# Data-driven nonsmooth optimization

Sebastian Banert[*1], Axel Ringh[*1], Jonas Adler[1,2],
Johan Karlsson[1], and Ozan Öktem[1]

[1]KTH Royal Institute of Technology, 100 44 Stockholm, Sweden.
[2]Elekta, Box 7593, 103 93 Stockholm, Sweden.
Email: {banert, aringh, jonasadl, ozan}@kth.se,
johan.karlsson@math.kth.se

August 3, 2018

### Abstract

In this work, we consider methods for solving large-scale optimization problems with a possibly nonsmooth objective function. The key idea is to first specify a class of optimization algorithms using a generic iterative scheme involving only linear operations and applications of proximal operators. This scheme contains many modern primal-dual first-order solvers like the Douglas–Rachford and hybrid gradient methods as special cases. Moreover, we show convergence to an optimal point for a new method which also belongs to this class. Next, we interpret the generic scheme as a neural network and use unsupervised training to learn the best set of parameters for a specific class of objective functions while imposing a fixed number of iterations. In contrast to other approaches of "learning to optimize", we present an approach which learns parameters only in the set of convergent schemes. As use cases, we consider optimization problems arising in tomographic reconstruction and image deconvolution, and in particular a family of total variation regularization problems.

## 1 Introduction

Many problems in science and engineering can be formulated as convex optimization problems which then need to be solved accurately and efficiently. In this paper we focus on methods for solving such problems, namely of the form

$$\min_{x \in \mathcal{X}} \Big[ F(x) + \sum_{i=1}^{m} G_i(L_i x) \Big]. \tag{1.1}$$

Here, $L_i \colon \mathcal{X} \to \mathcal{Y}_i$, $i = 1, \ldots, m$, are linear operators, where $\mathcal{X}, \mathcal{Y}_1, \ldots, \mathcal{Y}_m$ are Hilbert spaces, and $F \colon \mathcal{X} \to \overline{\mathbb{R}}$ and $G_i \colon \mathcal{Y}_i \to \overline{\mathbb{R}}$, $i = 1, \ldots, m$, are proper, convex and lower semicontinuous functions. This class of optimization problems appears for example in variational regularization of inverse problems in imaging, such as

---

[*]equal contribution

1

X-ray computed tomography (CT) [40, 41], magnetic resonance imaging (MRI) [18], and electron tomography [42].

A key challenge is to handle the computational burden. In imaging, and especially so for three-dimensional imaging, the resulting optimization problem is very high-dimensional even after clever digitization and might involve more than one billion variables. Moreover, many regularizers that are popular in imaging (see Section 5), like those associated with sparsity, result in a nonsmooth objective function. These issues prevent usage of variational methods in time-critical applications, such as medical imaging in a clinical setting. Modern methods which aim at overcoming these obstacles are typically based on the proximal point algorithm [46] and operator splitting techniques, see e.g., [10, 12, 14–16, 20–22, 25, 29, 33, 34] and references therein.

The main objective of the paper is to offer a computationally tractable approach for minimizing large-scale nondifferentiable, convex functions. The key idea is to "learn" how to optimize from training data, resulting in an iterative scheme that is optimal given a fixed number of steps, while its convergence properties can be analyzed. We will make this precise in Section 4.

Similar ideas have been proposed previously in [8, 27, 35], but these approaches are either limited to specific classes of iterative schemes, like gradient-descent-like schemes [8, 35] that are not applicable for nonsmooth optimization, or specialized to a specific class of regularizers as in [27], which limits the possible choices of regularizers and forward operators. The approach taken here leverages upon these ideas and yields a general framework for learning optimization algorithms that are applicable to solving optimization problems of the type (1.1), inspired by the proximal-type methods mentioned above.

A key feature is to present a general formulation that includes several existing algorithms, among them the primal-dual hybrid gradient (PDHG) algorithm (also called the Chambolle–Pock algorithm) [20] and the primal-dual Douglas–Rachford algorithm [15] as a special case. This means that the learning can be done in a space of schemes that includes these solvers as special cases. Moreover, from the proposed parametrization we also derive a new optimization algorithm. We demonstrate the performance of a solver based on this general formulation by training in an unsupervised manner for two inverse problems: image reconstruction in CT and deconvolution, both through TV regularization. In particular, we present a method to learn the parameters of a convergent solver and demonstrate the improvement to the ad-hoc parameter choice. Moreover, empirical results indicate that by using additional parameters we can achieve improved performance.

The paper is organized as follows: In Section 2 we recall elements of monotone operator theory and convex optimization, while setting up the notation. In Section 3, we present and analyze a new solver for monotone inclusions, which in Section 3.3 is specialized to convex optimization problems of the form (1.1). Section 4 deals with the notion of "learning" an optimization solver, and in Section 5 we present numerical experiments for variational regularization of inverse problems in imaging.

2

# 2 Background

Solving optimization problems of the type in (1.1) are often addressed using *variable splitting* techniques, which work well if the different terms are "simple" [10, 22, 24]. To keep the discussion as general as possible and since it does not add complexity to the proof of convergence, we will carry it out for monotone inclusions instead of convex optimization problems. The following subsections present necessary background material on monotone operators, convex optimization, and variable splitting.

## 2.1 Fundamental notions

Let $\mathcal{H}$ be a real Hilbert space with the inner product $\langle \cdot, \cdot \rangle$. We denote convergence in norm (or strong convergence) and weak convergence by $\to$ and $\rightharpoonup$, respectively. A set-valued operator $S \colon \mathcal{H} \rightrightarrows \mathcal{H}$ is *monotone* if

$$\langle z - z', w - w' \rangle \geq 0 \quad \text{for all } z, z' \in \mathcal{H}, \ w \in S(z), \text{ and } w' \in S(z').$$

A monotone operator $S$ is called *maximally monotone* if, in addition, the graph of $S$, defined by $\operatorname{graph}(S) := \{(z, w) \in \mathcal{H} \times \mathcal{H} \mid w \in S(z)\}$, is not properly contained in the graph of any other monotone operator, i.e.,

$$(z, w) \in \operatorname{graph}(S) \iff \langle z - z', w - w' \rangle \geq 0 \text{ for all } (z', w') \in \operatorname{graph}(S).$$

A monotone operator is called *strongly monotone* if there exists a $\mu > 0$ such that

$$\langle z - z', w - w' \rangle \geq \mu \|z - z'\|^2 \quad \text{for all } z, z' \in \mathcal{H}, \ w \in S(z), \text{ and } w' \in S(z').$$

Next, for any scalar $\sigma > 0$, the operator $J_S^\sigma = (\operatorname{Id} + \sigma S)^{-1}$ is called the *resolvent operator* or *proximal mapping* [46]. It can be shown that $J_S^\sigma$ is a single-valued operator $\mathcal{H} \to \mathcal{H}$ [10, Proposition 23.8]. Note that an efficient routine to evaluate $J_S^\sigma$ for all $\sigma > 0$ also enables to evaluate the resolvent operator of $S^{-1}$ via

$$J_{S^{-1}}^\sigma(z) = z - \sigma J_S^{1/\sigma}(z/\sigma) \tag{2.1}$$

for $z \in \mathcal{H}$ (see [10, Proposition 23.20]).

A *maximally monotone inclusion problem* is defined as the problem of finding a point $z \in \mathcal{H}$ such that $0 \in S(z)$, which we henceforth denote $z \in \operatorname{zer}(S)$. In fact, it is easily seen that $z \in \operatorname{zer}(S)$ is equivalent with $z$ being a fixed-point for the resolvent operator, i.e., $z = J_S^\sigma(z)$.

One reason for the interest in maximally monotone inclusion problems is that the *subdifferential* $\partial F$ of a proper, convex and lower semicontinuous function $F \colon \mathcal{H} \to \overline{\mathbb{R}}$ is a maximally monotone operator [39]. Here, $\partial F \colon \mathcal{H} \rightrightarrows \mathcal{H}$ is defined to be

$$\partial F(x) := \{y \in \mathcal{H} \mid \forall \tilde{x} \in \mathcal{H} \colon F(\tilde{x}) \geq F(x) + \langle y, \tilde{x} - x \rangle\}$$

if $F(x) \in \mathbb{R}$ and $\partial F(x) = \emptyset$ if $F(x) \in \{\pm\infty\}$. Moreover, the subdifferential at any minimizer of such a function contains zero, so $F$ can be minimized by solving a maximally monotone inclusion problem [10, Theorem 16.3]. Note that we do not distinguish between local and global minimizers, since any local minimizer of a convex function is global [10, Proposition 11.4].

*Remark* 2.1. A continuous linear operator $A \colon \mathcal{H} \to \mathcal{H}$ of a Hilbert space $\mathcal{H}$ into itself is maximally monotone if and only if it is accretive, i.e., if $\langle x, Ax \rangle \geq 0$ for all $x \in \mathcal{H}$ [10, Corollary 20.28, see also Definition 2.23], and it is the subdifferential $\partial f$ of a function $f \colon \mathcal{H} \to \overline{\mathbb{R}}$ if and only if it is additionally symmetric [9, Proposition 2.51]. In particular the Volterra integral operator [11, Example 4.4]

$$(Af)(t) = \int_0^t f(s)\,\mathrm{d}s$$

and its inverse are maximally monotone, but not the subdifferential of a proper, convex and lower semicontinuous function.

For $F \colon \mathcal{H} \to \overline{\mathbb{R}}$, the Fenchel dual (convex conjugate) function $F^* \colon \mathcal{H} \to \overline{\mathbb{R}}$ is defined by [10, Chapter 13]

$$F^*(y) := \sup_{x \in \mathcal{H}} \Big[ \langle x, y \rangle - F(x) \Big] \quad \text{for } y \in \mathcal{H}.$$

If $F$ is proper, convex and lower semicontinuous, then $\partial F^* = (\partial F)^{-1}$ [10, Corollary 16.30].

The *proximal point algorithm* is a fixed-point iterative scheme for solving the maximally monotone inclusion problem. It is given by repeatedly applying the resolvent operator:

$$z^{k+1} = J_S^\sigma\big(z^k\big).$$

It can now be shown that if $\mathrm{zer}(S) \neq \emptyset$ then $z^k$ converges weakly to a point $z^\infty \in \mathrm{zer}(S)$ [46] for all starting points $z^0 \in \mathcal{H}$. The special case when $S := \partial F$, i.e., the case of the resolvent of a subdifferential of $F$, is called the *proximal operator*. One can express the proximal as [39]

$$J_{\partial F}^\sigma(x) = \mathrm{Prox}_F^\sigma(x) = \arg\min_{x' \in \mathcal{H}} \left\{ F(x') + \frac{1}{2\sigma} \|x' - x\|^2 \right\}. \qquad (2.2)$$

To see this, we simply note that if $x'$ is a minimizing argument then

$$0 \in \partial F(x') + \frac{1}{\sigma}(x' - x) \quad \Longleftrightarrow \quad x' = J_{\partial F}^\sigma(x).$$

It is thus interesting to note that the fixed-point iteration

$$x^{k+1} = \mathrm{Prox}_F^\sigma\big(x^k\big) = \arg\min_{x' \in \mathcal{H}} \left\{ F(x') + \frac{1}{2\sigma} \|x' - x^k\|^2 \right\}$$

generates a sequence $\big(x^k\big)$ that converges weakly to a minimizer of $F$. In this setting, the parameter $\sigma$ can be interpreted as a step length. This can give rise to methods for solving the optimization problems if the proximal operator can be efficiently computed, e.g., through a closed-form expression. Note that (2.1) gives a method to obtain the proximal points of $F^*$ from those of $F$, namely

$$\mathrm{Prox}_{F^*}^\tau(x) = x - \tau\,\mathrm{Prox}_F^{1/\tau}(x/\tau) \qquad \text{for all } \tau > 0.$$

Sometimes the resolvent of the maximally monotone operator $S$ is not easy to evaluate, but $S$ is of the form $S = A + B$ where $A$ and $B$ are maximally monotone and the resolvents of $A$ and $B$ can be evaluated efficiently. One may

then consider approximating $J^\sigma_{A+B}$ with $J^\sigma_A$ and $J^\sigma_B$ (splitting) [24]. An example when this arises is in convex minimization of an objective that is a sum of two (or more) functions $F + G$, like in (1.1). In these cases it is often not possible to compute a closed-form expression for the proximal operator $\text{Prox}^\sigma_{F+G}$. Such problems can be addressed using operator splitting techniques that allow for solving the problem by only evaluating $\text{Prox}^\sigma_F$ and $\text{Prox}^\sigma_G$ [22].*

## 2.2 Convex optimization

Next, we will consider duality and optimality conditions for the problem (1.1). To simplify the notation, we consider the case $m = 1$ in (1.1), i.e., let $\mathcal{X}$ and $\mathcal{Y}$ be two Hilbert spaces and consider the model problem

$$\min_{x \in \mathcal{X}} \Big[ F(x) + G(Lx) \Big], \tag{2.3}$$

where $L \colon \mathcal{X} \to \mathcal{Y}$ is a continuous linear operator and $F \colon \mathcal{X} \to \overline{\mathbb{R}}$ and $G \colon \mathcal{Y} \to \overline{\mathbb{R}}$ are proper, convex and lower semicontinuous functions. Note that (1.1) is recovered by setting

$$G(y) := \sum_{i=1}^m G_i(y_i) \quad \text{for } y = (y_1, \dots, y_m) \in \mathcal{Y} := \mathcal{Y}_1 \times \dots \times \mathcal{Y}_m \tag{2.4}$$

and $Lx := (L_1 x, \dots, L_m x)$ for $x \in \mathcal{X}$ in (2.3).

The dual formulation of the primal problem (2.3) is

$$\max_{y \in \mathcal{Y}} \Big[ -F^*(L^* y) - G^*(-y) \Big]. \tag{2.5}$$

Under suitable conditions the two optimization problems (2.3) and (2.5) have the same optimal value [10, Chapter 15.3]. Also note that, since both $F$ and $G$ are proper, convex and lower semicontinuous functions, $F^{**} = F$ and $G^{**} = G$ by the Fenchel–Moreau theorem [10, Theorem 13.37]. Hence, the following primal-dual formulation

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \ \mathcal{L}(x; y) \quad \text{with} \quad \mathcal{L}(x; y) := \langle Lx, y \rangle + F(x) - G^*(y) \tag{2.6}$$

(the mapping $\mathcal{L}(\cdot; \cdot)$ is called the *Lagrangian*) is equivalent to the primal problem.[†] In fact, under suitable assumptions it can be shown that if $(\bar{x}, \bar{y})$ is a saddle point to (2.6), then $\bar{x}$ is a solution to the primal problem (2.3) and $\bar{y}$ is a solution to the dual problem (2.5) [10, Proposition 19.20].

A necessary optimality condition for the primal-dual formulation (2.6) is that the corresponding point $(\bar{x}, \bar{y}) \in \mathcal{X} \times \mathcal{Y}$ be stationarity with respect to both variables, i.e., that

$$L\bar{x} \in \partial G^*(\bar{y}) \quad \text{and} \quad -L^* \bar{y} \in \partial F(\bar{x}). \tag{2.7}$$

For later use we note that the first of these conditions can be reformulated as

$$L\bar{x} \in \partial G^*(\bar{y}) \quad \Longleftrightarrow \quad \bar{y} + \sigma L\bar{x} \in \bar{y} + \sigma \partial G^*(\bar{y}) = (I + \sigma \partial G^*)(\bar{y})$$
$$\Longleftrightarrow \quad \bar{y} = J^\sigma_{\partial G^*}(\bar{y} + \sigma L\bar{x}) = \text{Prox}^\sigma_{G^*}(\bar{y} + \sigma L\bar{x}),$$

---

*In optimization, this operator splitting is sometimes referred to as *variable splitting*. The reason for this can be understood by comparing equations (2.3) and (2.8) below.

†To see this, note that $\max_{y \in \mathcal{Y}} \big[ \langle Lx, y \rangle - G^*(y) \big] = G^{**}(Lx) = G(Lx)$.

5

and the second as

$$-L^*\bar{y} \in \partial F(\bar{x}) \quad \Longleftrightarrow \quad \bar{x} - \tau L^*\bar{y} \in \bar{x} + \tau\partial F(\bar{x}) = (I + \tau\partial F)(\bar{x})$$
$$\Longleftrightarrow \quad \bar{x} = J^\tau_{\partial F}(\bar{x} - \tau L^*\bar{y}) = \mathrm{Prox}^\tau_F(\bar{x} - \tau L^*\bar{y}).$$

Therefore, an equivalent condition to (2.7) is

$$\bar{y} = \mathrm{Prox}^\sigma_{G^*}(\bar{y} + \sigma L\bar{x}) \quad \text{and} \quad \bar{x} = \mathrm{Prox}^\tau_F(\bar{x} - \tau L^*\bar{y}). \tag{2.8}$$

## 2.3 Two splitting algorithms

As mentioned before, there are many different splitting methods available to solve problems of the form (1.1). For ease of reference, we here mention two popular choices. The first one, given in (2.9), is PDHG [20]

$$\begin{aligned}
y_{n+1} &= \mathrm{Prox}^\sigma_{G^*}(y_n + \sigma L v_n), \\
x_{n+1} &= \mathrm{Prox}^\tau_F(x_n - \tau L^* y_{n+1}), \\
v_{n+1} &= x_{n+1} + \theta(x_{n+1} - x_n).
\end{aligned} \tag{2.9}$$

The second one is the Douglas–Rachford type primal-dual algorithm [15], presented in (2.10)

$$\begin{aligned}
p_n &= \mathrm{Prox}^\tau_F(x_n - \tau L^* y_n), \\
x_{n+1} &= x_n + \lambda_n(p_n - x_n), \\
q_n &= \mathrm{Prox}^\sigma_{G^*}(y_n + \sigma L(2p_n - x_n)), \\
y_{n+1} &= y_n + \lambda_n(q_n - y_n).
\end{aligned} \tag{2.10}$$

# 3 A new family of optimization solvers

In this section we introduce a new family of optimization algorithms and prove convergence for a subfamily. For ease of notation we will consider the simplified optimization problem (2.3), but results easily extend to the general case (1.1).

To this end, consider the two algorithms (2.9) and (2.10). Note that they can both be written as

$$q_n = \mathrm{Prox}^\sigma_{G^*}(b_{12}y_n + b_{11}L(c_{11}p_{n-1} + c_{12}x_{n-1})), \tag{3.1a}$$
$$y_{n+1} = a_{21}q_n + a_{22}y_n, \tag{3.1b}$$
$$p_n = \mathrm{Prox}^\tau_F(d_{12}x_n + d_{11}L^*(a_{11}q_n + a_{12}y_n)), \tag{3.1c}$$
$$x_{n+1} = c_{21}p_n + c_{22}x_n, \tag{3.1d}$$

for suitable values of the coefficients. More precisely, the PDHG algorithm (2.9) is obtained by setting

| | | | | | |
|---|---|---|---|---|---|
| $a_{11} = 1$ | $a_{12} = 0$ | $a_{21} = 1$ | $a_{22} = 0$ | $b_{11} = \sigma$ | $b_{12} = 1$ |
| $c_{11} = 1 + \theta$ | $c_{12} = -\theta$ | $c_{21} = 1$ | $c_{22} = 0$ | $d_{11} = -\tau$ | $d_{12} = 1$ |

and the Douglas-Rachford algorithm (2.10) by setting

| | | | | | |
|---|---|---|---|---|---|
| $a_{11} = \lambda_n$ | $a_{12} = 1 - \lambda_n$ | $a_{21} = \lambda_n$ | $a_{22} = 1 - \lambda_n$ | $b_{11} = \sigma$ | $b_{12} = 1$ |
| $c_{11} = 2$ | $c_{12} = -1$ | $c_{21} = \lambda_n$ | $c_{22} = 1 - \lambda_n$ | $d_{11} = -\tau$ | $d_{12} = 1.$ |

We now go on to analyze the scheme (3.1). To state our results as generally as possible, we formulate them for a monotone inclusion problem that in particular specializes to the optimality conditions in (2.7) when the operators are subdifferentials. The monotone inclusion problem we seek to solve reads as follows: Let $\mathcal{X}$ and $\mathcal{Y}$ be two (not necessarily finite-dimensional) Hilbert spaces, and let $L\colon \mathcal{X} \to \mathcal{Y}$ be a continuous linear operator. Let $A\colon \mathcal{X} \rightrightarrows \mathcal{X}$ and $B\colon \mathcal{Y} \rightrightarrows \mathcal{Y}$ be maximally monotone operators. Find a pair $(\bar{x}, \bar{y}) \in \mathcal{X} \times \mathcal{Y}$ such that

$$L\bar{x} \in B^{-1}\bar{y} \quad \text{and} \quad -L^*\bar{y} \in A\bar{x}. \tag{3.2}$$

In this setting, the scheme (3.1) generalizes to

$$q_n = J_{B^{-1}}^{\sigma}(b_{12}y_n + b_{11}L(c_{11}p_{n-1} + c_{12}x_{n-1})), \tag{3.3a}$$

$$y_{n+1} = a_{21}q_n + a_{22}y_n, \tag{3.3b}$$

$$p_n = J_A^{\tau}(d_{12}x_n + d_{11}L^*(a_{11}q_n + a_{12}y_n)), \tag{3.3c}$$

$$x_{n+1} = c_{21}p_n + c_{22}x_n. \tag{3.3d}$$

We first note that if $a_{21} = 0$ or $c_{21} = 0$ the update for either $y_{n+1}$ or $x_{n+1}$ becomes trivial, and the algorithm will not be globally convergent to a point fulfilling (3.2) in general. Henceforth we will therefore assume that $a_{21}$ and $c_{21}$ are not equal to 0, unless the opposite is explicitly stated.

## 3.1 Fixed-point analysis

In this section, we give necessary and sufficient conditions for the solution set of (3.2) and the fixed point set of (3.3) to coincide for any choice of $A$, $B$, and $L$. To this end, let $(\bar{q}, \bar{y}, \bar{p}, \bar{x}) \in \mathcal{Y} \times \mathcal{Y} \times \mathcal{X} \times \mathcal{X}$ be a fixed point of the iterative scheme (3.3) and note that (3.3b) and (3.3d) gives

$$\bar{q} = \frac{1 - a_{22}}{a_{21}}\bar{y} \quad \text{and} \quad \bar{p} = \frac{1 - c_{22}}{c_{21}}\bar{x}.$$

Using this, we further get that

$$\frac{1 - a_{22}}{a_{21}}\bar{y} = J_{B^{-1}}^{\sigma}\left(b_{12}\bar{y} + b_{11}L\left(c_{11}\frac{1 - c_{22}}{c_{21}}\bar{x} + c_{12}\bar{x}\right)\right)$$

$$\frac{1 - c_{22}}{c_{21}}\bar{x} = J_A^{\tau}\left(d_{12}\bar{x} + d_{11}L^*\left(a_{11}\frac{1 - a_{22}}{a_{21}}\bar{y} + a_{12}\bar{y}\right)\right)$$

The conditions in (3.2) can now be re-phrased as

$$\bar{y} = J_{B^{-1}}^{\sigma}(\bar{y} + \sigma L\bar{x}) \quad \text{and} \quad \bar{x} = J_A^{\tau}(\bar{x} - \tau L^*\bar{y}),$$

and combining the above two equations yields

$$\begin{aligned} a_{21} + a_{22} = 1, \quad b_{12} = 1, \quad b_{11}(c_{11} + c_{12}) = \sigma, \\ c_{21} + c_{22} = 1, \quad d_{12} = 1, \quad d_{11}(a_{11} + a_{12}) = -\tau. \end{aligned} \tag{3.4}$$

The conditions in (3.4) are necessary and sufficient, however, due to the linearity of $L$, the algorithm does not change if we agree to the normalization

$$\begin{aligned} b_{11} &= \sigma, & c_{11} + c_{12} &= 1, \\ d_{11} &= -\tau, & a_{11} + a_{12} &= 1. \end{aligned}$$

If we fix all these conditions, the iteration (3.3) takes the form

$$q_n = J_{B^{-1}}^{\sigma}(y_n + \sigma L(x_{n-1} + c_{11}(p_{n-1} - x_{n-1}))), \tag{3.5a}$$

$$y_{n+1} = y_n + a_{21}(q_n - y_n), \tag{3.5b}$$

$$p_n = J_A^{\tau}(x_n - \tau L^*(y_n + a_{11}(q_n - y_n))), \tag{3.5c}$$

$$x_{n+1} = x_n + c_{21}(p_n - x_n). \tag{3.5d}$$

## 3.2 Convergence analysis

The following theorem gives sufficient conditions for the weak convergence of the sequence $(x_n, y_n)$ generated by (3.5) to a point that satisfies (3.2), i.e., a point that solves the monotone inclusion problem.

**Theorem 3.1.** *Assume that there is a point that satisfies (3.2), i.e., the monotone inclusion problem has a solution. Moreover, let*

$$a_{11} = a_{21} \quad and \quad c_{11} = 1 + \frac{c_{21}}{a_{21}}. \tag{3.6}$$

*Assume furthermore that $0 < a_{21} < 2$, $0 < c_{21} < 2$ and*

$$\sigma\tau\|L\|^2 < \frac{a_{21}^2(2 - a_{21})(2 - c_{21})}{(a_{21} + c_{21} - a_{21}c_{21})^2} \quad with \ \sigma, \tau > 0. \tag{3.7}$$

*Finally, let $(q_n, y_n, p_n, x_n)$ be the sequence generated by scheme (3.5). Then the following holds:*

*(a)* $\displaystyle\sum_{n\geq 0}\|x_n - p_n\|^2 < +\infty$ *and* $\displaystyle\sum_{n\geq 0}\|x_n - x_{n+1}\|^2 < +\infty.$

*(b)* $\displaystyle\sum_{n\geq 0}\|y_n - q_n\|^2 < +\infty$ *and* $\displaystyle\sum_{n\geq 0}\|y_n - y_{n+1}\|^2 < +\infty.$

*(c) The sequence $(x_n, y_n)_n$ converges weakly to a point that satisfies (3.2).*

*(d) If $A$ is strongly monotone, then there is a unique $\bar{x} \in \mathcal{X}$ such that all solutions of (3.2) are of the form $(\bar{x}, y)$ with some $y \in \mathcal{Y}$. Moreover, $\displaystyle\sum_{n=1}^{\infty}\|p_n - \bar{x}\|^2 < +\infty$, in particular $p_n \to \bar{x}$ strongly.*
*If $B^{-1}$ is strongly monotone, then there is a unique $\bar{y} \in \mathcal{Y}$ such that all solutions of (3.2) are of the form $(x, \bar{y})$ with some $x \in \mathcal{X}$. Moreover, $\displaystyle\sum_{n=1}^{\infty}\|q_{n+1} - \bar{y}\|^2 < +\infty$, in particular $q_n \to \bar{y}$ strongly.*

By rewriting with (3.6), the iteration (3.5) takes the following form:

**Algorithm 3.2.** Choose parameters $\sigma > 0$, $\tau > 0$ and $a_{21} \in \mathbb{R}$, $c_{21} \in \mathbb{R}$ and starting points $x_0 \in \mathcal{X}$, $x_1 \in \mathcal{X}$, $p_0 \in \mathcal{X}$, $y_1 \in \mathcal{Y}$. For all $n = 1, 2, \ldots$, calculate

$$q_n = J_{B^{-1}}^{\sigma}\left(y_n + \sigma L\left(p_{n-1} + \frac{c_{21}}{a_{21}}(p_{n-1} - x_{n-1})\right)\right), \tag{3.8a}$$

$$y_{n+1} = y_n + a_{21}(q_n - y_n), \tag{3.8b}$$

$$p_n = J_A^{\tau}(x_n - \tau L^* y_{n+1}), \tag{3.8c}$$

$$x_{n+1} = x_n + c_{21}(p_n - x_n). \tag{3.8d}$$

Then, $x_n \rightharpoonup \bar{x}$, $p_n \rightharpoonup \bar{x}$, $y_n \rightharpoonup \bar{y}$, and $q_n \rightharpoonup \bar{y}$, where $(\bar{x}, \bar{y})$ is a solution of (3.2), provided that $0 < a_{21} < 2$, $0 < c_{21} < 2$ and (3.7) are satisfied.

The remainder of the convergence analysis will therefore refer to scheme (3.8). The proof of Theorem 3.1 rests upon a number of technical results and is given in Section 3.2.1. An immediate corollary is the convergence of the primal-dual Douglas–Rachford method with constant relaxation [15].

**Corollary 3.3.** *Let $\sigma\tau\|L\| < 1$ and $0 < \lambda < 2$. Then, for the iteration*

$$
\begin{aligned}
q_n &= J_{B^{-1}}^{\sigma}(y_n + \sigma L(2p_{n-1} - x_{n-1})), \\
y_{n+1} &= y_n + \lambda(q_n - y_n), \\
p_n &= J_A^{\tau}(x_n - \tau L^* y_{n+1}), \\
x_{n+1} &= x_n + \lambda(p_n - x_n),
\end{aligned}
$$

*the sequence $(x_n, y_n)_n$ converges weakly to a point that satisfies (3.2).*

*Proof.* Set $a_{21} = c_{21} = \lambda$ in Theorem 3.1 and observe that (3.7) reduces to $\sigma\tau\|L\|^2 < 1$. $\qquad\square$

### 3.2.1 Proof of Theorem 3.1

For the proof, we define notions of distance $Q_1$ and $Q_2$ on the space $\mathcal{X} \times \mathcal{Y}$ of pairs of primal and dual variables (Lemma 3.5). Next, we show that the distance (in terms of $Q_1$) between the iterates and the set of solutions of (3.2) decreases (Proposition 3.6). This property is also known as *Fejér monotonicity* [10, Chapter 5]. Proposition 3.7 improves the statement of Proposition 3.6 for strongly monotone operators. The proof of Theorem 3.1 is completed by showing that any weak sequential cluster point of the iteration sequence is a solution to (3.2).

We start with some simple inequalities between real numbers. In particular, Lemma 3.4 (a) shows that we do not divide by zero in (3.7).

**Lemma 3.4.** *Let $0 < a_{21} < 2$ and $0 < c_{21} < 2$. Then*

*(a) $a_{21} + c_{21} > a_{21}c_{21}$ and*

*(b) $\dfrac{a_{21}c_{21}(2 - a_{21})(2 - c_{21})}{(a_{21} + c_{21} - a_{21}c_{21})^2} \le 1$.*

*Proof.* By assumption, $a_{21}(2 - a_{21}) > 0$, i.e., $a_{21} > \frac{1}{2}a_{21}^2$, and the same holds for $c_{21}$. Therefore,

$$
a_{21} + c_{21} > \frac{1}{2}a_{21}^2 + \frac{1}{2}c_{21}^2 \ge a_{21}c_{21},
$$

whence (a).

For (b), use the inequality $2a_{21}c_{21} \le a_{21}^2 + c_{21}^2$ in

$$
\begin{aligned}
a_{21}c_{21}(2 - a_{21})(2 - c_{21}) &= 4a_{21}c_{21} - 2a_{21}^2 c_{21} - 2a_{21}c_{21}^2 + a_{21}^2 c_{21}^2 \\
&\le a_{21}^2 + c_{21}^2 + 2a_{21}c_{21} - 2a_{21}^2 c_{21} - 2a_{21}c_{21}^2 + a_{21}^2 c_{21}^2 \\
&= (a_{21} + c_{21} - a_{21}c_{21})^2. \qquad\square
\end{aligned}
$$

**Lemma 3.5.** *Define the quadratic forms $Q_1, Q_2 \colon \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ by*

$$Q_1(x, y) = \frac{1}{2\tau c_{21}}\|x\|^2 + \frac{1}{2\sigma a_{21}}\|y\|^2 - \frac{1}{a_{21}}\langle y, Lx\rangle,$$

$$Q_2(x, y) = \frac{2 - c_{21}}{2\tau}\|x\|^2 + \frac{2 - a_{21}}{2\sigma}\|y\|^2 - \frac{a_{21} + c_{21} - a_{21}c_{21}}{a_{21}}\langle y, Lx\rangle$$

*for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Under the assumptions in Theorem 3.1, there exist $C_1, C_2, D_1, D_2 > 0$ such that*

$$Q_i(x, y) \geq C_i\|x\|^2 \quad and \quad Q_i(x, y) \geq D_i\|y\|^2$$

*for all $x \in \mathcal{X}$, $y \in \mathcal{Y}$ and $i = 1, 2$.*

*Proof.* We can rewrite

$$Q_1(x, y) = \frac{1}{2\sigma a_{21}}\|y - \sigma Lx\|^2 + \frac{1}{2\tau c_{21}}\|x\|^2 - \frac{\sigma}{2a_{21}}\|Lx\|^2,$$

$$Q_1(x, y) = \frac{1}{2\tau c_{21}}\left\|x - \frac{c_{21}\tau}{a_{21}}L^*y\right\|^2 + \frac{1}{2\sigma a_{21}}\|y\|^2 - \frac{c_{21}\tau}{2a_{21}^2}\|L^*y\|^2$$

and

$$Q_2(x, y) = \frac{2 - a_{21}}{2\sigma}\left\|y - \frac{\sigma(a_{21} + c_{21} - a_{21}c_{21})}{a_{21}(2 - a_{21})}Lx\right\|^2$$

$$+ \frac{2 - c_{21}}{2\tau}\|x\|^2 - \frac{\sigma(a_{21} + c_{21} - a_{21}c_{21})^2}{2a_{21}^2(2 - a_{21})}\|Lx\|^2,$$

$$Q_2(x, y) = \frac{2 - c_{21}}{2\tau}\left\|x - \frac{\tau(a_{21} + c_{21} - a_{21}c_{21})}{a_{21}(2 - c_{21})}L^*y\right\|^2$$

$$+ \frac{2 - a_{21}}{2\sigma}\|y\|^2 - \frac{\tau(a_{21} + c_{21} - a_{21}c_{21})^2}{2a_{21}^2(2 - c_{21})}\|L^*y\|^2.$$

From this, the assertion of the lemma is clear with the quantities

$$C_1 = \frac{1}{2\tau c_{21}} - \frac{\sigma}{2a_{21}}\|L\|^2 = \frac{a_{21} - c_{21}\sigma\tau\|L\|^2}{2\tau a_{21}c_{21}},$$

$$D_1 = \frac{1}{2\sigma a_{21}} - \frac{c_{21}\tau}{2a_{21}^2}\|L\|^2 = \frac{a_{21} - c_{21}\sigma\tau\|L\|^2}{2\sigma a_{21}^2},$$

$$C_2 = \frac{2 - c_{21}}{2\tau} - \frac{\sigma(a_{21} + c_{21} - a_{21}c_{21})^2}{2a_{21}^2(2 - a_{21})}\|L\|^2$$

$$= \frac{a_{21}^2(2 - a_{21})(2 - c_{21}) - (a_{21} + c_{21} - a_{21}c_{21})^2\sigma\tau\|L\|^2}{2\tau a_{21}^2(2 - a_{21})},$$

$$D_2 = \frac{2 - a_{21}}{2\sigma} - \frac{\tau(a_{21} + c_{21} - a_{21}c_{21})^2}{2a_{21}^2(2 - c_{21})}\|L\|^2$$

$$= \frac{a_{21}^2(2 - a_{21})(2 - c_{21}) - (a_{21} + c_{21} - a_{21}c_{21})^2\sigma\tau\|L\|^2}{2\sigma a_{21}^2(2 - c_{21})}$$

provided that the numerators are positive, i.e.,

$$\sigma\tau\|L\|^2 < \min\left\{\frac{a_{21}}{c_{21}}, \frac{a_{21}^2(2 - a_{21})(2 - c_{21})}{(a_{21} + c_{21} - a_{21}c_{21})^2}\right\}.$$

Now, by Lemma 3.4, the minimum is always attained by the second value, and positivity is guaranteed by (3.7). $\qquad\square$

**Proposition 3.6.** *Define $Q_1$ and $Q_2$ as in Lemma 3.5, let $(\bar{x}, \bar{y}) \in \mathcal{X} \times \mathcal{Y}$ satisfy (3.2), and let the sequence $(q_n, y_n, p_n, x_n)$ be generated by scheme (3.8). Under the assumptions in Theorem 3.1, we have for all $n \geq 1$*

$$Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}) \leq -Q_2(p_n - x_n, q_{n+1} - y_{n+1}).$$

*Proof.* Let $(\bar{x}, \bar{y})$ satisfy (3.2). Then

$$\begin{aligned}
&Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}) \\
&= \frac{1}{2\tau c_{21}} \left( \|x_{n+1} - \bar{x}\|^2 - \|x_n - \bar{x}\|^2 \right) \\
&\quad + \frac{1}{2\sigma a_{21}} \left( \|y_{n+2} - \bar{y}\|^2 - \|y_{n+1} - \bar{y}\|^2 \right) \\
&\quad + \frac{1}{a_{21}} (\langle y_{n+1} - \bar{y}, Lx_n - L\bar{x} \rangle - \langle y_{n+2} - \bar{y}, Lx_{n+1} - L\bar{x} \rangle) \\
&= \frac{1}{2\tau c_{21}} \left( \|x_n - \bar{x} + c_{21}(p_n - x_n)\|^2 - \|x_n - \bar{x}\|^2 \right) \\
&\quad + \frac{1}{2\sigma a_{21}} \left( \|y_{n+1} - \bar{y} + a_{21}(q_{n+1} - y_{n+1})\|^2 - \|y_{n+1} - \bar{y}\|^2 \right) \\
&\quad + \frac{1}{a_{21}} \Big( \langle y_{n+1} - \bar{y}, Lx_n - L\bar{x} \rangle \\
&\qquad - \langle y_{n+1} - \bar{y} + a_{21}(q_{n+1} - y_{n+1}), Lx_n - L\bar{x} + c_{21}(Lp_n - Lx_n) \rangle \Big) \\
&= \frac{c_{21}}{2\tau} \|p_n - x_n\|^2 + \frac{1}{\tau} \langle x_n - \bar{x}, p_n - x_n \rangle + \frac{a_{21}}{2\sigma} \|q_{n+1} - y_{n+1}\|^2 \\
&\quad + \frac{1}{\sigma} \langle y_{n+1} - \bar{y}, q_{n+1} - y_{n+1} \rangle + \frac{c_{21}}{a_{21}} \langle \bar{y} - y_{n+1}, Lp_n - Lx_n \rangle \\
&\quad + \langle q_{n+1} - y_{n+1}, L\bar{x} - Lx_n \rangle + c_{21} \langle y_{n+1} - q_{n+1}, Lp_n - Lx_n \rangle. \qquad (3.9)
\end{aligned}$$

To estimate the above, we use the monotonicity of the operator $B^{-1}$ together with the inclusions $L\bar{x} \in B^{-1}\bar{y}$ from (3.2) and

$$\frac{y_{n+1} - q_{n+1}}{\sigma} + Lp_n + \frac{c_{21}}{a_{21}}(Lp_n - Lx_n) \in B^{-1}q_{n+1}, \qquad (3.10)$$

which is a reformulation of (3.8a) with $n$ replaced by $n+1$. This yields the inequality

$$\begin{aligned}
0 &\leq \left\langle \frac{y_{n+1} - q_{n+1}}{\sigma} + Lp_n + \frac{c_{21}}{a_{21}}(Lp_n - Lx_n) - L\bar{x}, q_{n+1} - \bar{y} \right\rangle \\
&= \frac{1}{\sigma} \langle y_{n+1} - q_{n+1}, q_{n+1} - \bar{y} \rangle + \langle Lp_n - L\bar{x}, q_{n+1} - \bar{y} \rangle \\
&\quad + \frac{c_{21}}{a_{21}} \langle Lp_n - Lx_n, q_{n+1} - \bar{y} \rangle \qquad (3.11)
\end{aligned}$$

Analogously, we can rewrite (3.8c) as

$$\frac{x_n - p_n}{\tau} - L^* y_{n+1} \in Ap_n. \qquad (3.12)$$

The monotonicity of $A$ together with the inclusion $-L^*\bar{y} \in A\bar{x}$ from (3.2) now yields

$$0 \le \left\langle \frac{x_n - p_n}{\tau} - L^*y_{n+1} + L^*\bar{y}, p_n - \bar{x} \right\rangle$$
$$= \frac{1}{\tau}\langle x_n - p_n, p_n - \bar{x}\rangle + \langle \bar{y} - y_{n+1}, Lp_n - L\bar{x}\rangle. \qquad (3.13)$$

Adding (3.11) and (3.13) yields

$$0 \le \frac{1}{\sigma}\langle y_{n+1} - q_{n+1}, q_{n+1} - \bar{y}\rangle + \langle Lp_n - L\bar{x}, q_{n+1} - y_{n+1}\rangle$$
$$+ \frac{c_{21}}{a_{21}}\langle Lp_n - Lx_n, q_{n+1} - \bar{y}\rangle + \frac{1}{\tau}\langle x_n - p_n, p_n - \bar{x}\rangle, \quad (3.14)$$

which, combined with (3.9), gives

$$Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y})$$
$$\le \frac{c_{21}}{2\tau}\|p_n - x_n\|^2 + \frac{1}{\tau}\langle x_n - \bar{x}, p_n - x_n\rangle + \frac{a_{21}}{2\sigma}\|q_{n+1} - y_{n+1}\|^2$$
$$+ \frac{1}{\sigma}\langle y_{n+1} - \bar{y}, q_{n+1} - y_{n+1}\rangle + \frac{c_{21}}{a_{21}}\langle \bar{y} - y_{n+1}, Lp_n - Lx_n\rangle$$
$$+ \langle q_{n+1} - y_{n+1}, L\bar{x} - Lx_n\rangle + c_{21}\langle y_{n+1} - q_{n+1}, Lp_n - Lx_n\rangle$$
$$+ \frac{1}{\sigma}\langle y_{n+1} - q_{n+1}, q_{n+1} - \bar{y}\rangle + \langle Lp_n - L\bar{x}, q_{n+1} - y_{n+1}\rangle$$
$$+ \frac{c_{21}}{a_{21}}\langle Lp_n - Lx_n, q_{n+1} - \bar{y}\rangle + \frac{1}{\tau}\langle x_n - p_n, p_n - \bar{x}\rangle$$
$$= \left(\frac{c_{21}}{2\tau} - \frac{1}{\tau}\right)\|p_n - x_n\|^2 + \left(\frac{a_{21}}{2\sigma} - \frac{1}{\sigma}\right)\|q_{n+1} - y_{n+1}\|^2$$
$$+ \left(\frac{c_{21}}{a_{21}} + 1 - c_{21}\right)\langle q_{n+1} - y_{n+1}, Lp_n - Lx_n\rangle$$
$$= -Q_2(p_n - x_n, q_{n+1} - y_{n+1}). \qquad (3.15)$$

This concludes the proof. $\qquad \square$

**Proposition 3.7.** *Let $Q_1$ and $Q_2$ be defined as in Lemma 3.5 and assume the conditions stated in Theorem 3.1 hold.*

*1. If $A$ is $\mu_1$-strongly monotone for some $\mu_1 > 0$, then*

$$Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}) + \mu_1\|p_n - \bar{x}\|^2$$
$$\le -Q_2(p_n - x_n, q_{n+1} - y_{n+1}).$$

*2. If $B^{-1}$ is $\mu_2$-strongly monotone for some $\mu_2 > 0$, then*

$$Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}) + \mu_2\|q_{n+1} - \bar{y}\|^2$$
$$\le -Q_2(p_n - x_n, q_{n+1} - y_{n+1}).$$

*Proof.* If $A$ is $\mu_1$-strongly monotone, we obtain from (3.12) and $-L^*\bar{y} \in A\bar{x}$ (3.2) the estimation

$$\mu_1 \|\bar{x} - p_n\|^2 \leq \left\langle \frac{x_n - p_n}{\tau} - L^* y_{n+1} + L^* \bar{y}, p_n - \bar{x} \right\rangle,$$

which is a sharpened version of (3.13). By modifying (3.14) and (3.15) accordingly, we get the assumption. The case of a strongly monotone $B^{-1}$ is analogously shown by improving (3.11). $\qquad\square$

Having stated and proved the necessary estimations, we are now ready to prove Theorem 3.1.

*Proof of Theorem 3.1.* Let $(\bar{x}, \bar{y})$ satisfy (3.2). By Proposition 3.6, we get the estimation

$$Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}) \leq -Q_2(p_n - x_n, q_{n+1} - y_{n+1}).$$

Considering Lemma 3.5, we see that the real sequence

$$(Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}))_{n \geq 1}$$

is monotonically nonincreasing and therefore has a limit for each primal-dual solution $(\bar{x}, \bar{y})$. Furthermore, for all $N \geq 1$,

$$Q_1(x_N - \bar{x}, y_{N+1} - \bar{y}) - Q_1(x_0 - \bar{x}, y_1 - \bar{y})$$
$$\leq -\sum_{n=0}^{N-1} Q_2(p_n - x_n, q_{n+1} - y_{n+1}).$$

By Lemma 3.5, we have $Q_1(x_N - \bar{x}, y_{N+1} - \bar{y}) \geq 0$ and

$$Q_1(x_0 - \bar{x}, y_1 - \bar{y}) \geq \sum_{n=0}^{N-1} Q_2(p_n - x_n, q_{n+1} - y_{n+1})$$
$$\geq \sum_{n=0}^{N-1} C_2 \|p_n - x_n\|^2$$

as well as

$$Q_1(x_0 - \bar{x}, y_1 - \bar{y}) \geq \sum_{n=0}^{N-1} D_2 \|q_{n+1} - y_{n+1}\|^2.$$

Since this holds for arbitrary $N \geq 1$, this proves parts (a) and (b) of the theorem. On the other hand, we have

$$Q_1(x_0 - \bar{x}, y_1 - \bar{y}) \geq Q_1(x_N - \bar{x}, y_{N+1} - \bar{y}) \geq C_1 \|x_N - \bar{x}\|^2$$

and

$$Q_1(x_0 - \bar{x}, y_1 - \bar{y}) \geq Q_1(x_N - \bar{x}, y_{N+1} - \bar{y}) \geq D_1 \|y_{N+1} - \bar{y}\|^2$$

for all $N \geq 1$, so the sequences $(x_n)_n$ and $(y_n)$ are bounded in $\mathcal{X}$ and $\mathcal{Y}$, respectively. Let $(n_k)_k$ be a subsequence with $x_{n_k} \rightharpoonup x_\infty \in \mathcal{X}$ and $y_{n_k+1} \rightharpoonup y_\infty \in \mathcal{Y}$. By (3.12) and (3.10), we obtain

$$\frac{x_{n_k} - p_{n_k}}{\tau} - L^* y_{n_k+1} \in A p_{n_k},$$

$$\frac{y_{n_k+1} - q_{n_k+1}}{\sigma} + L p_{n_k} + \frac{c_{21}}{a_{21}}(L p_{n_k} - L x_{n_k}) \in B^{-1} q_{n_k+1}.$$

Now apply [7, Proposition 2.4] with

$$a_k = p_{n_k},$$

$$a_k^* = \frac{x_{n_k} - p_{n_k}}{\tau} - L^* y_{n_k+1},$$

$$b_k = \frac{y_{n_k+1} - q_{n_k+1}}{\sigma} + L p_{n_k} + \frac{c_{21}}{a_{21}}(L p_{n_k} - L x_{n_k}),$$

$$b_k^* = q_{n_k+1}$$

and observe that

$$a_k = x_{n_k} + (p_{n_k} - x_{n_k}) \rightharpoonup x_\infty,$$

$$b_k^* = y_{n_k+1} + (q_{n_k+1} - y_{n_k+1}) \rightharpoonup y_\infty,$$

$$a_k^* + L^* b_k^* = \frac{x_{n_k} - p_{n_k}}{\tau} + L^*(q_{n_k+1} - y_{n_k+1}) \to 0,$$

$$L a_k - b_k = -\frac{y_{n_k+1} - q_{n_k+1}}{\sigma} - \frac{c_{21}}{a_{21}}(L p_{n_k} - L x_{n_k}) \to 0$$

because parts (a) and (b) imply that $x_{n_k} - p_{n_k} \to 0$ and $y_{n_k+1} - q_{n_k+1} \to 0$ as $k \to +\infty$. This gives $L x_\infty \in B^{-1} y_\infty$ and $-L^* y_\infty \in A x_\infty$, i.e., $(x_\infty, y_\infty)$ satisfies (3.2). Since the choice of the weakly convergent subsequence was arbitrary, each weak sequential cluster point satisfies (3.2). Claim (c) now follows from [10, Lemma 2.47] applied to the norm $\sqrt{Q_1(\cdot)}$ on the product space $\mathcal{X} \times \mathcal{Y}$ and to the solution set of (3.2).

Now assume that $A$ is $\mu_1$-strongly monotone for some $\mu_1 > 0$. By Proposition 3.7, we get the estimation

$$Q_1(x_{n+1} - \bar{x}, y_{n+2} - \bar{y}) - Q_1(x_n - \bar{x}, y_{n+1} - \bar{y}) + \mu_1 \|p_n - \bar{x}\|^2$$
$$\leq -Q_2(p_n - x_n, q_{n+1} - y_{n+1})$$

for all $n \geq 0$. Choose $N \geq 1$ and sum up this inequality for $n = 0, \ldots, N - 1$ to obtain

$$Q_1(x_N - \bar{x}, y_{N+1} - \bar{y}) - Q_1(x_0 - \bar{x}, y_1 - \bar{y}) + \mu_1 \sum_{n=0}^{N-1} \|p_n - \bar{x}\|^2$$
$$\leq -\sum_{n=0}^{N-1} Q_2(p_n - x_n, q_{n+1} - y_{n+1})$$

Since the terms $Q_1(x_N - \bar{x}, y_{N+1} - \bar{y})$ and $\sum_{n=0}^{N-1} Q_2(p_n - x_n, q_{n+1} - y_{n+1})$ are nonnegative by Lemma 3.5, we obtain

$$\mu_1 \sum_{n=0}^{N-1} \|p_n - \bar{x}\|^2 \leq Q_1(x_0 - \bar{x}, y_1 - \bar{y}).$$

14

Analogously, one gets

$$\mu_2 \sum_{n=0}^{N-1} \|q_{n+1} - \bar{y}\|^2 \le Q_1(x_0 - \bar{x}, y_1 - \bar{y}),$$

if $B^{-1}$ is $\mu_2$-strongly monotone, and since $N$ is arbitrary, both sums

$$\sum_{n=0}^{\infty} \|p_n - \bar{x}\|^2 \qquad \text{and} \sum_{n=0}^{\infty} \|q_{n+1} - \bar{y}\|^2$$

are finite in the respective cases. The uniqueness of the point $\bar{x}$ under the assumption of strong monotonicity of $A$ holds by the fact that we have shown $p_n \to \bar{x}$ for *any* solution $(\bar{x}, \bar{y})$ of (3.2). An analogous argument for $\bar{y}$ concludes the proof of Claim (d). □

*Remark* 3.8. We were not able to show the weak convergence of PDHG (2.9) for $\theta \ne 1$ with this proof method. Indeed, by a straightforward calculation it can be shown that from Fejér monotonicity with respect to any quadratic form of the sequence $(x_n, y_{n+1})_n$ the conditions (3.6) can be derived, which implies $\theta = 1$.

## 3.3 Application to convex optimization

In this section, we specialize the scheme (3.8) to the case where the monotone operators $A$ and $B$ are subdifferentials $\partial F$ and $\partial G$ of proper, convex and lower semicontinuous functions $F : \mathcal{X} \to \overline{\mathbb{R}}$ and $G : \mathcal{Y} \to \overline{\mathbb{R}}$, resectively. Algorithm 3.2 then reads as follows:

**Algorithm 3.9.** Choose parameters $\sigma > 0$, $\tau > 0$ and $a_{21} \in \mathbb{R}$, $c_{21} \in \mathbb{R}$ and starting points $x_0 \in \mathcal{X}$, $x_1 \in \mathcal{X}$, $p_0 \in \mathcal{X}$, $y_1 \in \mathcal{Y}$. For all $n = 1, 2, \ldots$, calculate

$$q_n = \text{Prox}_{G^*}^{\sigma} \left( y_n + \sigma L \left( p_{n-1} + \frac{c_{21}}{a_{21}} (p_{n-1} - x_{n-1}) \right) \right), \qquad (3.16\text{a})$$

$$y_{n+1} = y_n + a_{21}(q_n - y_n), \qquad (3.16\text{b})$$

$$p_n = \text{Prox}_F^{\tau} (x_n - \tau L^* y_{n+1}), \qquad (3.16\text{c})$$

$$x_{n+1} = x_n + c_{21}(p_n - x_n). \qquad (3.16\text{d})$$

Then, $x_n \rightharpoonup \bar{x}$, $p_n \rightharpoonup \bar{x}$, $y_n \rightharpoonup \bar{y}$, and $q_n \rightharpoonup \bar{y}$, where $(\bar{x}, \bar{y})$ is a solution of (2.7), provided that $0 < a_{21} < 2$, $0 < c_{21} < 2$, and (3.7) are satisfied.

In this case, it is possible to get estimations for the Lagrangian, which is defined in (2.6).

**Theorem 3.10.** *Given the assumptions in Theorem 3.1, let $F : \mathcal{X} \to \overline{\mathbb{R}}$ and $G : \mathcal{Y} \to \overline{\mathbb{R}}$ be two proper, convex and lower semicontinuous functions. Let $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ be arbitrary. Then the sequence $(q_n, y_n, p_n x_n)$ generated by (3.16) satisfies*

$$\min_{n=0,\ldots,N-1} \left( \mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1}) \right) \le \frac{1}{N} Q_1(x_0 - x, y_1 - y),$$

$$\mathcal{L}\left( \frac{1}{N} \sum_{n=0}^{N-1} p_n; y \right) - \mathcal{L}\left( x; \frac{1}{N} \sum_{n=0}^{N-1} q_{n+1} \right) \le \frac{1}{N} Q_1(x_0 - x, y_1 - y).$$

This theorem is proved using the following proposition, which bounds the Lagrangian in terms of the quadratic forms defined in Lemma 3.5.

**Proposition 3.11.** *Given the assumptions in Theorem 3.1, let $F: \mathcal{X} \to \overline{\mathbb{R}}$ and $G: \mathcal{Y} \to \overline{\mathbb{R}}$ be two proper, convex and lower semicontinuous functions. Let $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ be arbitrary. Then the sequence $(q_n, y_n, p_n, x_n)$ generated by (3.16) satisfies*

$$\mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1}) \leq Q_1(x_n - x, y_{n+1} - y) - Q_1(x_{n+1} - x, y_{n+2} - y)$$
$$- Q_2(p_n - x_n, q_{n+1} - y_{n+1})$$

*for all $n \geq 1$, $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.*

*Proof.* Since $B^{-1} = \partial G^*$ and $A = \partial F$, the inclusions (3.10) and (3.12) provide certain subgradients, which imply the inequalities

$$G^*(y) \geq G^*(q_{n+1}) + \frac{1}{\sigma}\langle y_{n+1} - q_{n+1}, y - q_{n+1}\rangle + \langle Lp_n, y - q_{n+1}\rangle$$
$$+ \frac{c_{21}}{a_{21}}\langle Lp_n - Lx_n, y - q_{n+1}\rangle,$$
$$F(x) \geq F(p_n) + \frac{1}{\tau}\langle x_n - p_n, x - p_n\rangle - \langle L^*y_{n+1}, x - p_n\rangle.$$

Therefore, we have

$$\mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1})$$
$$= \langle Lp_n, y\rangle + F(p_n) - G^*(y) - \langle Lx, q_{n+1}\rangle - F(x) + G^*(q_{n+1})$$
$$\leq \frac{1}{\tau}\langle x_n - p_n, p_n - x\rangle + \langle Lp_n - Lx, q_{n+1} - y_{n+1}\rangle$$
$$+ \frac{1}{\sigma}\langle y_{n+1} - q_{n+1}, q_{n+1} - y\rangle + \frac{c_{21}}{a_{21}}\langle Lp_n - Lx_n, q_{n+1} - y\rangle.$$

The right-hand side is now (except for the replacement of $\bar{x}$ and $\bar{y}$ by $x$ and $y$, respectively) equal to the one in (3.14), and one easily checks by an analogous calculation, that it equals the expression in the assertion. $\square$

*Proof of Theorem 3.10.* By summing the inequality in Proposition 3.11 for $n = 0, \ldots, N-1$ and dividing by $N$ for some $N \geq 1$, we get

$$\frac{1}{N}\sum_{n=0}^{N-1}(\mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1})) \leq \frac{1}{N}Q_1(x_0 - x, y_1 - y)$$

for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, where we dropped nonpositive terms on the right-hand side.

We have two possibilities to further estimate the left-hand side: First, we notice that it is the arithmetic mean of numbers, which is always greater than the minimum, i.e.,

$$\frac{1}{N}\sum_{n=0}^{N-1}(\mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1})) \geq \min_{n=0,\ldots,N-1}(\mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1})).$$

On the other hand, the Lagrangian is convex in its first and concave in its second component, so

$$\frac{1}{N}\sum_{n=0}^{N-1}(\mathcal{L}(p_n; y) - \mathcal{L}(x; q_{n+1})) \geq \mathcal{L}\left(\frac{1}{N}\sum_{n=0}^{N-1}p_n; y\right) - \mathcal{L}\left(x; \frac{1}{N}\sum_{n=0}^{N-1}q_{n+1}\right). \square$$

# 4 Learning an optimization solver

Most optimization problems are solved using iterative methods, akin to the ones presented in Sections 2 and 3. However, the number of iterations it takes in order for the algorithm to converge is in general hard to predict, which creates problems in time-critical applications. In these situations one could instead consider only doing a predefined fixed number $n$ of iterations. A natural question that arises in response to this is: *what parameter values in the optimization solver give the best improvement of the objective function in $n$ iterations?* This question leads to a meta-optimization over optimization solvers. Moreover, in general we are not only interested in optimizing one single cost function, but rather a (potentially infinite) family $\{F_\theta\}_{\theta \in \Theta}$ of cost functions, each with a minimizer $\bar{x}_\theta$. Hence, to make the question precise one needs to specify which family of optimization solvers one is considering, which is the family of cost functions of interested, and what is meant with "best improvement".

One such question was raised in [23], where the authors consider the worst-case performance $\sup_{\theta \in \Theta} \big[ F_\theta(x_n) - F_\theta(\bar{x}_\theta) \big]$ of gradient-based algorithms over the set of continuously differentiable functions with Lipschitz-continuous gradients, and with a uniform upper bound on the Lipschitz constants. Subsequent work along the same lines can be found in [31, 48].

The idea of optimizing over optimization solvers has also been considered from a machine learning perspective. This has for example been done using *reinforcement learning* [35], and using *unsupervised learning* [8, 27]. In the latter category, one looks for algorithm parameters which minimize the expected value of the difference in objective function value,

$$\mathbb{E}_\theta \big[ F_\theta(x_n) - F_\theta(\bar{x}_\theta) \big] = \mathbb{E}_\theta \big[ F_\theta(x_n) \big] - \mathbb{E}_\theta \big[ F_\theta(\bar{x}_\theta) \big] \tag{4.1}$$

where $\Theta$ is endowed with a probability measure and $x_n$ is the output of the algorithm after $n$ iterations. However, optimizing (4.1) with respect to the parameters of the method is independent of the optimal points $\{\bar{x}_\theta\}_{\theta \in \Theta}$, thus, this translates into unsupervised learning, i.e., the cost function $\mathbb{E}_\theta \big[ F_\theta(x_n) \big]$ does not depend on $\bar{x}_\theta$. In this setting, [8] restricts attention to an architecture that operates individually on each coordinate of $x$. This is done in order to limit the number of parameters in the algorithm, which otherwise would grow exponentially with the dimension of $x$. To overcome this, we use an approach similar to [27], where the network architecture is inspired by modern first-order optimization solvers for nonsmooth problems, as presented in Sections 2 and 3. Similar ideas have also recently been explored for supervised learning in inverse problems in [4–6, 28, 37, 45, 49].

## 4.1 Unrolled gradient descent as a neural network

Before we define the architecture considered in this work, we first present an illustrative example. To this end, consider the optimization problem

$$\min_x \quad F(x).$$

We assume that $F$ is smooth, which means that the problem can be solved using a standard gradient descent algorithm, i.e., by performing the updates

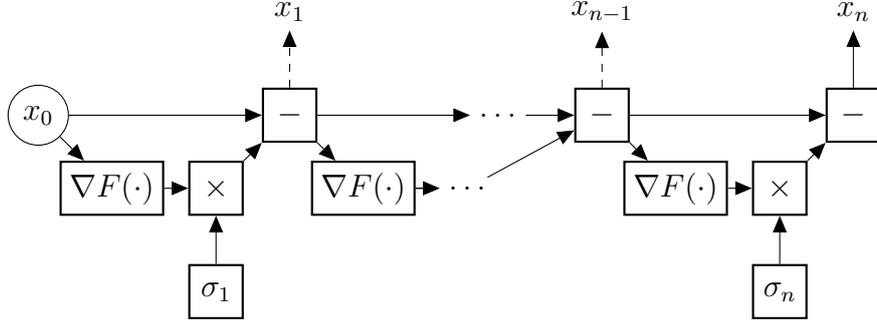$$x_k = x_{k-1} - \sigma_k \nabla F(x_k).$$

17

Figure 1: Gradient descent.

The gradient descent algorithm contains a set of parameters that need to be selected, namely the step length for each iteration, $\sigma_k$. This is normally done via the Goldstein rule or backtracking line search (Armijo rule) [13], which under suitable conditions ensures convergence to the optimal point $\bar{x}$.

However, if we only run the algorithm for a fixed number $n$ of steps, the gradient descent algorithm can be seen as a *feedforward neural network*, as shown in Figure 1. Each layer in the network performs the computation $x_{k-1} - \sigma_k \nabla F(x_{k-1})$ and the parameters of the network are $[\sigma_1, \ldots, \sigma_n]$. Moreover, if the step length is fixed to be the same in all iterations, i.e., $\sigma_1 = \ldots = \sigma_n = \sigma$ for some $\sigma$, the gradient descent algorithm can in fact be interpreted as a *recurrent neural network*. In both cases, for a given family $\{F_\theta\}_{\theta \in \Theta}$ of cost functions the network parameter(s) can be trained (optimized) by minimizing $\mathbb{E}_\theta \left[ F_\theta(x_n) \right]$, where $x_n$ is the output of the network in Figure 1. For simple cases this can be done analytically.

*Example* 4.1. Consider the family $(F_b)_b$ of functions $F_b : \mathbb{R}^n \to \mathbb{R}$ given by $F_b(x) = \frac{1}{2} x^\top A x - b^\top x$, where $A \in \mathbb{R}^{n \times n}$ is a (fixed) symmetric and positive definite matrix. The minimum of $F_b$ is given by $\bar{x}_b = A^{-1} b$. Denote by $\Lambda_\sigma$ the result of taking a gradient step of length $\sigma > 0$, i.e.,

$$\Lambda_\sigma(x) = x - \sigma \nabla F_b(x) = x - \sigma(Ax - b), \qquad x \in \mathbb{R}^n.$$

Let $x_0 \in \mathbb{R}^n$ be an arbitrary starting point of the iteration. This gives

$$
\begin{aligned}
F_b(\Lambda_\sigma(x_0)) &= F_b(x_0 - \sigma(Ax_0 - b)) \\
&= \frac{1}{2}(x_0 - \sigma(Ax_0 - b))^\top A(x_0 - \sigma(Ax_0 - b)) - b^\top(x_0 - \sigma(Ax_0 - b)) \\
&= \frac{\sigma^2}{2}(Ax_0 - b)^\top A(Ax_0 - b) - \sigma\|Ax_0 - b\|^2 + F_b(x_0).
\end{aligned}
$$

Let $\mathsf{b}$ be a random variable distributed according to $\mathsf{b} \sim \mathcal{P}$ for some probability distribution $\mathcal{P}$ with finite first and second moments. Finding a $\sigma$ that minimizes the expectation

$$
\begin{aligned}
\mathbb{E}_{\mathsf{b} \sim \mathcal{P}} \left[ F_{\mathsf{b}}(\Lambda_\sigma(x_0)) \right] = \frac{\sigma^2}{2} \mathbb{E}_{\mathsf{b} \sim \mathcal{P}} \left[ (Ax_0 - \mathsf{b})^\top A(Ax_0 - \mathsf{b}) \right] \\
- \sigma \, \mathbb{E}_{\mathsf{b} \sim \mathcal{P}} \left[ \|Ax_0 - \mathsf{b}\|^2 \right] + \mathbb{E}_{\mathsf{b} \sim \mathcal{P}} \left[ F_{\mathsf{b}}(x_0) \right],
\end{aligned}
$$

is a quadratic problem in one variable, and the optimal value of $\sigma$ is thus

$$\sigma = \frac{\mathbb{E}_{\mathsf{b}\sim\mathcal{P}}\left[\|Ax_0 - \mathsf{b}\|^2\right]}{\mathbb{E}_{\mathsf{b}\sim\mathcal{P}}\left[(Ax_0 - \mathsf{b})^\top A(Ax_0 - \mathsf{b})\right]}$$

$$= \frac{\|Ax_0\|^2 - 2(Ax_0)^\top \mathbb{E}_{\mathsf{b}\sim\mathcal{P}}\left[\mathsf{b}\right] + \mathbb{E}_{\mathsf{b}\sim\mathcal{P}}\left[\|\mathsf{b}\|^2\right]}{x_0^\top A^3 x_0 - 2(A^2 x_0)^\top \mathbb{E}_{\mathsf{b}\sim\mathcal{P}}\left[\mathsf{b}\right] + \mathbb{E}_{\mathsf{b}\sim\mathcal{P}}\left[\mathsf{b}^\top A\mathsf{b}\right]}.$$

In some particular cases this expression can be simplified. For example if $A = I$, then $\sigma = 1$ as expected. Or if $x_0 = 0$, then $\sigma = \mathbb{E}_{\mathsf{b}\sim\mathcal{P}}[\|\mathsf{b}\|^2]/\mathbb{E}_{\mathsf{b}\sim\mathcal{P}}[\mathsf{b}^\top A\mathsf{b}]$.

## 4.2 Parametrizing a family of optimization algorithms

Similarly to the considerations in Section 4.1, for a fixed number of iterations one can consider the optimization algorithms (2.9), (2.10) and (3.16) as neural networks, where the variables we want to train are the parameters of the optimization methods. Optimizing these parameters with respect to the constraints corresponding to each algorithm is effectively trying to find optimal parameters for the corresponding algorithm for a given family of cost functions. However, if one only intends to do a finite number of iterations one could also remove this constraint, and thereby enlarge the space of schemes one is optimizing over.

As noted in Section 3, all of the above mentioned optimization algorithms can be written on the form (3.1). That means that optimizing over the parameters in (3.1) can be seen as optimizing over a space of schemes that includes all three algorithms. Now, introducing the intermediate states $w_n = a_{11}q_n + a_{12}y_n$ and $v_{n+1} = c_{11}p_n + c_{12}x_n$, and the $2 \times 2$ matrices $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}, \boldsymbol{D}$, the scheme (3.1) can be written as

$$
\begin{aligned}
\begin{bmatrix} w_n \\ y_{n+1} \end{bmatrix} &= (\boldsymbol{A} \otimes \mathrm{Id}) \;\; \mathrm{diag}(\mathrm{Prox}_{G^*}^\sigma, \mathrm{Id}) \;\; (\boldsymbol{B} \otimes \mathrm{Id}) \begin{bmatrix} Lv_n \\ y_n \end{bmatrix} \\
\begin{bmatrix} v_{n+1} \\ x_{n+1} \end{bmatrix} &= (\boldsymbol{C} \otimes \mathrm{Id}) \;\; \mathrm{diag}(\mathrm{Prox}_F^\tau, \mathrm{Id}) \;\; (\boldsymbol{D} \otimes \mathrm{Id}) \begin{bmatrix} L^* w_n \\ x_n \end{bmatrix},
\end{aligned}
\tag{4.2}
$$

where the parameters of the scheme are the elements of the matrices. Here, by $\otimes$ we denote the Kronecker product, and by $\mathrm{diag}\,(A, B, \dots)$ we denote the diagonal operator with the operators $A, B, \dots$ on the diagonal. Connecting this with the previous optimization algorithms, the PDHG algorithm (2.9) is obtained by setting

$$\boldsymbol{A} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix}, \quad \boldsymbol{B} = \begin{bmatrix} \sigma & 1 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{C} = \begin{bmatrix} 1+\theta & -\theta \\ 1 & 0 \end{bmatrix}, \quad \boldsymbol{D} = \begin{bmatrix} -\tau & 1 \\ 0 & 1 \end{bmatrix},$$

the primal-dual Douglas-Rachford algorithm (2.10) by taking

$$\boldsymbol{A} = \begin{bmatrix} \lambda_n & 1-\lambda_n \\ \lambda_n & 1-\lambda_n \end{bmatrix}, \quad \boldsymbol{B} = \begin{bmatrix} \sigma & 1 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{C} = \begin{bmatrix} 2 & -1 \\ \lambda_n & 1-\lambda_n \end{bmatrix}, \quad \boldsymbol{D} = \begin{bmatrix} -\tau & 1 \\ 0 & 1 \end{bmatrix},$$

and the proposed algorithm from Section 3 by setting

$$\boldsymbol{A} = \begin{bmatrix} a_{21} & 1-a_{21} \\ a_{21} & 1-a_{21} \end{bmatrix}, \quad \boldsymbol{B} = \begin{bmatrix} \sigma & 1 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{C} = \begin{bmatrix} 1+\frac{c_{21}}{a_{21}} & -\frac{c_{21}}{a_{21}} \\ c_{21} & 1-c_{21} \end{bmatrix}, \quad \boldsymbol{D} = \begin{bmatrix} -\tau & 1 \\ 0 & 1 \end{bmatrix}.$$

Considering (4.2) as a neural network, the structure can easily be extended in order to incorporate more memory in the network. In this work we assume that the computationally expensive part of the algorithm is the evaluation of the operator $L$ and its adjoint, which is typically the case in inverse problems in imaging, e.g., in three-dimensional CT [40, 41]. Therefore, the extension presented here thus keeps one evaluation $L$ and one evaluation of $L^*$ in each iteration.

To this end, let $N$ be the number of primal variables $x^1, \ldots, x^N \in \mathcal{X}$ and $M$ be the number of dual variables $y^1, \ldots, y^M \in \mathcal{Y}$. Introducing the four sequences of matrices $\boldsymbol{A}_n, \boldsymbol{B}_n \in \mathbb{R}^{M \times M}$ and $\boldsymbol{C}_n, \boldsymbol{D}_n \in \mathbb{R}^{N \times N}$, the iterations in (4.2) can be extended to yield the following algorithm.

**Algorithm 4.2.** Choose parameters $\boldsymbol{A}_n, \boldsymbol{B}_n \in \mathbb{R}^{M \times M}$ and $\boldsymbol{C}_n, \boldsymbol{D}_n \in \mathbb{R}^{N \times N}$, stepsizes $\sigma, \tau > 0$, and starting points $x_0^1, \ldots, x_0^N \in \mathcal{X}$, $y_0^2, \ldots, y_0^M \in \mathcal{Y}$. For all $n = 1, 2, \ldots$, calculate

$$
\begin{bmatrix} y_{n+1}^1 \\ y_{n+1}^2 \\ \vdots \\ y_{n+1}^M \end{bmatrix} = (\boldsymbol{A}_n \otimes \mathrm{Id}) \ \ \mathrm{diag}(\mathrm{Prox}_{G^*}^\sigma, \mathrm{Id}^{M-1}) \ (\boldsymbol{B}_n \otimes \mathrm{Id}) \begin{bmatrix} Lx_n^1 \\ y_n^2 \\ \vdots \\ y_n^M \end{bmatrix},
$$

$$
\begin{bmatrix} x_{n+1}^1 \\ x_{n+1}^2 \\ \vdots \\ x_{n+1}^N \end{bmatrix} = (\boldsymbol{C}_n \otimes \mathrm{Id}) \ \ \mathrm{diag}(\mathrm{Prox}_F^\tau, \mathrm{Id}^{N-1}) \ (\boldsymbol{D}_n \otimes \mathrm{Id}) \begin{bmatrix} L^* y_{n+1}^1 \\ x_n^2 \\ \vdots \\ x_n^N \end{bmatrix}.
$$

*Remark* 4.3. For the more general formulation of (1.1), more specialized network architectures than the one resulting from the choice (2.4) are possible, which handle the dual spaces separately instead of using the same stepsize $\sigma$ and matrices $\boldsymbol{A}_n$ and $\boldsymbol{B}_n$ for all of them. An alternative network in the spirit of, e.g., [14, Theorem 2], to solve (1.1) reads as follows.

**Algorithm 4.4.** Choose parameters $\boldsymbol{A}_{n,i}, \boldsymbol{B}_{n,i} \in \mathbb{R}^{M \times M}$, for $i = 1, \ldots, m$, and $\boldsymbol{C}_n, \boldsymbol{D}_n \in \mathbb{R}^{N \times N}$, stepsizes $\sigma_1, \ldots, \sigma_m, \tau > 0$, and starting points $x_0^1, \ldots, x_0^N \in \mathcal{X}$, $y_{0,i}^2, \ldots, y_{0,i}^M \in \mathcal{Y}_i$, $i = 1, \ldots, m$. For all $n = 1, 2, \ldots$, calculate

$$
\begin{bmatrix} y_{n+1,i}^1 \\ y_{n+1,i}^2 \\ \vdots \\ y_{n+1,i}^M \end{bmatrix} = (\boldsymbol{A}_{n,i} \otimes \mathrm{Id}) \ \ \mathrm{diag}(\mathrm{Prox}_{G_i^*}^{\sigma_i}, \mathrm{Id}^{M-1}) \ (\boldsymbol{B}_{n,i} \otimes \mathrm{Id}) \begin{bmatrix} L_i x_n^1 \\ y_{n,i}^2 \\ \vdots \\ y_{n,i}^M \end{bmatrix},
$$
$$
i = 1, \ldots, m,
$$

$$
\begin{bmatrix} x_{n+1}^1 \\ x_{n+1}^2 \\ \vdots \\ x_{n+1}^N \end{bmatrix} = (\boldsymbol{C}_n \otimes \mathrm{Id}) \ \ \mathrm{diag}(\mathrm{Prox}_F^\tau, \mathrm{Id}^{N-1}) \ (\boldsymbol{D}_n \otimes \mathrm{Id}) \begin{bmatrix} \sum_{i=1}^m L_i^* y_{n+1,i}^1 \\ x_n^2 \\ \vdots \\ x_n^N \end{bmatrix}.
$$
$$(4.3)$$

### 4.2.1 Extension to forward-backward-forward methods

Some methods in the literature, so called forward-backward-forward methods, include an extra evaluation of the operator and its adjoint per iteration, see, e.g.,

[15, 17, 21]. However, since the evaluation of the linear operator is assumed to be the expensive part in our setting we consider this as two iterations. Thus, if we start with the $x$-iterate and allow for two iterations in our framework to complete one iteration in such a framework, our proposed algorithm contains, e.g., [17, Equation (3.1)]. Letting $\cdot$ denote an element that can take any value, one such set of matrices is given by

$$\boldsymbol{A}_{2n} = \begin{bmatrix} 0 & 0 & 1 \\ \cdot & \cdot & \cdot \\ 1 & 0 & 0 \end{bmatrix}, \qquad \boldsymbol{B}_{2n} = \begin{bmatrix} \gamma_n & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

$$\boldsymbol{C}_{2n} = \begin{bmatrix} 1 & 0 & -1 \\ \cdot & \cdot & \cdot \\ 1 & 1 & 0 \end{bmatrix}, \qquad \boldsymbol{D}_{2n} = \begin{bmatrix} -\gamma_n & 0 & 1 \\ \gamma_n & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

for the even iterations and

$$\boldsymbol{A}_{2n+1} = \begin{bmatrix} 0 & 1 & 0 \\ \cdot & \cdot & \cdot \\ 0 & 0 & 1 \end{bmatrix}, \qquad \boldsymbol{B}_{2n+1} = \begin{bmatrix} \cdot & \cdot & \cdot \\ 0 & 0 & 1 \\ \gamma_n & 0 & 1 \end{bmatrix},$$

$$\boldsymbol{C}_{2n+1} = \begin{bmatrix} 0 & 0 & 1 \\ \cdot & \cdot & \cdot \\ 0 & 0 & 1 \end{bmatrix}, \qquad \boldsymbol{D}_{2n+1} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ -\gamma_n & 0 & 1 \end{bmatrix},$$

for the odd iterations.

*Remark* 4.5. Other forward-backward-forward methods have been proposed in the literature, some of which are general enough to include the PDHG as a special case [29], or both the PDHG and the Douglas-Rachford algorithm as special cases [34]. However, these methods include a step-length computation in their updates. This computation involves evaluating the norm of current iterates, which is not possible to achieve by only doing the linear operations we propose. Of course, allowing the matrix elements to be nonlinear functions of the states would allow us to incorporate also these methods, however, that is beyond the scope of this article.

# 5 Application to inverse problems and numerical experiments

As we briefly outline next, optimization problems of the type in (1.1) arise when solving ill-posed inverse problems by means of variational regularization.

The goal in an inverse problem is to recover parameters characterizing a system under investigation from indirect observations. This can be formalized as the task of estimating (reconstructing) model parameters, henceforth called signal, $x_{\text{true}} \in \mathcal{X}$ from indirect observations (data) $b \in \mathcal{Y}$ where

$$b = \text{T}\left(x_{\text{true}}\right) + \delta b. \tag{5.1}$$

In the above, $\mathcal{X}$ and $\mathcal{Y}$ are typically Hilbert or Banach spaces, and $\text{T} \colon \mathcal{X} \to \mathcal{Y}$ (forward operator) models how a given signal gives rise to data in absence of noise. Furthermore, $\delta b \in \mathcal{Y}$ is a single sample of a $\mathcal{Y}$-valued random element that represents the noise component of data.

A natural approach for solving (5.1) is to minimize a function $D\colon \mathcal{X} \to \mathbb{R}$ (data discrepancy functional) that quantifies the miss-fit in data space. Since this function needs to incorporate the aforementioned forward operator T and the data $b$, it is often of the form

$$D(x) := \ell\left(\mathrm{T}(x), b\right) \quad \text{for some } \ell\colon \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}.$$

If $\ell$ is the negative data log-likelihood, then minimizing $x \mapsto D(x)$ corresponds to finding a maximum likelihood solution to (5.1).

However, finding a minimizer to $D$ is an ill-posed problem, meaning that a solution (if it exists) is discontinuous with respect to the data $b$. Variational regularization addresses this issue by introducing an additional function $R\colon \mathcal{X} \to \overline{\mathbb{R}}$ (regularization functional) that encodes a priori information about $x_{\text{true}}$ and penalizes undesirable solutions [26]. This results in an optimization problem

$$\min_{x \in \mathcal{X}} \left[ \lambda D(x) + R(x) \right], \tag{5.2}$$

which from a statistical perspective can be interpreted as trying to find a maximum a posteriori estimate [30]. A common choice of regularization functional, especially for inverse problems in imaging, is the total variation (TV) regularization $R(x) := \|\nabla x\|_1$, but several more advanced regularizers have also been suggested in the literature, typically exploiting some kind of sparsity using an $L_1$-like norm [19].

In this section, we consider an inverse problem in computerized tomography. To this end, let T be the *Radon transform* and consider TV regularization. This means that we are interested in minimizing

$$H_b(x) = \|\mathrm{T}(x) - b\|_2^2 + \lambda \|\nabla x\|_1, \tag{5.3}$$

i.e., a family of objective functions that is parametrized by the data $b$. This means that we can apply the ideas from Section 4 on learning an optimization solver.

## 5.1 Implementation and specifications of the training

We train and evaluate several of the algorithms described in this article on a clinically realistic data set, namely simulated data from human abdomen CT scans as provided by Mayo Clinic for the AAPM Low Dose CT Grand Challenge [38]. Examples of two-dimensional phantoms from this data set are given in Figure 2. Throughout all examples, the size of the image $x$ is $512 \times 512$ pixels, and the regularization parameter $\lambda > 0$ is fixed. The Radon transform T used in this example is sampled according to a fan-beam geometry [40] and the data is generated by applying T to the phantoms and then adding 5% white Gaussian noise. Examples of such data (sinograms) are also shown in Figure 2.

Problem (5.3) is obtained from (1.1) by setting $F(x) := 0$ for all $x$,

$$L_1 x := \mathrm{T}(x), \qquad\qquad L_2 x := \nabla x,$$
$$G_1(y_1) := \|y_1 - b\|_2^2, \qquad\qquad G_2(y_2) := \|y_2\|_1,$$

and all the proximal operators are implemented in ODL [3]. If not stated otherwise, we use (2.4) to reduce (1.1) to (2.3).
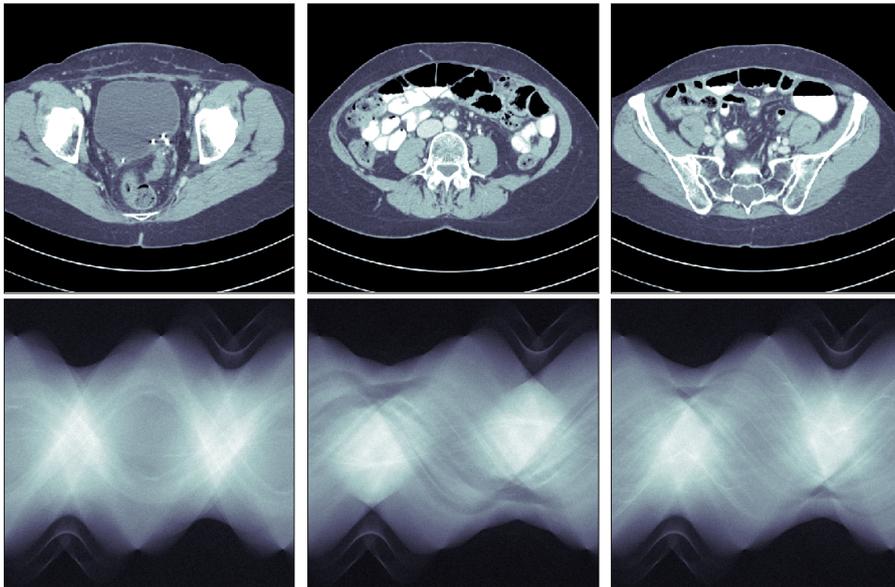
Figure 2: The top row shows three examples of phantoms used for generating data. These phantoms take values between $[0.0, 3.25]$, but all images are shown using a window set to $[0.8, 1.2]$ in order to enhance contrast of clinically more relevant details. The lower row shows corresponding simulated, noisy sinograms.

For each algorithm, the number of unrolled iterations, corresponding to the depth of the network, was set to $n_{\max} = 10$, and all evaluations have been done with this depth. However, in order to heuristically induce better stability of the general schemes, we have trained using a stochastic depth as follows: In each step of the training, the depth of the network has been set to the outcome of the heavy-tailed random variable $n_{\max} = \min\left[\texttt{round}(8 + \mathsf{Z}), 100\right]$, where $\mathsf{Z}$ is the exponential of a Gaussian random variable with standard deviation $1.25$ and mean value $\log(2) - 1.25^2/2$, so that $\mathbb{E}[\mathsf{Z}] = 2$. The limitation to 100 iterations is due to limits in computational resources.

In order to improve stability and generalization properties of the trained networks, we have normalized the operators before training, i.e., rescaled them so that $\|\mathsf{T}\|_2 = \|\nabla\|_2 = 1$. For the same reasons, we have used the zero vector as initial guess for all networks. Training has been done using the *Adam* solver [32], with standard parameter values except for $\beta_2 = 0.99$. Moreover, we have used gradient clipping to limit the norm of the gradient of the training cost function (4.1) to be less than or equal to one [44]. As step length (learning rate) we have used a cosine annealing scheme [36], i.e., a step length which in step $t$ takes the value

$$\eta_t = \frac{\eta_0}{2}\left(1 + \cos\left(\pi \frac{t}{t_{\max}}\right)\right),$$

where the initial step length $\eta_0$ has been set to $10^{-3}$. We have trained for $t_{\max} = 100\,000$ steps and have used 9 out of 10 phantoms from the AAPM Low Dose CT Grand Challenge for training and one for evaluation.

All algorithms have been implemented using ODL [3], the GPU accelerated

version of ASTRA [1, 43], and Tensorflow [2]. The source code to replicate the experiments is available online, where the weights of the trained networks are also explicitly given.[‡] We have used this setup to train the following methods.

**PDHG method.** This corresponds to optimal selection of the parameters $\theta$, $\tau$, and $\sigma$ for the PDHG method (2.9) on the family of cost functions (5.3). In order to achieve this, we need to enforce the constraints $\theta \in [0, 1]$ and $\sigma\tau\|L\|^2 < 1$. This has been done implicitly by a change of variables, namely by

$$\theta = \frac{e^{s_1}}{1 + e^{s_1}}, \qquad \tau = \frac{1}{\|L\|} \cdot \frac{e^{s_2 + s_3}}{1 + e^{s_2}}, \qquad \sigma = \frac{1}{\|L\|} \cdot \frac{e^{s_2 - s_3}}{1 + e^{s_2}} \qquad (5.4)$$

with $s_1, s_2, s_3 \in \mathbb{R}$. Here, $s_2$ determines how close the parameters $\sigma$ and $\tau$ are to the constraint $\sigma\tau\|L\|^2 < 1$, while $s_3$ determines the trade-off between $\tau$ and $\sigma$.

**PDHG method without constraints on the parameters.** Here we train the same parameters $\theta, \tau, \sigma$ as in the PDHG method. However, we do not make the change of variables (5.4), therefore, no constraints on $\theta$, $\tau$, and $\sigma$ are enforced in the training. This means that the resulting scheme might not correspond to a globally convergent optimization algorithm.

**Proposed method from Section 3.** This corresponds to optimal parameter selection for the method (3.16) on the family of cost functions (5.3). To adhere to the constraints in the assumptions in Theorem 3.1, we have used the same kind of variable change as in (5.4), namely

$$a_{21} = \frac{2e^{s_1}}{1 + e^{s_1}}, \quad c_{21} = \frac{2e^{s_2}}{1 + e^{s_2}}, \quad \sigma = \frac{K}{\|L\|} \cdot \frac{e^{s_3 - s_4}}{1 + e^{s_3}}, \quad \tau = \frac{K}{\|L\|} \cdot \frac{e^{s_3 + s_4}}{1 + e^{s_3}}$$

with $s_1, \ldots, s_4 \in \mathbb{R}$, where $K = \frac{a_{21}^2(2 - a_{21})(2 - c_{21})}{(a_{21} + c_{21} - a_{21}c_{21})^2}$, as in (3.7).

**Parametrization proposed in Section 4.2.** Here, we have trained schemes of the form (4.3). We have done this for constant sequences of matrices $\boldsymbol{A}_1$, $\boldsymbol{A}_2$, $\boldsymbol{B}_1$, $\boldsymbol{B}_2$, $\boldsymbol{C}$, and $\boldsymbol{D}$. We restricted ourselves to the sizes $N = M = 2$ and $N = M = 3$.
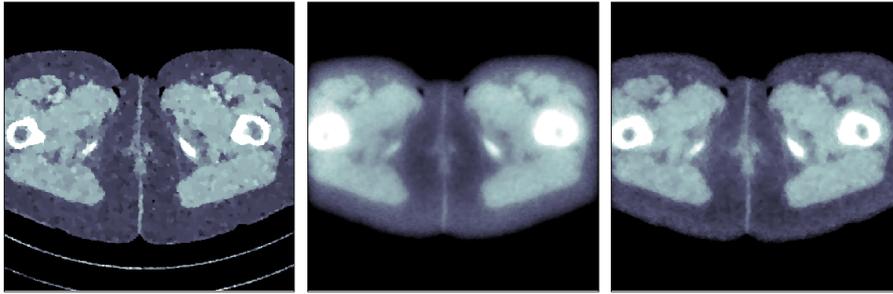
## 5.2 Performance of the trained methods

To obtain an estimation of the true optimal value of (5.3), we have run $1\,000$ iterations of PDHG with parameters as in [47]. In Table 1 we show the difference between the obtained objective function value and the minimal objective function value, averaged over 100 samples. As can be seen, the scheme proposed in Section 4 with $N = M = 3$ performs best at 10 iterations. Moreover, a general trend seems to be that more parameters in the algorithms improve the performance. Finally, the results from one specific phantom are presented as reconstructions in Figure 3. Note that the reconstruction by PDHG with parameters as in [47] is left out due to the page layout.
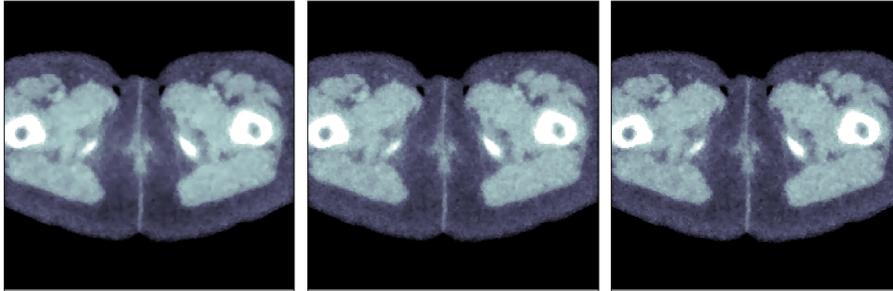
---

[‡]https://github.com/aringh/data-driven_nonsmooth_optimization

Table 1: Loss function values for the CT reconstruction after 10 iterations. The values given are of the form $\frac{1}{100}\sum_{i=1}^{100} H_{b_i}(x_{10}) - H_{b_i}(x_i^*)$, i.e., the difference of the obtained objective function value and an estimate of the true minimum objective function value $H_{b_i}(x_i^*)$ corresponding to data $b_i$, averaged over 100 samples.

| Method | Loss function values |
| --- | --- |
| PDHG with parameters from [47] | 109.93 |
| Trained PDHG with constraints on parameters | 82.381 |
| Trained solver (3.16) | 24.183 |
| Trained PDHG without constraints on parameters | 27.761 |
| Trained scheme of type (4.3) with $N = M = 2$ | 20.024 |
| Trained scheme of type (4.3) with $N = M = 3$ | **14.905** |



(a) TV reconstruction.     (b) Trained PDHG with constraints on parameters.     (c) Trained solver (3.16).



(d) Trained PDHG without constraints on parameters.     (e) Trained scheme of type (4.3) with $N = M = 2$.     (f) Trained scheme of type (4.3) with $N = M = 3$.

Figure 3: Reconstruction with data from a phantom that was not used in the training. The TV reconstruction, to which they should be compared, is shown in (a). All reconstructions use 10 steps. The phantom takes values between $[0.0, 2.33]$, but all images are shown using a window set to $[0.8, 1.2]$ in order to enhance contrast of clinically more relevant details.
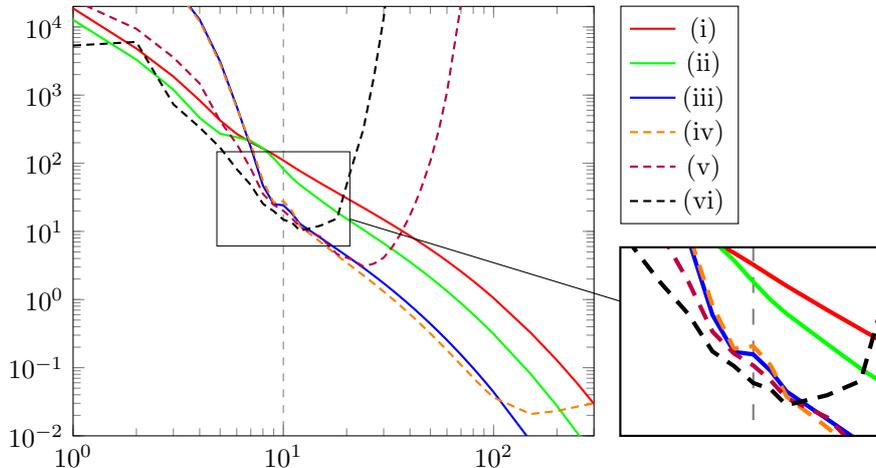
Figure 4: The figure shows the values $\frac{1}{100}\sum_{i=1}^{100} H_{b_i}(x_n) - H_{b_i}(x_i^*)$, where $H_{b_i}(x_i^*)$ is an estimate of the true minimum objective function value corresponding to data $b_i$, of several reconstruction methods as a function of the iteration number $n$. Solid lines are real optimization solvers, dotted lines are schemes that might not converge to the true optimal solution. (i) PDHG with parameters as in [47], (ii) PDHG with trained parameters with constraints, (iii) proposed solver (3.16) with trained parameters, (iv) PDHG with trained free parameters, (v) proposed scheme (4.3) with $N = M = 2$, and (vi) proposed scheme (4.3) with $N = M = 3$.

### 5.2.1 Generalization to other iteration numbers

Figure 4 shows the objective function value (5.3) as a function of the iteration number, i.e., how well the learned algorithms generalize to iteration numbers they are not trained for. For the trained, convergent solvers, the objective function value keeps decreasing as expected. Furthermore, the solver proposed in (3.16) performs better than the others also when the number of iterations are increased, but poorer in the beginning. For the other schemes, it can be noted that, while training more parameters seems to increase the performance after 10 iterations, it also seems to decrease the generalizability of the algorithm with respect to an increase in the number of iterations.

### 5.2.2 Generalization to deblurring

Next, we investigate the generalizability of the trained networks to other optimization problems by replacing the forward operator T in (5.3) with a convolution. This corresponds to another TV problem in imaging, namely image deblurring.

Clearly, the trained networks that correspond to optimization solvers with convergence guarantees can be applied to other convex optimization problems. (Note that we still normalize the operators to have operator norm one so that the assumptions in Theorem 3.1 do not change.) However, nothing guarantees that parameters that give fast convergence on one type of problems will also give fast convergence on another one.

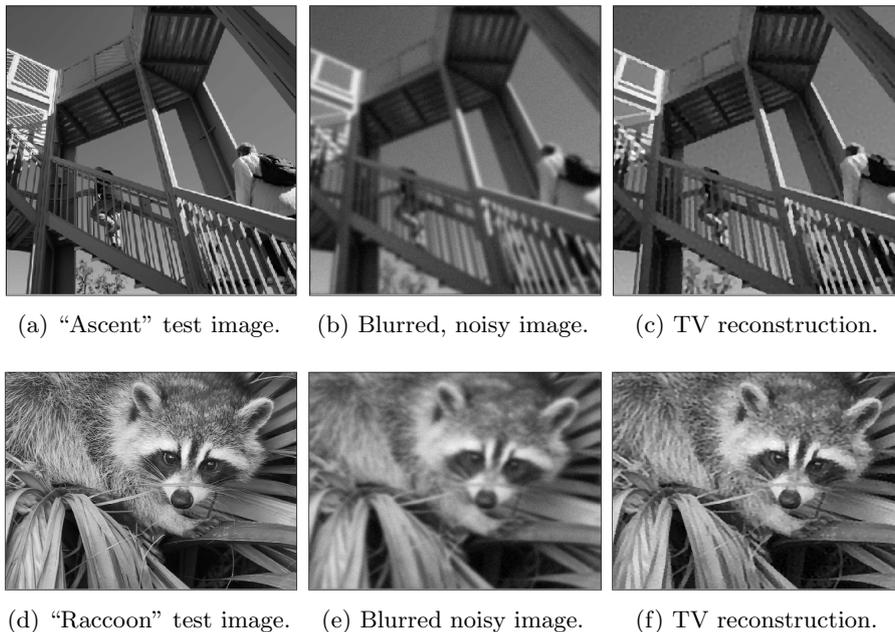Two example images are shown in Figure 5. The images in Figure 5d–5f,

| (a) "Ascent" test image. | (b) Blurred, noisy image. | (c) TV reconstruction. |



| (d) "Raccoon" test image. | (e) Blurred noisy image. | (f) TV reconstruction. |

Figure 5: Example images used for the deblurring problem in Section 5.2.2.

corresponding to the "Raccoon" test image, are of size $1024 \times 768$ and use a different regularization parameter. Blurring has been done with Gaussian kernels. For the "Ascent" test image, the kernel has a standard deviation of approximately three pixels in each direction, whereas for the "Raccoon" test image, the kernel has a standard deviation of approximately four pixels in the up-down and six pixels in the left-right direction. As for the sinograms in the CT example, 5% white noise has been added to the blurred images. Again, to obtain an estimation of the true optimal value of we have run 1 000 iterations of PDHG with parameters as in [47]. For each algorithm, the difference between the obtained objective function value and minimal objective function value is presented in Table 2, and the deblurred images are shown in Figures 6 and 7. Again, the reconstruction by PDHG with parameters as in [47] is left out due to the page layout.

The method with $N = M = 3$ does not generalize well. However, the method with $N = M = 2$ generalizes, and the optimization algorithm from Section 3, with trained parameters, is one of the best on these two test problems.

# 6    Conclusions and future work

In this work, we have first proposed a new solver for maximally monotone inclusion problems and proved convergence guarantees. In particular, we have also proposed a new convergent primal-dual proximal solver for convex optimization problems. Further, we have investigated new aspects of learning an optimization solver. This is particularly relevant in inverse problems where one can parametrize the objective function by data, leaving the other parts unchanged. This can, in fact, also be interpreted as learning a pseudo-inverse of the forward operator in

Table 2: Loss function values for the deblurring problem in Section 5.2.2. Here, $H_{b_i}(x_i^*)$ is an estimate of the true minimum objective function value corresponding to data $b_i$.

| Method | $H_{b_i}(x_{10}) - H_{b_i}(x_i^*)$ | |
|---|---|---|
| | Ascent | Raccoon |
| PDHG with parameters from [47] | 5.514 | 11.475 |
| Trained PDHG with constraints on parameters | 4.256 | 8.5126 |
| Trained solver (3.16) | **2.173** | 4.5898 |
| Trained PDHG without constraints on parameters | 2.204 | **4.4790** |
| Trained scheme of type (4.3) with $N = M = 2$ | 3.514 | 9.9139 |
| Trained scheme of type (4.3) with $N = M = 3$ | 208.37 | 873.33 |



(a) TV reconstruction.  (b) Trained PDHG with constraints on parameters.  (c) Trained solver (3.16).

(d) Trained PDHG without constraints on parameters.  (e) Trained scheme of type (4.3) with $N = M = 2$.  (f) Trained scheme of type (4.3) with $N = M = 3$.

Figure 6: Reconstructions with the trained algorithms for the "Ascent" image.

(a) TV reconstruction.    (b) Trained PDHG with constraints on parameters.    (c) Trained solver (3.16).

(d) Trained PDHG without constraints on parameters.    (e) Trained scheme of type (4.3) with $N = M = 2$.    (f) Trained scheme of type (4.3) with $N = M = 3$.
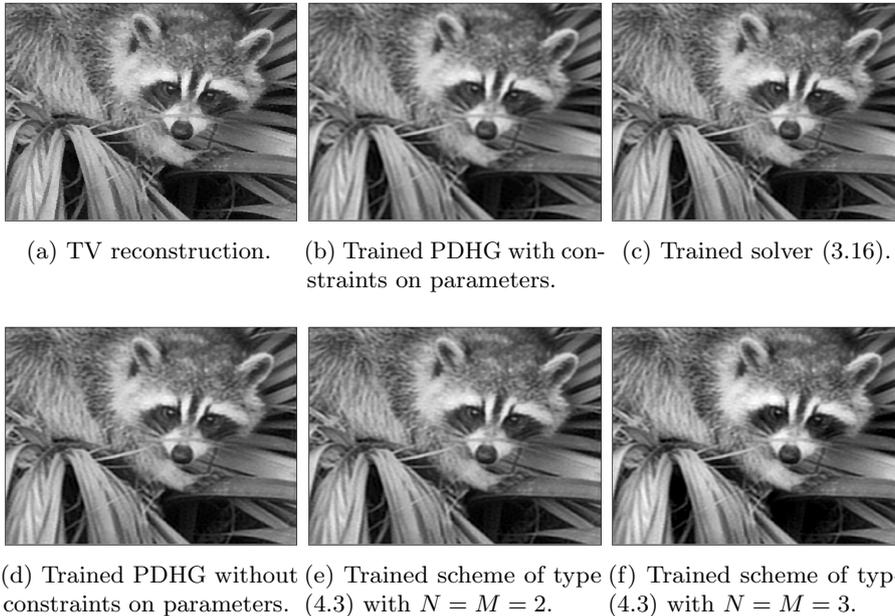
Figure 7: Reconstructions with the trained algorithms for the "Raccoon" image.

an unsupervised fashion. Moreover, the framework admits enforcing convergence and stability properties in the learning. We should emphasize that this implies a form of generalizability to other data, and even other forward operators, since the scheme cannot diverge.

There are several different directions in which the work from this article can be extended: Regarding the optimization perspective, one could investigate whether (3.8) can be further relaxed to introduce more free parameters while retaining convergence, e.g. by relaxing (3.6) or letting parameters vary in each iteration.

Also from a machine learning perspective, there are aspects to be further investigated:

- Since accelerated first-order algorithms like FISTA [12] can be parametrized by (4.3), does the learning result in a scheme with $\mathcal{O}(1/n^2)$ convergence rate for the objective function values when trained for $n$ iterations?

- Our numerical experiments suggest that training without "convergence constraints" gives the network more freedom and thereby improves accuracy. However, the resulting schemes seem to be unstable beyond the fixed number of iterates used for training. Is it true that, in general, convergence cannot be enforced by training alone?

- Is it possible to state and prove a time accuracy trade-off theorem, i.e., to estimate the error between the trained solver and the true solution to the optimization? If so, which properties of the underlying family of objective functions (training data) does this require?

# Acknowledgments

# References

[1] W. van Aarle, W.J. Palenstijn, J. Cant, E. Janssens, F. Bleichrodt, A. Dabravolski, J. De Beenhouwer, K.J. Batenburg, and J. Sijbers. "Fast and flexible X-ray tomography using the ASTRA toolbox." In: *Optics express* 24.22 (2016), pp. 25129–25147. DOI: `10.1364/OE.24.025129`.

[2] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mane, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viegas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. *TensorFlow: Large-scale machine learning on heterogeneous distributed systems.* 2016. arXiv: `1603.04467v2 [cs.DC]`.

[3] J. Adler, H. Kohr, and O. Öktem. *ODL 0.6.0.* Apr. 2017. DOI: `10.5281/zenodo.556409`. URL: `https://github.com/odlgroup/odl`.

[4] J. Adler and O. Öktem. "Learned primal-dual reconstruction." In: *IEEE Transactions on Medical Imaging* 37.6 (2018), pp. 1322–1332. DOI: `10.1109/TMI.2018.2799231`.

[5] J. Adler and O. Öktem. "Solving ill-posed inverse problems using iterative deep neural networks." In: *Inverse Problems* 33.12 (2017), p. 124007. DOI: `10.1088/1361-6420/aa9581`.

[6] J. Adler, A. Ringh, O. Öktem, and J. Karlsson. *Learning to solve inverse problems using Wasserstein loss.* 2017. arXiv: `1710.10898 [cs.CV]`.

[7] A. Alotaibi, P.L. Combettes, and N. Shahzad. "Solving coupled composite monotone inclusions by successive Fejér approximations of their Kuhn–Tucker set." In: *SIAM Journal on Optimization* 24.4 (2014), pp. 2076–2095. DOI: `10.1137/130950616`.

[8] M. Andrychowicz, M. Denil, S. Gomez, M. Hoffman, D. Pfau, T. Schaul, and N. de Freitas. "Learning to learn by gradient descent by gradient descent." In: *Advances in Neural Information Processing Systems.* Vol. 29. 2016, pp. 3981–3989. URL: `https://papers.nips.cc/paper/6461-learning-to-learn-by-gradient-descent-by-gradient-descent`.

[9] V. Barbu and T. Precupanu. *Convexity and optimization in Banach spaces.* 4th ed. Springer Monographs in Mathematics. Dordrecht: Springer, 2012. DOI: `10.1007/978-94-007-2247-7`.

[10]   H.H. Bauschke and P.L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces.* 2nd ed. CMS Books in Mathematics. New York: Springer, 2017. DOI: 10.1007/978-3-319-48311-5.

[11]   H.H. Bauschke, X. Wang, and L. Yao. "Examples of discontinuous maximal monotone linear operators and the solution to a recent problem posed by B.F. Svaiter." In: *Journal of Mathematical Analysis and Applications* 370.1 (2010), pp. 224–241. DOI: 10.1016/j.jmaa.2010.04.029.

[12]   A. Beck and M. Teboulle. "A fast iterative shrinkage-thresholding algorithm for linear inverse problems." In: *SIAM Journal on Imaging Sciences* 2.1 (2009), pp. 183–202. DOI: 10.1137/080716542.

[13]   D. Bertsekas. *Nonlinear programming.* 2nd ed. Belmont: Athena Scientific, 1999.

[14]   R.I. Boţ and E.R. Csetnek. "On the convergence rate of a forward-backward type primal-dual splitting algorithm for convex optimization problems." In: *Optimization* 64.1 (2015), pp. 5–23. DOI: 10.1080/02331934.2014.966306.

[15]   R.I. Boţ and C. Hendrich. "A Douglas–Rachford type primal-dual method for solving inclusions with mixtures of composite and parallel-sum type monotone operators." In: *SIAM Journal on Optimization* 23.4 (2013), pp. 2541–2565. DOI: 10.1137/120901106.

[16]   S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. "Distributed optimization and statistical learning via the alternating direction method of multipliers." In: *Foundations and Trends in Machine Learning* 3.1 (2011), pp. 1–122. DOI: 10.1561/2200000016.

[17]   L.M. Briceño-Arias and P.L. Combettes. "A monotone+skew splitting model for composite monotone inclusions in duality." In: *SIAM Journal on Optimization* 21.4 (2011), pp. 1230–1250. DOI: 10.1137/10081602X.

[18]   R.W. Brown, Y.-C. N. Cheng, E.M. Haacke, M.R. Thompson, and R. Venkatesan. *Magnetic resonance imaging: physical principles and sequence design.* John Wiley & Sons Ltd, 2014. DOI: 10.1002/9781118633953.

[19]   A.M. Bruckstein, D.L. Donoho, and M. Elad. "From sparse solutions of systems of equations to sparse modeling of signals and images." In: *SIAM review* 51.1 (2009), pp. 34–81. DOI: 10.1137/060657704.

[20]   A. Chambolle and T. Pock. "A first-order primal-dual algorithm for convex problems with applications to imaging." In: *Journal of Mathematical Imaging and Vision* 40.1 (2011), pp. 120–145. DOI: 10.1007/s10851-010-0251-1.

[21]   P.L. Combettes and J.-C. Pesquet. "Primal-dual splitting algorithm for solving inclusions with mixtures of composite, Lipschitzian, and parallel-sum type monotone operators." In: *Set-Valued and Variational Analysis* 20.2 (2012), pp. 307–330. DOI: 10.1007/s11228-011-0191-y.

[22]   P.L. Combettes and J.-C. Pesquet. "Proximal splitting methods in signal processing." In: *Fixed-point algorithms for inverse problems in science and engineering.* Ed. by H.H. Bauschke, R. Burachik, P.L. Combettes, V. Elser, D.R. Luke, and H. Wolkowicz. Vol. 49. Springer Optimization and Its Applications. New York: Springer, 2011, pp. 185–212. DOI: 10.1007/978-1-4419-9569-8_10.

[23] Y. Drori and M. Teboulle. "Performance of first-order methods for smooth convex minimization: a novel approach." In: *Mathematical Programming* 145.1–2 (2014), pp. 451–482. DOI: 10.1007/s10107-013-0653-0.

[24] J. Eckstein. "Splitting methods for monotone operators with applications to parallel optimization." PhD thesis. Department of Civil Engineering, Massachusetts Institute of Technology, 1989. URL: http://hdl.handle.net/1721.1/14356.

[25] J. Eckstein and D.P. Bertsekas. "On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators." In: *Mathematical Programming* 55.1-3 (1992), pp. 293–318. DOI: 10.1007/BF01581204.

[26] H.W. Engl, M. Hanke, and A. Neubauer. *Regularization of inverse problems.* Vol. 375. Mathematics and Its Applications. Kluwer Academic Publisher, 2000.

[27] K. Gregor and Y. LeCun. "Learning fast approximations of sparse coding." In: *Proceedings of the 27th International Conference on Machine Learning (ICML-10).* 2010, pp. 399–406. URL: http://yann.lecun.com/exdb/publis/pdf/gregor-icml-10.pdf.

[28] K. Hammernik, T. Klatzer, E. Kobler, M.P. Recht, D.K. Sodickson, T. Pock, and F. Knoll. "Learning a variational network for reconstruction of accelerated MRI data." In: *Magnetic resonance in medicine* 79.6 (2018), pp. 3055–3071. DOI: 10.1002/mrm.26977.

[29] B. He and X. Yuan. "Convergence analysis of primal-dual algorithms for a saddle-point problem: from contraction perspective." In: *SIAM Journal on Imaging Sciences* 5.1 (2012), pp. 119–149. DOI: 10.1137/100814494.

[30] J. Kaipio and E. Somersalo. *Statistical and computational inverse problems.* Vol. 160. Applied Mathematical Sciences. New York: Springer, 2005. DOI: 10.1007/b138659.

[31] D. Kim and J.A. Fessler. "Optimized first-order methods for smooth convex minimization." In: *Mathematical Programming* 159.1-2 (2016), pp. 81–107. DOI: 10.1007/s10107-015-0949-3.

[32] D. Kingma and J. Ba. *Adam: A method for stochastic optimization.* Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015. 2014. arXiv: 1412.6980v9 [cs.LG].

[33] S. Ko, D. Yu, and J.-H. Won. *On a Class of First-order Primal-Dual Algorithms for Composite Convex Minimization Problems.* 2017. arXiv: 1702.06234v2 [stat.ML].

[34] P. Latafat and P. Patrinos. "Asymmetric forward–backward–adjoint splitting for solving monotone inclusions involving three operators." In: *Computational Optimization and Applications* 68.1 (2017), pp. 57–93. DOI: 10.1007/s10589-017-9909-6.

[35] K. Li and J. Malik. *Learning to optimize.* 2016. arXiv: 1606.01885v1 [cs.LG].

[36] I. Loshchilov and F. Hutter. *SGDR: Stochastic gradient descent with warm restarts.* 2016. arXiv: 1608.03983v5 [math.OC].

[37]  M. Mardani, E. Gong, J.Y. Cheng, S. Vasanawala, G. Zaharchuk, M. Alley, N. Thakur, S. Han, W. Dally, J.M. Pauly, and L. Xing. "Deep generative adversarial networks for compressed sensing automates MRI." In: (2017). arXiv: 1706.00051 [cs.CV].

[38]  C. McCollough. "TU-FG-207A-04: Overview of the Low Dose CT Grand Challenge." In: *Medical Physics* 43.6Part35 (2016), pp. 3759–3760. DOI: 10.1118/1.4957556.

[39]  J.-J. Moreau. "Proximité et dualité dans un espace hilbertien." In: *Bulletin de la Société Mathématique de France* 93.2 (1965), pp. 273–299. DOI: 10.24033/bsmf.1625.

[40]  F. Natterer. *The mathematics of computerized tomography.* Vol. 32. Classics in Applied Mathematics. SIAM, 2001. DOI: 10.1137/1.9780898719284.

[41]  F. Natterer and F. Wübbeling. *Mathematical methods in image reconstruction.* Mathematical Modelling and Computation. SIAM, 2001. DOI: 10.1137/1.9780898718324.

[42]  O. Öktem. "Mathematics of electron tomography." In: *Handbook of mathematical methods in imaging.* Ed. by O. Scherzer. New York: Springer, 2015, pp. 937–1031. DOI: 10.1007/978-1-4939-0790-8_43.

[43]  W.J. Palenstijn, K.J. Batenburg, and J. Sijbers. "Performance improvements for iterative electron tomography reconstruction using graphics processing units (GPUs)." In: *Journal of structural biology* 176.2 (2011), pp. 250–253. DOI: 10.1016/j.jsb.2011.07.017.

[44]  R. Pascanu, T. Mikolov, and Y. Bengio. *On the difficulty of training recurrent neural networks.* 2012. arXiv: 1211.5063v2 [cs.LG].

[45]  P. Putzky and M. Welling. "Recurrent inference machines for solving inverse problems." In: (2017). arXiv: 1706.04008 [cs.NE].

[46]  R.T. Rockafellar. "Monotone operators and the proximal point algorithm." In: *SIAM Journal on Control and Optimization* 14.5 (1976), pp. 877–898. DOI: 10.1137/0314056.

[47]  E.Y. Sidky, J.H. Jørgensen, and X. Pan. "Convex optimization problem prototyping for image reconstruction in computed tomography with the Chambolle–Pock algorithm." In: *Physics in Medicine and Biology* 57.10 (2012), pp. 3065–3091. DOI: 10.1088/0031-9155/57/10/3065.

[48]  A.B. Taylor, J.M. Hendrickx, and F. Glineur. "Smooth strongly convex interpolation and exact worst-case performance of first-order methods." In: *Mathematical Programming* 161.1–2 (2017), pp. 307–345. DOI: 10.1007/s10107-016-1009-3.

[49]  Y. Yang, J. Sun, H. Li, and Z. Xu. "Deep ADMM-Net for Compressive Sensing MRI." In: *Advances in Neural Information Processing Systems.* Ed. by D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett. Vol. 29. Curran Associates, 2016, pp. 10–18.