based on bounds on the initial step lengths. Since our main derivations hold for general operators, one potential extension of the present work is to combine its approach with that of [23] to derive block-coordinate methods for nonconvex problems.

## ACKNOWLEDGMENTS

T. Valkonen and S. Mazurenko have been supported by the EPSRC First Grant EP/P021298/1, “PARTIAL Analysis of Relations in Tasks of Inversion for Algorithmic Leverage”. C. Clason is supported by the German Science Foundation (DFG) under grant Cl 487/1-1.

All data and source codes will be publicly deposited when the final accepted version of the manuscript is submitted.

## APPENDIX A A SMALL IMPROVEMENT OF OPIAL'S LEMMA

The earliest version of the next lemma is contained in the proof of [16, Theorem 1].

**Lemma A.1** ([6, Lemma 6]). *On a Hilbert space  $X$ , let  $\hat{X} \subset X$  be closed and convex, and  $\{x^i\}_{i \in \mathbb{N}} \subset X$ . Then  $x^i \rightharpoonup \bar{x}$  weakly in  $X$  for some  $\bar{x} \in \hat{X}$  if:*

- (i)  $i \mapsto \|x^i - \bar{x}\|$  is nonincreasing for all  $\bar{x} \in \hat{X}$ .
- (ii) All weak limit points of  $\{x^i\}_{i \in \mathbb{N}}$  belong to  $\hat{X}$ .

We can improve this result to the following

**Lemma A.2.** *Let  $X$  be a Hilbert space,  $\hat{X} \subset X$  (not necessarily closed or convex), and  $\{x^i\}_{i \in \mathbb{N}} \subset X$ . Also let  $A_i \in \mathbb{L}(X; X)$  be self-adjoint and  $A_i \geq \hat{\varepsilon}^2 I$  for some  $\hat{\varepsilon} \neq 0$  for all  $i \in \mathbb{N}$ . If the following conditions hold, then  $x^i \rightharpoonup \bar{x}$  weakly in  $X$  for some  $\bar{x} \in \hat{X}$ :*

- (i)  $i \mapsto \|x^i - \hat{x}\|_{A_i}$  is nonincreasing for some  $\hat{x} \in \hat{X}$ .
- (ii) All weak limit points of  $\{x^i\}_{i \in \mathbb{N}}$  belong to  $\hat{X}$ .
- (iii) There exists  $C$  such that  $\|A_i\| \leq C^2$  for all  $i$ , and for any weakly convergent subsequence  $x_{i_k}$  there exists  $A_\infty \in \mathbb{L}(X; X)$  such that  $A_{i_k} x \rightarrow A_\infty x$  strongly in  $X$  for all  $x \in X$ .

*Proof.* For  $x \in \text{cl conv } \hat{X}$ , define  $p(x) := \liminf_{i \rightarrow \infty} \|x - x^i\|_{A_i}$ . Clearly (i) yields

$$p(\hat{x}) = \lim_{i \rightarrow \infty} \|\hat{x} - x^i\|_{A_i} \in [0, \infty).$$

Using the triangle inequality and (iii), for any  $x, x' \in \text{cl conv } \hat{X}$  moreover

$$(A.1) \quad 0 \leq p(x) \leq p(x') + \limsup_{i \rightarrow \infty} \|x' - x\|_{A_i} \leq p(x') + C\|x' - x\|.$$

Choosing  $x' = \hat{x}$  we see from (A.1) that  $p$  is well-defined and finite. It is moreover bounded from below. Given  $\varepsilon > 0$ , we can therefore find  $x_\varepsilon^* \in \text{cl conv } \hat{X}$  such that  $p(x_\varepsilon^*)^2 - \varepsilon^2 \leq \inf_{\text{cl conv } \hat{X}} p^2$ . The norm  $\|x_\varepsilon^*\|$  is bounded from above for small values of  $\varepsilon$ : for the subsequence  $\{x_{i_k}\}$  realizing the limes inferior in  $p(x_\varepsilon^*)$ ,

$$\|x_\varepsilon^*\|_{A_{i_k}} \leq \|x_\varepsilon^* - x^{i_k}\|_{A_{i_k}} + \|x^{i_k} - \hat{x}\|_{A_{i_k}} + \|\hat{x}\|_{A_{i_k}},$$

and consequently

$$\hat{\varepsilon}\|x_\varepsilon^*\| \leq \left( \inf_{\text{cl conv } \hat{X}} p \right) + \varepsilon + \|x^0 - \hat{x}\|_{A_0} + C\|\hat{x}\|,$$

so there is a subsequence of  $\|x_\varepsilon^*\|$  weakly converging to some  $\bar{x}$  when  $\varepsilon \searrow 0$ . Without loss of generality, by restricting the allowed values of  $\varepsilon$ , we may assume that  $\bar{x}$  is unique.

Let  $\bar{x}'$  be some weak limit of  $\{x^i\}$ . By (ii),  $\bar{x}' \in \hat{X}$ . We have to show that  $\bar{x} = \bar{x}'$ . For simplicity of notation, we may assume that the whole sequence  $\{x^i\}$  converges weakly to  $\bar{x}'$ . By (iii), for any  $x \in X$ , we have

$$(A.2) \quad \lim_{i \rightarrow \infty} \langle x, \bar{x}_\varepsilon - x^i \rangle_{A_i} = \lim_{i \rightarrow \infty} \left( \langle x, \bar{x}_\varepsilon - x^i \rangle_{A_\infty} + \langle (A_i - A_\infty)x, \bar{x}_\varepsilon - x^i \rangle \right) = \langle x, \bar{x}_\varepsilon - \bar{x}' \rangle_{A_\infty}.$$

Moreover, for any  $\lambda \in (0, 1)$ , we have  $\bar{x}_{\varepsilon, \lambda} := (1 - \lambda)\bar{x}_\varepsilon + \lambda\bar{x}' \in \text{cl conv } \hat{X}$ . Now, since  $\bar{x}$  is a minimizer of  $p$  on  $\text{cl conv } \hat{X}$ , we can estimate

$$(A.3) \quad \begin{aligned} p(\bar{x}_\varepsilon)^2 - \varepsilon^2 &\leq p(\bar{x}_{\varepsilon, \lambda})^2 = p(\bar{x}_\varepsilon)^2 + \lim_{i \rightarrow \infty} \left( \lambda^2 \|\bar{x}_\varepsilon - \bar{x}'\|_{A_i}^2 - 2\lambda \langle \bar{x}_\varepsilon - \bar{x}', \bar{x}_\varepsilon - x^i \rangle_{A_i} \right) \\ &= p(\bar{x}_\varepsilon)^2 + (\lambda^2 - 2\lambda) \|\bar{x}_\varepsilon - \bar{x}'\|_{A_\infty}^2. \end{aligned}$$

In the second equality we have used (iii) and (A.2). Now, since  $\lambda^2 \leq 2\lambda$ , we obtain

$$0 \leq (2\lambda - \lambda^2) \|\bar{x}_\varepsilon - \bar{x}'\|_{A_\infty}^2 \leq \varepsilon^2.$$

This implies  $\bar{x}_\varepsilon \rightarrow \bar{x}'$  strongly as  $\varepsilon \searrow 0$ . But also  $\bar{x}_\varepsilon \rightarrow \bar{x}$ . Therefore  $\bar{x}' = \bar{x}$ .

Finally, by  $A_i \geq \hat{\varepsilon}I$  and (i), the sequence  $\{x^i\}$  is bounded, so any subsequence contains a weakly convergent subsequence. Since the limit is always  $\bar{x}$ , the whole sequence converges weakly to  $\bar{x}$ .  $\square$

**Remark A.3.** *The condition  $A_i \geq \hat{\varepsilon}^2 I$  is automatically satisfied if we replace (iii) by  $A_i \rightarrow A_\infty$  in the operator topology with  $A_\infty \geq 2\hat{\varepsilon}^2 I$ .*

## APPENDIX B RECONSTRUCTION OF THE PHASE AND AMPLITUDE OF A COMPLEX NUMBER

The purpose of this appendix is to verify [Assumption 3.3](#) for a simplified example related to the MRI reconstruction examples from [\[22\]](#). Consider

$$\min_{t, v \in \mathbb{R}} \frac{1}{2} |z - te^{iv}|^2 + G_0(t), \quad \text{where } G_0(t) := \begin{cases} \alpha t, & t \geq 0, \\ \infty, & t < 0. \end{cases}$$

for some data  $z \in \mathbb{C}$  and a regularization parameter  $\alpha > 0$ . We point out that the following does not depend on the specific structure of  $G_0$  for  $t \geq 0$ , as long as it is convex and increasing. In terms of real variables, this can be written in general saddle point form as

$$(B.1) \quad K(t, v) := \begin{pmatrix} t \cos v - \Re z \\ t \sin v - \Im z \end{pmatrix}, \quad G(t, v) := G_0(t), \quad \text{and } F^*(\lambda, \mu) := \frac{1}{2}(\lambda^2 + \mu^2).$$

To simplify the notation, let  $x = (t, v)$  and  $y = (\lambda, \mu)$ .

We now make a case distinction based on the sign of the optimal  $\hat{t} \geq 0$ . We first consider the case  $\hat{t} > 0$ .



**Lemma B.1.** Let  $\widehat{u} \in H^{-1}(0)$ , where  $H(u)$  is defined in (2.1) for  $K$ ,  $G$ , and  $F^*$  given by (B.1), and suppose  $\widehat{t} > 0$ . Let  $L > \widehat{\alpha}\widehat{t}/4$  as well as  $\theta > 0$  be arbitrary. Then Assumption 3.3 holds with  $p = 2$ , i.e., there exists  $\varepsilon > 0$  such that for all  $x, x' \in \mathbb{B}(\widehat{x}, \varepsilon)$ ,

$$(B.2) \quad \langle [\nabla K(x') - \nabla K(\widehat{x})]^* \widehat{y}, x - \widehat{x} \rangle \geq \theta \|K(\widehat{x}) - K(x) - \nabla K(x)(\widehat{x} - x)\|^2 - L(v - v')^2.$$

*Proof.* The saddle point condition  $0 \in H(\widehat{u})$  expands as  $K(\widehat{t}, \widehat{v}) \in \partial F^*(\widehat{\lambda}, \widehat{\mu})$  and  $-[\nabla K(\widehat{t}, \widehat{v})]^* \begin{pmatrix} \widehat{\lambda} \\ \widehat{\mu} \end{pmatrix} \in \partial G(\widehat{t}, \widehat{v})$ . Since

$$\nabla K(t, v) = \begin{pmatrix} \cos v & -t \sin v \\ \sin v & t \cos v \end{pmatrix} \quad \text{and} \quad [\nabla K(t, v)]^* \begin{pmatrix} \lambda \\ \mu \end{pmatrix} = \begin{pmatrix} \lambda \cos v + \mu \sin v \\ \mu t \cos v - \lambda t \sin v \end{pmatrix},$$

the latter further expands as

$$(B.3) \quad -(\widehat{\lambda} \cos \widehat{v} + \widehat{\mu} \sin \widehat{v}) \in \partial G_0(\widehat{t}), \quad \text{and} \quad \widehat{\mu} \widehat{t} \cos \widehat{v} = \widehat{\lambda} \widehat{t} \sin \widehat{v}.$$

From the second equality,  $\widehat{\mu} \cos \widehat{v} = \widehat{\lambda} \sin \widehat{v}$ . Since  $\partial G_0(t) = \alpha$  for  $t > 0$ , multiplying the first equality by  $\cos \widehat{v}$  and  $\sin \widehat{v}$  results in

$$(B.4) \quad \widehat{\lambda} = -\alpha \cos \widehat{v}, \quad \text{and} \quad \widehat{\mu} = -\alpha \sin \widehat{v}.$$

We can thus write the left-hand side of (B.2) as

$$\begin{aligned} d_1 &:= \langle [\nabla K(x') - \nabla K(\widehat{x})]^* \widehat{y}, x - \widehat{x} \rangle \\ &= \widehat{\lambda}(\cos v' - \cos \widehat{v})(t - \widehat{t}) + \widehat{\mu}(\sin v' - \sin \widehat{v})(t - \widehat{t}) \\ &\quad + \widehat{\lambda}(-t' \sin v' + \widehat{t} \sin \widehat{v})(v - \widehat{v}) + \widehat{\mu}(t' \cos v' - \widehat{t} \cos \widehat{v})(v - \widehat{v}) \\ &= -\alpha \left[ \cos \widehat{v}(\cos v' - \cos \widehat{v})(t - \widehat{t}) + \sin \widehat{v}(\sin v' - \sin \widehat{v})(t - \widehat{t}) \right. \\ &\quad \left. + \cos \widehat{v}(-t' \sin v' + \widehat{t} \sin \widehat{v})(v - \widehat{v}) + \sin \widehat{v}(t' \cos v' - \widehat{t} \cos \widehat{v})(v - \widehat{v}) \right] \\ &= -\alpha \left[ \cos \widehat{v}(\cos v' - \cos \widehat{v})(t - \widehat{t}) + \sin \widehat{v}(\sin v' - \sin \widehat{v})(t - \widehat{t}) \right. \\ &\quad \left. + t'[\sin \widehat{v} \cos v' - \cos \widehat{v} \sin v'](v - \widehat{v}) \right]. \end{aligned}$$

Using the standard trigonometric identities

$$(B.5A) \quad 2 \cos \widehat{v} \cos v' = \cos(\widehat{v} - v') + \cos(\widehat{v} + v'), \quad 2 \sin \widehat{v} \sin v' = \cos(\widehat{v} - v') - \cos(\widehat{v} + v'),$$

$$(B.5B) \quad 2 \sin \widehat{v} \cos v' = \sin(\widehat{v} + v') + \sin(\widehat{v} - v'), \quad 2 \cos \widehat{v} \sin v' = \sin(\widehat{v} + v') - \sin(\widehat{v} - v'),$$

as well as  $\cos^2 \widehat{v} + \sin^2 \widehat{v} = 1$ , this becomes

$$d_1 = -\alpha \left[ [\cos(\widehat{v} - v') - 1](t - \widehat{t}) + t' \sin(\widehat{v} - v')(v - \widehat{v}) \right].$$

Using Taylor expansion, we obtain for some  $\eta_1$  and  $\eta_2$  between 0 and  $\widehat{v} - v'$  that

$$\begin{aligned} d_1 &= \alpha \left[ \frac{\cos \eta_1}{2} (\widehat{v} - v')^2 (t - \widehat{t}) - t' \cos \eta_2 (\widehat{v} - v')(v - \widehat{v}) \right] \\ &= \alpha \left[ \frac{\cos \eta_1}{2} (\widehat{v} - v')^2 (t - \widehat{t}) + t' \cos \eta_2 (v - \widehat{v})^2 - t' \cos \eta_2 (v - v')(v - \widehat{v}) \right]. \end{aligned}$$

Note that  $\cos \eta_1, \cos \eta_2 \approx 1$  for  $v'$  close to  $\widehat{v}$ . Using Cauchy's inequality, we have for any  $\beta > 0$  that

$$(B.6) \quad d_1 \geq \alpha \left[ \frac{\cos \eta_1}{2} (\widehat{v} - v')^2 (t - \widehat{t}) + (1 - \beta) t' \cos \eta_2 (v - \widehat{v})^2 - \frac{\cos \eta_2}{4\beta} t' (v - v')^2 \right].$$

We also have  $\frac{\cos \eta_1}{2} (\widehat{v} - v')^2 (t - \widehat{t}) \geq -|t - \widehat{t}| [(v - \widehat{v})^2 + (v - v')^2]$  and hence

$$d_1 \geq \alpha \left[ (1 - \beta) t' \cos \eta_2 - |t - \widehat{t}| \right] (v - \widehat{v})^2 - \alpha \left[ \frac{\cos \eta_2}{4\beta} t' + |t - \widehat{t}| \right] (v - v')^2.$$

Choosing  $\varepsilon, \delta > 0$  small enough,  $\beta < 1$  large enough, and  $t' \in \mathbb{B}(\widehat{t}, \varepsilon)$ , we can thus ensure that

$$(B.7) \quad d_1 \geq \delta \theta (v - \widehat{v})^2 - L (v - v')^2.$$

We now turn to the right-hand side of (B.2), which we write as

$$\begin{aligned} D_2 &:= K(\widehat{x}) - K(x) - \nabla K(x)(\widehat{x} - x) \\ &= \begin{pmatrix} \widehat{t} \cos \widehat{v} - t \cos v - \cos v (\widehat{t} - t) + t \sin v (\widehat{v} - v) \\ \widehat{t} \sin \widehat{v} - t \sin v - \sin v (\widehat{t} - t) - t \cos v (\widehat{v} - v) \end{pmatrix} \\ &= \begin{pmatrix} \widehat{t} (\cos \widehat{v} - \cos v) + t \sin v (\widehat{v} - v) \\ \widehat{t} (\sin \widehat{v} - \sin v) - t \cos v (\widehat{v} - v) \end{pmatrix}. \end{aligned}$$

Thus

$$\begin{aligned} \|D_2\|^2 &= 2\widehat{t}^2 (1 - \cos \widehat{v} \cos v - \sin \widehat{v} \sin v) + t^2 (v - \widehat{v})^2 \\ &\quad + 2t\widehat{t}(\widehat{v} - v) [(\cos \widehat{v} - \cos v) \sin v - (\sin \widehat{v} - \sin v) \cos v]. \end{aligned}$$

Using the trigonometric identities (B.5) and Taylor expansion, it follows that

$$\begin{aligned} \|D_2\|^2 &= 2\widehat{t}^2 [1 - \cos(\widehat{v} - v)] + t^2 (v - \widehat{v})^2 - 2t\widehat{t}(\widehat{v} - v) \sin(\widehat{v} - v) \\ &\leq \widehat{t}^2 (\widehat{v} - v)^2 + t^2 (v - \widehat{v})^2 - 2t\widehat{t}(\widehat{v} - v)^2 + 2t\widehat{t}(\widehat{v} - v)^4 \\ &= (\widehat{t} - t)^2 (\widehat{v} - v)^2 + 2t\widehat{t}(\widehat{v} - v)^4. \end{aligned}$$

By taking  $\varepsilon > 0$  small enough and  $x = (t, v) \in \mathbb{B}(\widehat{x}, \varepsilon)$ , we thus obtain for any  $\delta > 0$  that  $\|D_2\|^2 \leq \delta (\widehat{v} - v)^2$ . We now obtain from (B.7) that

$$d_1 \geq \theta \|D_2\|^2 - L (v - v')^2,$$

which is exactly (B.2).  $\square$

The case of  $\widehat{t} = 0$  is complicated by the fact that  $\widehat{v}$  is then no longer unique. We therefore cannot expect convergence of  $x^i = (t^i, v^i)$  in the sense studied in this work; we would instead need to consider convergence to the entire solution set; cf. [24] for such an abstract approach for convex problems. However, under additional assumptions on the data  $z$ , we can proceed as before. The next lemma lays the groundwork for showing that the algorithm actually converges locally to  $\widehat{v} = v_z$  if  $t_z > 0$ .

**Lemma B.2.** Suppose  $\widehat{t} = 0$  and  $\widehat{v} = v_z$  for  $z = t_z e^{iv_z}$  with  $t_z > 0$ . Then the conclusions of [Lemma B.1](#) hold for some  $\varepsilon > 0$  and any  $\theta, L > 0$ .

*Proof.* First we deduce from the optimality condition  $K(\widehat{t}, \widehat{v}) \in \partial F^*(\widehat{\lambda}, \widehat{\mu})$  that

$$\widehat{\lambda} = -\Re z = -t_z \cos v_z \quad \text{and} \quad \widehat{\mu} = -\Im z = -t_z \sin v_z,$$

which is analogous to [\(B.4\)](#). Using the assumption that  $\widehat{v} = v_z$ , we can then proceed as in the proof of [Lemma B.1](#) to derive the estimate [\(B.6\)](#) with  $t_z$  in place of  $\alpha$ , which for  $\beta = 1$  reads

$$d_1 \geq t_z \left[ \frac{\cos \eta_1}{2} t (\widehat{v} - v')^2 - \frac{\cos \eta_2}{4} t' (v - v')^2 \right].$$

Now we have for any  $\zeta > 0$  that

$$(\widehat{v} - v')^2 = (v - \widehat{v})^2 + (v - v')^2 - 2(v - \widehat{v})(v - v') \geq (1 - \zeta)(v - \widehat{v})^2 - (\zeta^{-1} - 1)(v - v')^2$$

and therefore

$$d_1 \geq \frac{\cos \eta_1}{2} (1 - \zeta) t_z t (v - \widehat{v})^2 - t_z \left[ \frac{\cos \eta_2}{4} t' + \frac{\cos \eta_1}{2} (\zeta^{-1} - 1) t \right] (v - v')^2.$$

As in the proof of [Lemma B.1](#), we also have that  $\|D_2\|^2 = t^2 (\widehat{v} - v)^2$ . Now taking  $\zeta \in (0, 1)$  and  $\varepsilon > 0$  small enough, we can force  $0 = \widehat{t} \leq t \leq \varepsilon$  sufficiently small that  $\frac{1}{2}(1 - \zeta) \cos \eta_1 t_z t > \theta t^2$  for any given  $\theta > 0$ . Likewise, we can guarantee

$$t_z \left[ \frac{1}{4} t' \cos \eta_2 + \frac{\zeta^{-1} - 1}{2} t \cos \eta_1 \right] \leq L$$

for any  $\zeta > 0$  provided  $t', t \geq 0$  are small enough. We therefore obtain that  $d_1 \geq \theta \|D_2\|^2 - L(v - v')^2$  and hence the claim.  $\square$

## REFERENCES

- [1] BACHMAYR & BURGER, Iterative total variation schemes for nonlinear inverse problems, *Inverse Problems* 25 (2009), DOI: [10.1088/0266-5611/25/10/105004](https://doi.org/10.1088/0266-5611/25/10/105004).
- [2] BAUSCHKE & COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, 2017, DOI: [10.1007/978-3-319-48311-5](https://doi.org/10.1007/978-3-319-48311-5).
- [3] BECK & TEOULLE, 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](https://doi.org/10.1109/TIP.2009.2028250).
- [4] BENNING, KNOLL, SCHÖNLIEB & 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*, ed. by BOCIU, DÉSIDÉRI & HABBAL, Springer International Publishing, 2016, 117–126, DOI: [10.1007/978-3-319-55795-3\\_10](https://doi.org/10.1007/978-3-319-55795-3_10), ARXIV: [1511.00425](https://arxiv.org/abs/1511.00425),

- [5] BREDIES, KUNISCH & POCK, Total generalized variation, *SIAM Journal on Imaging Sciences* 3 (2011), 492–526, DOI: [10.1137/090769521](https://doi.org/10.1137/090769521).
- [6] BROWDER, Convergence theorems for sequences of nonlinear operators in Banach spaces, *Mathematische Zeitschrift* 100 (1967), 201–225, DOI: [10.1007/BF01109805](https://doi.org/10.1007/BF01109805).
- [7] CHAMBOLLE & 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).
- [8] CHAMBOLLE & LIONS, Image recovery via total variation minimization and related problems, *Numerische Mathematik* 76 (1997), 167–188, DOI: [10.1007/s002110050258](https://doi.org/10.1007/s002110050258).
- [9] CHAMBOLLE & 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).
- [10] CLASON & VALKONEN, Codes supporting: Primal-dual extragradient methods for nonlinear nonsmooth PDE-constrained optimization [nlpdegm], 2017, DOI: [10.5281/zenodo.398822](https://doi.org/10.5281/zenodo.398822).
- [11] CLASON & VALKONEN, Primal-dual extragradient methods for nonlinear nonsmooth PDE-constrained optimization, *SIAM Journal on Optimization* 27 (2017), 1313–1339, DOI: [10.1137/16M1080859](https://doi.org/10.1137/16M1080859), ARXIV: [1606.06219](https://arxiv.org/abs/1606.06219),
- [12] ESSER, ZHANG & CHAN, A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science, *SIAM Journal on Imaging Sciences* 3 (2010), 1015–1046, DOI: [10.1137/09076934X](https://doi.org/10.1137/09076934X).
- [13] KRÖNER & VEXLER, A priori error estimates for elliptic optimal control problems with a bilinear state equation, *J. Comput. Appl. Math.* 230 (2009), 781–802, DOI: [10.1016/j.cam.2009.01.023](https://doi.org/10.1016/j.cam.2009.01.023).
- [14] MÖLLENHOFF, STREKALOVSKIY, MOELLER & CREMERS, The primal-dual hybrid gradient method for semiconvex splittings, *SIAM Journal on Imaging Sciences* 8 (2015), 827–857, DOI: [10.1137/140976601](https://doi.org/10.1137/140976601).
- [15] OCHS, CHEN, BROX & POCK, iPiano: inertial proximal algorithm for nonconvex optimization, *SIAM Journal on Imaging Sciences* 7 (2014), 1388–1419, DOI: [10.1137/130942954](https://doi.org/10.1137/130942954).
- [16] 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](https://doi.org/10.1090/S0002-9904-1967-11761-0).
- [17] POCK, CREMERS, BISCHOF & CHAMBOLLE, An algorithm for minimizing the Mumford-Shah functional, in: *12th IEEE Conference on Computer Vision*, 2009, 1133–1140, DOI: [10.1109/ICCV.2009.5459348](https://doi.org/10.1109/ICCV.2009.5459348).
- [18] POCK & SABACH, Inertial proximal alternating linearized minimization (iPALM) for non-convex and nonsmooth problems, *SIAM Journal on Imaging Sciences* 9 (2016), 1756–1787, DOI: [10.1137/16M1064064](https://doi.org/10.1137/16M1064064).
- [19] ROCKAFELLAR, Monotone operators and the proximal point algorithm, *SIAM Journal on Optimization* 14 (1976), 877–898, DOI: [10.1137/0314056](https://doi.org/10.1137/0314056).

- [20] ROCKAFELLAR & WETS, *Variational Analysis*, Springer, 1998, DOI: [10.1007/978-3-642-02431-3](https://doi.org/10.1007/978-3-642-02431-3).
- [21] ULBRICH & ULBRICH, Non-monotone trust region methods for nonlinear equality constrained optimization without a penalty function, *Mathematical Programming* 95 (2003), 103–135, DOI: [10.1007/s10107-002-0343-9](https://doi.org/10.1007/s10107-002-0343-9).
- [22] 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](https://doi.org/10.1088/0266-5611/30/5/055012), ARXIV: [1309.5032](https://arxiv.org/abs/1309.5032),
- [23] VALKONEN, Block-proximal methods with spatially adapted acceleration (2017), ARXIV: [1609.07373](https://arxiv.org/abs/1609.07373), URL: [tuomov.iki.fi/m/blockcp.pdf](https://tuomov.iki.fi/m/blockcp.pdf).
- [24] VALKONEN, Preconditioned proximal point methods and notions of partial subregularity, 2017, ARXIV: [1711.05123](https://arxiv.org/abs/1711.05123), URL: [tuomov.iki.fi/m/subreg.pdf](https://tuomov.iki.fi/m/subreg.pdf).
- [25] VALKONEN, Testing and non-linear preconditioning of the proximal point method, 2017, ARXIV: [1703.05705](https://arxiv.org/abs/1703.05705), URL: [tuomov.iki.fi/m/proxtest.pdf](https://tuomov.iki.fi/m/proxtest.pdf).
- [26] WANG, YIN & ZENG, Global convergence of ADMM in nonconvex nonsmooth optimization, *Journal of Scientific Computing* (2018), DOI: [10.1007/s10915-018-0757-z](https://doi.org/10.1007/s10915-018-0757-z).