

# **Finite Element Methods and Iterative Algorithms for Total-Variation Problems**

MAXIMILIAN BOCHMANN

MASTER'S THESIS

Fakultät für Mathematik

TU CHEMNITZ

Juni 2019

Supervisors:

Prof. Dr. Roland Herzog  
Dr. Max Winkler

# Declaration

I hereby declare and confirm that this thesis is entirely the result of my own original work. Where other sources of information have been used, they have been indicated as such and properly acknowledged. I further declare that this or similar work has not been submitted for credit elsewhere.

Chemnitz, June 10, 2019

Maximilian Bochmann

# Contents

<b>1</b>	<b>Introduction and Motivation</b>	<b>1</b>
<b>2</b>	<b>Preliminaries and Problem Statement</b>	<b>2</b>
2.1	Preliminaries . . . . .	2
2.2	The proximal Point Algorithm . . . . .	5
2.3	Problem Statement and Existence of a Solution . . . . .	7
<b>3</b>	<b>Bregman’s Method and the ADMM Method</b>	<b>12</b>
3.1	Bregman’s Method . . . . .	12
3.2	Convergence Properties of Bregman’s Method . . . . .	14
3.3	The ADMM Method . . . . .	16
3.4	Choosing the Penalty Parameter and different Error Measures	19
3.5	ADMM Approach to solve the TV denoising Problem . . . . .	20
<b>4</b>	<b>Chambolle’s and Pock’s Algorithm</b>	<b>25</b>
4.1	The Algorithm . . . . .	25
4.2	Application to the Mosolov Model . . . . .	27
4.3	Connections between the Algorithms . . . . .	29
<b>5</b>	<b>Discretization</b>	<b>31</b>
5.1	Nédélec’s Edge Elements . . . . .	31
5.2	Global $H_{\text{curl}}$ Conformity, Orientation of the Edges and FEM Matrices . . . . .	34
5.3	The discrete Problems . . . . .	36
<b>6</b>	<b>Numerical Results</b>	<b>44</b>
6.1	Two-dimensional Flow Problem . . . . .	44
6.2	Image denoising Problem . . . . .	48
<b>7</b>	<b>Summary and Future Work</b>	<b>51</b>
7.1	Summary . . . . .	51
7.2	Future Work . . . . .	51
	<b>References</b>	<b>52</b>

# Chapter 1

## Introduction and Motivation

During the last decades total variation minimization has evolved into an important tool for a range of different problems in fields like engineering or image processing. Important applications include problems such as optimal control or image denoising and restoration.

Total variation terms however pose difficulties for numerical optimization since they are non-differentiable. A common remedy hence is to regularize the corresponding terms to ensure differentiability. However this smoothing comes at the price of damping key features of the total variation, because of which it is used in the first place, e.g. preserving sharp edges for images.

A simple but very efficient method to solve such problems is the Alternating direction method of multipliers (ADMM), which was established earlier in the 1970s. It is a first order method and can be seen as an application of Bregman's method. At the same time it is closely related to the Augmented Lagrangian method and likewise a variety of other methods turn out to be related or even equivalent to it.

The key feature of this method is the splitting of the objective function into multiple independent parts which then allows iterative minimization. Similar to ADMM there is Chambolle and Pock's method, a primal-dual first order method using results from Fenchel duality.

We want to apply these methods to the usual image denoising problem and as well to a simplified Bingham flow problem, called the Mosolov model. A Bingham fluid is a viscoplastic substance, which acts as a rigid body when exposed to low shear stresses, but flows as a viscous fluid at high shear stresses. We use finite element methods to solve the corresponding iterative subproblems and we discuss different strategies for choosing corresponding finite element spaces and the corresponding impact on convergence of the iterative minimization methods.

All the methods have been implemented using MATLAB.

## Chapter 2

# Preliminaries and Problem Statement

### 2.1 Preliminaries

In this section we want to mention some basic notions and results necessary for the following outlines.

Alongside the common notion of convexity we will also need further notions of convexity.

**Definition 2.1.1.** (*Convexity*)

Let  $X$  be a normed linear space  $f : X \rightarrow \mathbb{R} \cup \{+\infty\}$ . We call  $f$  convex if for all  $x, y \in X$  and  $\alpha \in [0, 1]$  we have

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y).$$

The functional  $f$  is called strictly convex if for all  $x \neq y \in X$  and  $\alpha \in (0, 1)$

$$f(\alpha x + (1 - \alpha)y) < \alpha f(x) + (1 - \alpha)f(y).$$

Let  $f$  be differentiable on an open set  $U \subset X$ . We call  $f$  strongly convex with parameter  $\alpha > 0$  if for all  $x, y \in U$

$$\langle f'(x) - f'(y), x - y \rangle_{X^* \times X} \geq \alpha \|x - y\|_X^2.$$

**Remark 2.1.1.** A function which is strongly convex for some  $\alpha > 0$  is also strictly convex but not vice versa.

We are mainly interested in convex functionals which possess the following properties.

**Definition 2.1.2.** Let  $(X, \|\cdot\|_X)$  be a normed vector space. For a given functional  $f : X \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$  we define the following properties.

- (i) We call  $f$  proper if  $f(x) > -\infty$  and there is at least one  $x_0 \in X$  such that  $f(x_0) < \infty$ .
- (ii) We say  $f$  is lower semicontinuous in  $x_0 \in X$  if

$$f(x_0) \leq \liminf_{x \rightarrow x_0} f(x).$$

- (iii) We say  $f$  is coercive if for every sequence  $(x_n)_{n \in \mathbb{N}} \subset X$ , which possesses the property  $\|x_n\|_X \rightarrow \infty$ , we have

$$\lim_{n \rightarrow \infty} f(x_n) = +\infty.$$

In this work we will always assume that the functionals we are working with are proper.

**Definition 2.1.3.** (Differentiability)

Let  $X$  and  $Y$  be normed linear spaces,  $U \subset X$  open and  $f : X \rightarrow Y$ . The directional derivative of  $f$  at  $u \in U$  in direction  $v \in X$  is defined as

$$df(u, v) := \left. \frac{d}{dt} f(u + tv) \right|_{t=0} = \lim_{t \rightarrow 0} \frac{f(u + tv) - f(u)}{t}$$

If  $df(u, v)$  exists for all  $v \in X$  and  $df(u, \cdot) \in \mathcal{L}(X, Y)$  then  $f$  is said to be Gâteaux-differentiable in  $u$ . If the limit exists uniformly in  $v$ , i.e.

$$f(u + v) = f(u) + df(u, v) + o(\|v\|_X)$$

for  $v \rightarrow 0$  then  $f$  is said to be Fréchet-differentiable in  $u$  and we write  $f'(u)$  instead of  $df(u, \cdot)$ .

As we will deal with non-smooth convex optimization problems, the concept of the subdifferential is indispensable when formulating necessary and sufficient optimality conditions respectively.

**Definition 2.1.4.** (Subdifferential)

Let  $(X, \|\cdot\|_X)$  be a normed vector space and  $f : X \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex functional then we call the set

$$\partial f(u) := \{s \in X^* : f(x) \geq f(u) + \langle s, x - u \rangle_{X^* \times X}, x \in X\}$$

the subdifferential of  $f$  evaluated at  $u \in X$ . An element  $s \in \partial f(u)$  is called subgradient.

The subdifferential is a set-valued mapping and has the property to be a maximally monotone operator, (see e.g. [1, Chapter 2]), if  $f$  is a convex and

lower semicontinuous functional. This property is essential for the forthcoming algorithms.

If  $f$  is Fréchet-differentiable in  $u$  the subdifferential is a singleton and we have

$$\partial f(u) = \{f'(u)\}.$$

**Remark 2.1.2.** (*Subdifferential of a norm*)

For a norm on a linear space the subdifferential is equal to

$$\partial \|u\|_X = \{s \in X^* : \langle s, u \rangle_{X^* \times X} = \|u\|_X, \|s\|_{X^*} \leq 1\}.$$

Considering the squared norm we find that the subdifferential is equal to the duality mapping  $J_X$ , i.e.

$$\partial(\tfrac{1}{2}\|u\|_X^2) = J_X(u) = \{s \in X^* : \langle s, u \rangle_{X^* \times X} = \|u\|_X^2 = \|s\|_{X^*}^2\}.$$

**Theorem 2.1.1.** (*First order condition*)

Let  $f : X \rightarrow \mathbb{R} \cup \{+\infty\}$  be a functional and let  $u^* \in X$  be a minimizer of  $f$ , i.e.  $f(u^*) = \min_{u \in X} f(u)$ . Then the condition

$$0 \in \partial f(u^*)$$

is necessary. If  $f$  is convex the above condition is also sufficient.

To clarify the question whether a functional  $f$  has a (unique) minimizer we have the following theorem.

**Theorem 2.1.2.** (*Direct method of variational calculus*)

Let  $f : X \rightarrow \mathbb{R} \cup \{+\infty\}$  be a functional on a reflexive Banach space  $X$  which is convex, coercive and weakly lower-semicontinuous. Then there exists an element  $x^* \in X$  such that

$$E(x^*) = \min_{x \in X} E(x).$$

If  $E$  even is strictly convex then the minimizer  $x^*$  is unique.

*Proof.* Since  $f$  is coercive it must attain a minimizer in some ball  $B_r(0)$  with sufficiently large  $r$ .

As the unit ball of  $X^*$  is compact with respect to the weak\*-topology by Banach-Alaogulu's theorem and  $X$  is reflexive, we know that every minimizing sequence  $\{x_n\} \subset B_r(0)$  has a subsequence that converges to some element  $x^*$  in the weak topology.

Since  $f$  is also lower semicontinuous we have

$$\inf_x f(x) \leq f(x^*) \leq \liminf_{k \rightarrow \infty} f(x_{n_k}) = \inf_{x \in X} f(x),$$

i.e.  $x^*$  is a minimizer.

The uniqueness claim in case of strict convexity follows by the standard assumption for contradiction that there is another minimizer.  $\square$

For a convex function lower-semicontinuity with respect to the weak topology even implies lower-semicontinuity with respect to the strong norm topology, as is stated in the next Lemma.

**Lemma 2.1.1.** *For a convex function  $f : X \rightarrow \mathbb{R} \cup \{\infty\}$  lower semicontinuity and weak lower-semicontinuity are equivalent.*

*Proof.* See for example [3, section 6.2.2]. □

## 2.2 The proximal Point Algorithm

All the methods we are going to discuss exploit proximal operators, which map a given element  $z \in X$  in some Banach space  $X$  to the minimizer of a convex functional  $f$  which is closest to the specified element. This tradeoff between minimization and being close to some element can further be influenced by considering weights and grants more convenient properties than the original problem. Thus proximal operators can also be understood as a regularization technique.

In the following we are mainly interested in the set of functions

$$\Gamma_0(X) := \{f : X \rightarrow \mathbb{R} \cup \{\infty\} : f \text{ proper, lower semicontinuous and convex}\}.$$

Instead of minimizing  $f$  over  $X$  we consider the Moreau-Yosida regularization family given by

$$\text{mor}_{f,\gamma}(z) := \inf_{x \in X} \left\{ f(x) + \frac{1}{2\gamma} \|x - z\|_X^2 \right\} \quad (2.1)$$

**Lemma 2.2.1.** *Let  $X$  be a reflexive Banach space and  $f \in \Gamma_0(X)$  then the function  $g : X \rightarrow \mathbb{R} \cup \{\infty\}$  defined by*

$$g(x) = f(x) + \frac{1}{2\gamma} \|x - z\|_X^2$$

*has a unique minimizer for each  $z \in X$ .*

*Proof.* Since by duality a lower semicontinuous convex function  $f$  can be seen as the pointwise supremum of all affine functions  $h$  such that  $h \leq f$  we know that  $f$  can be bounded from below by an affine function. This means that  $g \rightarrow \infty$  if  $\|x\|_X \rightarrow \infty$  and hence  $g$  is coercive.

Clearly  $g$  is also convex since a norm on a linear space  $X$  is convex by the triangular inequality and  $g$  is also (weakly) lower-semicontinuous by virtue of Lemma 2.1.1. Thus by using the direct method it attains a minimizer.

This minimizer is unique since  $\|\cdot\|_X^2$  is strongly convex because

$$\langle J_X(u) - J_X(v), u - v \rangle_{X^* \times X} = \|u - v\|_X^2.$$

□



With the aid of Lemma 2.2.1 the following Definition is well-posed.

**Definition 2.2.1.** *Let  $f \in \Gamma_0(X)$  on a reflexive Banach space  $X$ . We call the mapping  $\text{prox}_{f,\gamma} : X \rightarrow X$ , defined by*

$$\text{prox}_{f,\gamma}(z) = \arg \min_{x \in X} \left\{ f(x) + \frac{1}{2\gamma} \|x - z\|_X^2 \right\}$$

*proximal operator of the functional  $f$ .*

The following result tells us that the fixed points of proximal operators are of particular interest.

**Lemma 2.2.2.** *Let  $X$  be a reflexive Banach space and  $f \in \Gamma_0(X)$ . Further let  $x^* \in X$  and  $\gamma > 0$  then the following statements are equivalent.*

- (i)  $x^* = \text{prox}_{f,\gamma}(x^*)$
- (ii)  $f(x^*) = \inf_{x \in X} f(x)$

*Proof.*

Due to the convexity of  $f$  we have that  $x^* = \text{prox}_{f,\gamma}(x^*)$  if and only if

$$\begin{aligned} 0 \in \partial \left( f(x^*) + \frac{1}{2\gamma} \|x^* - x^*\|_X^2 \right) &= \partial f(x^*) + \frac{1}{2\gamma} \partial (\|x^* - x^*\|_X^2) \\ &= \partial f(x^*) + \frac{1}{\gamma} J_X(x^* - x^*) = \partial f(x^*) \end{aligned}$$

where we were allowed to use the sum decomposition of the subdifferential since the norm is continuous. □

The proximal point-algorithm simply reads as

$$x_{k+1} = \text{prox}_{f,\gamma}(x_k).$$

This iteration can also be written using resolvent operators, i.e.

$$x_{k+1} = (I + \gamma \partial f)^{-1}(x_k).$$

Since the resolvent of a maximally monotone operator (e.g. the subdifferential) can be shown to be firmly nonexpansive, it is guaranteed that there exists a fixed point, c.f. [7]. These fixed points correspond to zeros of the monotone operator as we saw in Lemma 2.2.2. This is the basis of the proximal point algorithm and similar fixed point iterations such as the Douglas-Rachford splitting the latter of which the ADMM method is a special case of, as we will see soon.

The Douglas-Rachford splitting aims to minimize a convex functional consisting of two terms

$$\min f(x) + g(x).$$

If we start with an initial guess  $y_0$  for the dual variable  $y$  the Douglas-Rachford splitting is defined by the sequence

$$\begin{aligned} x_{k+1} &= \text{prox}_{f,\gamma}(y_k) \\ y_{k+1} &= y_k + \text{prox}_{g,\gamma}(2x_{k+1} - y_k) - x_{k+1}. \end{aligned} \quad (2.2)$$

## 2.3 Problem Statement and Existence of a Solution

In the following we always assume  $\Omega \subset \mathbb{R}^n$  to be a connected, open and bounded set with a Lipschitz boundary  $\partial\Omega$ .

For such domains the divergence theorem holds and thus allows us to use the generalized formula of partial integration. Assume we have a differentiable scalar field  $u$  and a differentiable vector field  $v$  then there holds

$$\int_{\Omega} u \operatorname{div}(v) \, dx = \int_{\partial\Omega} uv \cdot d\vec{s} - \int_{\Omega} v \cdot \nabla u \, dx \quad (2.3)$$

In the following we will use the total variation of a function as a mean of denoising.

**Definition 2.3.1.** (*Total variation*)

The total variation norm for  $u \in L^1(\Omega)$  is defined as

$$\|u\|_{\text{TV}} := \sup \left\{ \int_{\Omega} u(x) \operatorname{div} \phi(x) \, dx : \phi \in C_c^1(\Omega, \mathbb{R}^n), |\phi(x)|_{s^*} \leq 1 \text{ a.e.} \right\}$$

where  $|\cdot|_s$ ,  $s \in [1, \infty]$  denotes the usual  $s$ -norm in  $\mathbb{R}^d$  and  $s^*$  is the conjugate exponent of  $s$ , i.e.  $\frac{1}{s} + \frac{1}{s^*} = 1$ . Common choices are  $s = 2$  for isotropic denoising and  $s = 1$  for anisotropic denoising. We consider here isotropic denoising, i.e.  $s = 2$ .

Total variation is an indispensable tool for image denoising and in particular for the Rudin–Osher–Fatemi (ROF) model. Let  $f \in L^2(\Omega)$  and let  $\beta, \mu > 0$  be fixed constants then the ROF model can be written as

$$\min_{u \in \text{BV}(\Omega)} \frac{1}{2} \|u - f\|_{L^2(\Omega)}^2 + \beta \|u\|_{\text{TV}}, \quad (2.4)$$

where  $\text{BV}(\Omega)$  denotes the space of functions with bounded variation.

Similarly to the ROF model we consider the minimization of the following functional  $E : H^1(\Omega) \rightarrow \mathbb{R}$ , given by

$$E(u) := \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \beta \|u\|_{\text{TV}}. \quad (2.5)$$

The functional (2.5) can be seen to consist of two parts. The first part aims to minimize

$$\frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx,$$

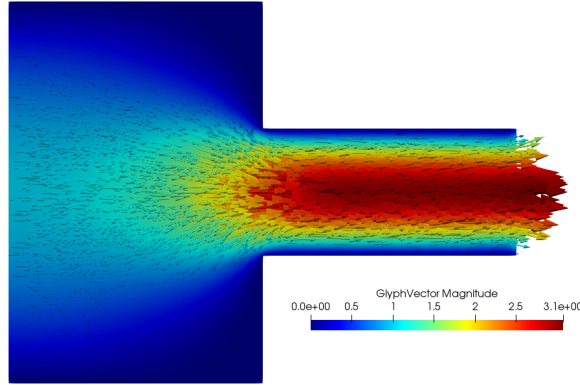
which is equivalent to solving a Poisson-equation. For simplicity we imply homogenous Dirichlet boundary conditions, i.e. we assume  $u \in H_0^1(\Omega)$ .

The second part, consisting of the total variation norm, forces the solution of the Poisson-problem to have few oscillations, i.e. piecewise-constant functions are preferred. We can control the influence of this property with the parameter  $\beta$ .

The model defined by problem (2.5) is called the Mosolov model, c.f. [9] or [8] respectively. It is a simplified case of the the general Bingham flow problem where the latter can be described using the Navier-Stokes equations.

### stationary Bingham Flow Problem

$$\begin{aligned}
 -\operatorname{div} \sigma + (\nabla u) \cdot u + \nabla p &= f && \text{in } \Omega \\
 \operatorname{div} u &= 0 && \text{in } \Omega \\
 \sigma &= 2\mu \varepsilon(u) + \beta \frac{\varepsilon(u)}{\|\varepsilon(u)\|} && \text{for } \|\varepsilon(u)\| \neq 0 \\
 \|\sigma\| &\leq \beta && \text{for } \|\varepsilon(u)\| = 0 \\
 u &= u_{\text{IN}} && \text{on } \Gamma_{\text{IN}} \\
 \sigma \cdot n &= 0 && \text{on } \Gamma_{\text{OUT}} \\
 u &= 0 && \text{on } \Gamma_0
 \end{aligned}$$



**Figure 2.1:** Flow velocity  $u(x_1, x_2)$  of a Bingham fluid through a bounded channel, consisting of two rectangles. The flow enters the channel from the left side and leaves on the right. For the remaining boundary parts we assume no slip conditions.

Here  $\varepsilon(u) = \frac{1}{2} (\nabla u + (\nabla u)^\top)$  is the rate-of-strain tensor,  $p$  is the pressure exerted on the fluid,  $\beta$  is a fluid-specific yield coefficient and  $\mu$  is the dynamic viscosity.

The Mosolov model considers a flow through the cross section  $\Omega$  of a pipe, where  $f$  represents the force along the pipe's axis and the desired quantity  $u$  is the flow velocity orthogonal to the cross section.

For the study of the iterative methods to solve the minimization problems we will confine to the Mosolov model (2.5). In this regard we want to consider the existence and uniqueness of a solution to problem (2.5). For this aspect we define the seminorm

$$|u|_{H^1(\Omega)}^2 := \int_{\Omega} |\nabla u|_2^2 \, dx$$

on  $H^1(\Omega)$ . By Poincaré's inequality there is a constant  $C$  such that for all functions  $u \in H_0^1(\Omega)$

$$\|u\|_{L^2(\Omega)} \leq C |u|_{H^1(\Omega)}.$$

Thus on the subspace  $H_0^1(\Omega)$  of  $H^1(\Omega)$  the  $H^1$ -norm and  $H^1$ -seminorm are equivalent by Poincaré's-inequality since we have

$$\|u\|_{H^1(\Omega)}^2 = \|u\|_{L^2(\Omega)}^2 + |u|_{H^1(\Omega)}^2 \leq (1 + C^2) |u|_{H^1(\Omega)}^2.$$

From now on we equip  $H_0^1(\Omega)$  with the  $H^1$ -seminorm.

As we assume  $u$  to be (weakly) differentiable we may express the total variation norm of  $u$  via  $\nabla u$ , as the following Lemma proposes, c.f. [3, Lemma 6.103].

**Lemma 2.3.1.** *Let  $u \in W^{1,1}(\Omega)$  then*

$$\|u\|_{\text{TV}} = \int_{\Omega} |\nabla u|_s \, dx.$$

*Proof.* (Only for the case  $s = 2$ )

Let  $\phi \in C_c^1(\Omega, \mathbb{R}^n)$  with  $|\phi(x)|_2 \leq 1$  almost everywhere. Considering that we have  $\phi = 0$  on  $\partial\Omega$  and using the formula of partial integration (2.3) we obtain

$$\int_{\Omega} u(x) \operatorname{div} \phi(x) \, dx = - \int_{\Omega} \nabla u \cdot \phi(x) \, dx \leq \int_{\Omega} |\nabla u|_2 \, dx$$

Thus

$$\|u\|_{\text{TV}} \leq \int_{\Omega} |\nabla u|_2 \, dx.$$

Let  $\nabla u \neq 0$  (otherwise the claim obviously holds). We can also assume without loss of generality that  $\nabla u \neq 0$  a.e. on  $\Omega$  since if  $\nabla u = 0$  on any

subset  $E \subset \Omega$  with positive measure then we can exclude this set from the integration domain.

By density of  $C_c^1(\Omega)$  in  $L^1(\Omega)$  and since  $\nabla u \in L^1(\Omega, \mathbb{R}^n)$ , there is a sequence  $\phi_n \subset C_c^1(\Omega, \mathbb{R}^n)$  with

$$\phi_n \rightarrow -\frac{\nabla u}{|\nabla u|_2}, \quad n \rightarrow \infty$$

in  $L^1(\Omega)$ .

Since the limit is a unit vector we may assume  $|\phi_n|_2 = 1$ .

Thus by Lebesgue's theorem of dominated convergence we have

$$\lim_{n \rightarrow \infty} \int_{\Omega} u(x) \operatorname{div} \phi_n(x) \, dx = - \lim_{n \rightarrow \infty} \int_{\Omega} \nabla u(x) \cdot \phi_n(x) \, dx = \int_{\Omega} |\nabla u|_2 \, dx$$

which proves the claim.  $\square$

With the aid of Theorem (2.1.2) we can now state the following Lemma.

**Lemma 2.3.2.** *The functional  $E : H_0^1(\Omega) \rightarrow \mathbb{R}$  defined by (2.5) has a unique minimizer  $u \in H_0^1(\Omega)$ .*

*Proof.* Let us define

$$F(u) := \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx.$$

The functional  $E$  is a sum of norms and a linear functional on  $H_0^1(\Omega)$  and thus convex. As  $F$  is clearly continuous with respect to the norm-topology of  $H_0^1$  and the total variation norm is a weaker norm than the  $H_0^1$ -norm we also have continuity of the functional  $E$ .

Now since  $E$  is convex and continuous we know that  $E$  is also lower semicontinuous.

To show strict convexity let  $u, w, v \in H_0^1(\Omega)$  then we have

$$F'(u)v = \mu (\nabla u, \nabla v)_{L^2(\Omega)} - \int_{\Omega} f v \, dx$$

which gives us

$$\langle F'(u) - F'(w), u - w \rangle = \mu |u - w|_{H^1(\Omega)}^2$$

which means that  $F$  is strongly convex and in particular that  $E$  is strictly convex. Finally we show that  $E$  is coercive. By Poincaré's-inequality we have

$$E(u) \geq \frac{\mu}{2} |u|_{H^1(\Omega)}^2 - C \|f\|_{L^2(\Omega)} |u|_{H^1(\Omega)} + \beta \|u\|_{\text{TV}}$$

and the right hand side goes to infinity if  $\|u\|_{H^1(\Omega)} \rightarrow \infty$  from which the coercivity follows. This concludes the proof.  $\square$

The minimization of  $E$  is equivalent to solving the following nonlinear elliptic boundary value problem

$$\begin{aligned} -\beta \operatorname{div} \Phi(\nabla u) - \mu \Delta u &= f \quad \text{in } \Omega \\ u &= 0 \quad \text{on } \partial\Omega. \end{aligned} \tag{2.6}$$

where  $\Phi(u) = \partial\|u\|_{L^1(\Omega)}$  is a maximally monotone operator.

**Definition 2.3.2.** *A function  $u \in H_0^1(\Omega)$  is said to be a weak solution of problem (2.6) if there  $w(x) \in L^2(\Omega)^d$  such that*

$$\begin{aligned} w(x) &\in \Phi(\nabla u(x)) \quad \text{a.e. for } x \in \Omega \\ \beta \int_{\Omega} w(x) \cdot \nabla v(x) \, dx + \mu \int_{\Omega} \nabla u(x) \cdot \nabla v(x) \, dx &= \int_{\Omega} f(x)v(x) \, dx \end{aligned}$$

for all  $v \in H_0^1(\Omega)$ .

Our next goal is to approximate weak solutions of (2.6) using the algorithms we are going to discuss in the next sections.

## Chapter 3

# Bregman's Method and the ADMM Method

We will now come to the techniques we will use for tackling the optimization problems.

### 3.1 Bregman's Method

Bregman's method was first mentioned in the 1960's and has since successfully been applicated in image denoising and other  $\ell_1$ -regularization problems. A fundamental concept for the method is the denotion of the Bregman distance.

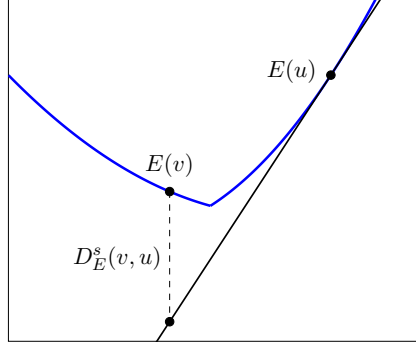
**Definition 3.1.1.** (*Bregman distance*)

Let  $E : X \rightarrow \mathbb{R}$  be a convex functional on a normed linear space  $X$  then the Bregman distance  $D_E^s(u, v)$  in  $u \in X$  is defined as

$$D_E^s(v, u) := E(v) - E(u) - \langle s, v - u \rangle_{X^* \times X}$$

where  $s \in \partial E(u)$ .

It describes the difference of the functional  $E$  and its linear model in  $u \in X$ , as depicted in the Figure below. Thus it is in some sense a measure for "closeness".



**Figure 3.1:** Bregman distance  $D_E^s(v, u)$

By convexity of the functional and the definition of the subdifferential  $D_E^s(v, u)$  is always non-negative. For a strictly convex functional  $E$  we have the equivalence  $D_E^s(v, u) = 0$  if and only if  $u = v$ . The Bregman distance however is not a distance in the sense of a metric in general as it does not satisfy the triangular inequality and it is not symmetric either.

For a given regularization functional  $J : X \rightarrow \mathbb{R} \cup \{\infty\}$  and a fitting functional  $H : X \rightarrow \mathbb{R} \cup \{\infty\}$  on a reflexive Banach space  $X$  we consider the following constrained minimization problem

$$\begin{aligned} \min_u \quad & J(u) \\ \text{subject to} \quad & H(u) = 0. \end{aligned} \tag{3.1}$$

This is a standard form of many regularization problems. It covers for example the basis pursuit problem for  $J(u) = \|u\|_1$  and  $H(u) = \frac{1}{2}\|Ax - b\|_2^2$  or the ROF model, if we set  $J(u) = \|u\|_{TV}$  and  $H(u) = \frac{1}{2}\|f - u\|_{L^2(\Omega)}^2$ .

We assume  $J$  and  $H$  to satisfy the following assumption.

**Assumption 3.1.1.**

- (i) *The regularization functional is convex, lower semicontinuous and bounded from below.*
- (ii) *The fitting functional is convex, Frechet-differentiable and bounded from below.*

We can transform (3.1) into an unconstrained problem by adding the constraint as a penalty and solving

$$\min_{u \in X} J(u) + \gamma H(u) \tag{3.2}$$

where,  $\gamma > 0$  is a penalty parameter. Then for  $\gamma \rightarrow \infty$  we get a solution of our original problem.



Bregman and later Goldstein and Osher proposed the following iterative strategy to solve (3.2), c.f. [5].

Given  $u_k \in X$  and  $s_k \in \partial J(u_k)$  solve

$$u_{k+1} = \arg \min_{u \in X} D_J^{s_k}(u, u_k) + \gamma H(u).$$

The optimality condition tells us that

$$0 \in \partial J(u_{k+1}) - s_k + \gamma \partial H(u_{k+1}).$$

As we assume that  $H$  is differentiable and considering that we require  $s_{k+1} \in \partial J(u_{k+1})$  we get

$$s_{k+1} = s_k - \gamma H'(u_{k+1}).$$

In the initial step we need to choose  $u_0$  and  $s_0 \in \partial J(u_0)$ . As  $J$  is bounded from below we may without loss of generality also assume that  $J$  is non-negative, otherwise we can define  $\hat{J} = J - \inf_u J$ , which will yield the same minimizers.

For this reason we can choose  $u_0 = s_0 = 0$  as

$$J(u) \geq J(0) = 0.$$

This leads to the following algorithm to solve problem (3.2).

---

**Algorithm 1:** Bregman's method

---

**Input** :  $u_0 = 0, s_0 = 0$

**Output:**  $u_k, s_k$

- 1 Set  $k \leftarrow 0$
  - 2  $u_{k+1} = \arg \min_{u \in X} D_J^{s_k}(u, u_k) + \gamma H(u)$
  - 3  $s_{k+1} = s_k - \gamma H'(u_{k+1})$
  - 4 Convergence test.
  - 5 Set  $k \leftarrow k + 1$  and go back to 2.
- 

As the minimization of  $D_J^{s_k}(u, u_k) + \gamma H(u)$  might be too hard or very expensive one can also consider a linearized version where one approximates  $H(u)$  by its first order Taylor expansion. To guarantee that the linear model is not too far off from  $H$  one can add a penalty

$$D_J^{s_k}(u, u_k) + \gamma \langle H'(u_k), u - u_k \rangle + \rho \|u - u_k\|_X^2$$

with  $\rho > 0$ .

## 3.2 Convergence Properties of Bregman's Method

Even if  $J$  and  $H$  satisfy Assumption 3.1.1 the iterates of Algorithm 1 might not exist. However if  $J$  or  $H$  are additionally coercive then by Theorem 2.1.2

the iterates are well defined. In any way if we assume that the iterates exist we can make some assertions about the convergence, c.f. [11], which we will state in the following.

For this purpose we define

$$E_k(u) := D_J^{s_k}(u, u_k) + \gamma H(u).$$

**Lemma 3.2.1.** (*Monotonic decrease*)

Assume that the iterates of Algorithm 1 exist. Then the sequence  $H(u_k)$  is monotonically nonincreasing.

$$\gamma H(u_{k+1}) \leq E_k(u_{k+1}) \leq \gamma H(u_k)$$

*Proof.* Since the Bregman distance is non-negative and  $u_{k+1}$  minimizes  $E_k(u)$  we have

$$\gamma H(u_{k+1}) \leq D_J^{s_k}(u_{k+1}, u_k) + \gamma H(u_{k+1}) \leq D_J^{s_k}(u_k, u_k) + \gamma H(u_k) = \gamma H(u_k).$$

□

**Theorem 3.2.1.** (*Convergence to a minimizer*)

Assume that the iterates of Algorithm 1 exist and that there is a minimizer  $u^*$  of  $H(u)$  such that  $J(u^*) < \infty$ . Then  $H(u_k)$  converges to a minimizer

$$H(u_k) \leq H(u^*) + \frac{J(u^*)}{\gamma k}.$$

*Proof.* We have

$$\begin{aligned} D_J^{s_{k+1}}(u, u_{k+1}) + D_J^{s_k}(u_{k+1}, u_k) - D_J^{s_k}(u, u_k) \\ = \langle s_k, u - u_k \rangle - \langle s_{k+1}, u - u_{k+1} \rangle - \langle s_k, u_{k+1} - u_k \rangle \\ = \langle s_{k+1} - s_k, u_{k+1} - u \rangle. \end{aligned}$$

Since  $s_{k+1} = s_k - \gamma H'(u_{k+1})$  we have that  $p_{k+1} := s_k - s_{k+1}$  is a subgradient of  $\gamma H$  at  $u_{k+1}$  and thus we have

$$\begin{aligned} D_J^{s_{k+1}}(u, u_{k+1}) + D_J^{s_k}(u_{k+1}, u_k) - D_J^{s_k}(u, u_k) &= \langle p_{k+1}, u - u_{k+1} \rangle \\ &\leq \gamma(H(u) - H(u_{k+1})). \end{aligned} \quad (3.3)$$

If we now sum from 0 to  $k-1$  in equation (3.3) and choose  $u = u^*$  we get

$$\sum_{j=0}^{k-1} D_J^{s_{j+1}}(u^*, u_{j+1}) + D_J^{s_j}(u_{j+1}^*, u_j) - D_J^{s_j}(u^*, u_j) \leq \gamma \sum_{j=0}^{k-1} (H(u^*) - H(u_{j+1}))$$

where the left hand side includes a telescoping series from which we get

$$D_J^{s_k}(u^*, u_k) + \sum_{j=0}^{k-1} \left[ D_J^{s_j}(u_{j+1}^*, u_j) - \gamma H(u^*) + \gamma H(u_{j+1}) \right] \leq D_J^{s_0}(u^*, u_0) = J(u^*).$$

With Lemma (3.2.1) we can further estimate

$$D_J^{s_k}(u^*, u_k) + \gamma k(H(u_k) - H(u^*)) \leq J(u^*)$$

and because  $D_J^{s_k}(u^*, u_k) \geq 0$  the assertion follows.  $\square$

### 3.3 The ADMM Method

Instead of (3.2) we now consider

$$\inf_u E(Bu) + F(u) \quad (3.4)$$

where  $B \in \mathcal{L}(X, Y)$  is a linear operator between Hilbert spaces  $X, Y$  such as functionals  $E$  and  $F$  which satisfy Assumption (3.1.1).

Since we are on Hilbert spaces the duality mappings  $J_X$  and  $J_Y$  respectively simply become the Riesz-isomorphism  $j_X$  and  $j_Y$ .

Instead of attempting to solve (3.4) we introduce  $d := Bu$  and consider the constrained optimization problem

$$\begin{aligned} \inf_{u,d} \quad & F(u) + E(d) \\ \text{subject to} \quad & d = Bu. \end{aligned} \quad (3.5)$$

We then transform this problem again into an unconstrained one with the aid of a quadratic penalty function

$$\min_{u,d} F(u) + E(d) + \frac{\gamma}{2} \|d - Bu\|_Y^2. \quad (3.6)$$

We can apply Bregman's method for this problem if we set

$$J(u, d) := F(u) + E(d) \quad \text{and} \quad H(u, d) := \frac{1}{2} \|d - Bu\|_Y^2.$$

This means we have to solve

$$(u_{k+1}, d_{k+1}) = \arg \min_{(u,d)} D_J^{s_k}(z, z_k) + \frac{\gamma}{2} \|d - Bu\|_Y^2 \quad (3.7)$$

where  $z = (u, d)$  and  $s_k = (s_k^u, s_k^d) \in \partial F(u) \times \partial E(d)$ .

For the subgradients the update is

$$s_{k+1}^u = s_k^u - \gamma B^* j_Y(Bu_{k+1} - d_{k+1}) \quad (3.8)$$

and

$$s_{k+1}^d = s_k^d - \gamma j_Y(d_{k+1} - Bu_{k+1}) \quad (3.9)$$

where  $B^* : Y^* \rightarrow X^*$  is the adjoint of  $B$ .

We can further simplify the method by introducing  $b_k$  as

$$b_0 = 0, \quad b_{k+1} = b_k + d_{k+1} - Bu_{k+1}. \quad (3.10)$$

Then we have  $s_k^u = \gamma B^* j_Y(b_k)$  and  $s_k^d = -\gamma j_Y(b_k)$ .

If we plug this into (3.7) we receive

$$\begin{aligned} (u_{k+1}, d_{k+1}) &= \arg \min_{(u,d)} D_f^{s_k}(z, z_k) + \frac{\gamma}{2} \|d - Bu\|_Y^2 \\ &= \arg \min_{(u,d)} J(z) - J(z_k) - \gamma \langle j_Y(b_k), B(u - u_k) \rangle_{Y^* \times Y} \\ &\quad + \gamma \langle j_Y(b_k), d - d_k \rangle_{Y^* \times Y} + \frac{\gamma}{2} \|d - Bu\|_Y^2 \\ &= \arg \min_{(u,d)} J(z) + \gamma \langle j_Y(b_k), d - Bu \rangle_{Y^* \times Y} + \frac{\gamma}{2} \|d - Bu\|_Y^2 \\ &= \arg \min_{(u,d)} J(z) + \frac{\gamma}{2} \|d - Bu + b_k\|_Y^2 \end{aligned}$$

At first glance this still seems even harder than the original problem since we have to minimize with respect to two variables, however since we have decoupled the functionals this structure allows us to split the problem into two subproblems, which are minimizing with respect to  $u$  on the one hand and then minimizing with respect to  $d$  on the other hand in each iteration. This leads to the following algorithm, called the split Bregman method or ADMM method.

---

**Algorithm 2:** ADMM method

---

**Input** :  $u_0 = 0, d_0 = 0, b_0 = 0$

**Output:**  $u_k, d_k, b_k$

1 Set  $k \leftarrow 0$

2  **$u$ -subproblem**

$$u_{k+1} = \arg \min_u F(u) + \frac{\gamma}{2} \|d_k - Bu + b_k\|_Y^2 \quad (3.11)$$

3  **$d$ -subproblem**

$$d_{k+1} = \arg \min_d E(d) + \frac{\gamma}{2} \|d - Bu_{k+1} + b_k\|_Y^2 \quad (3.12)$$

4  **$b$ -update**  $b_{k+1} = b_k + d_{k+1} - Bu_{k+1}$

5 Convergence test, e.g. check if  $\|u_{k+1} - u_k\|_X$  is less than a given tolerance.

6 Set  $k \leftarrow k + 1$  and go back to 2.

---

If  $E$  and  $F$  satisfy the corresponding conditions of Assumption 3.1.1 and  $F$  is additionally coercive the subproblems are well defined.

In that case we know from Theorem (3.2.1) that the fitting functional  $H(u_k, d_k)$  will converge to zero and thus  $(u_k, d_k)$  will converge to a minimizer of the regularization functional  $J(u, d)$  which together means a solution of the original problem (3.4).

**Remark 3.3.1.** *From Algorithm 2 we can see that any fixed point  $(u^*, d^*, b^*)$  of this algorithm is a solution to the original problem. Thus the condition to break the algorithm if  $\|u_{k+1} - u_k\|_X$  is sufficiently small is feasible if the subproblems are well defined.*

Alternatively we can derive the split Bregman method by considering the augmented Lagrangian for problem (3.4)

$$\mathcal{L}_\gamma(u, d, \lambda) = F(u) + E(d) + \langle j_Y(\lambda), d - Bu \rangle_{Y^* \times Y} + \frac{\gamma}{2} \|d - Bu\|_Y^2.$$

Minimizing  $\mathcal{L}_\gamma(u, d, \lambda)$  is equivalent to solving

$$\arg \min_{(u, d)} F(u) + E(d) + \frac{\gamma}{2} \left\| d - Bu + \frac{\lambda}{\gamma} \right\|_Y^2. \quad (3.13)$$

We solve this by iteratively minimizing first with respect to  $u$  for fixed  $d$  and then minimizing with respect to  $d$  for fixed  $u$ .

The update for the Lagrange multiplier is

$$\lambda_{k+1} = \lambda_k - \gamma(Bu_{k+1} - d_{k+1})$$

which, together with (3.13) yields the same procedure as Algorithm (2). Hence the method is called alternating direction method of multipliers (ADMM).

**Remark 3.3.2.** *If we reconsider equation (3.9) we see that the Lagrange multiplier is equal to*

$$\lambda_k = \gamma b_k = -j_Y^{-1}(s_k^d) = \gamma \sum_{i=1}^k (d_i - Bu_i).$$

since  $u_0 = d_0 = 0$ .

Note that for the ADMM method the penalty parameter  $\gamma$  may remain constant, i.e. we do not need to let  $\gamma \rightarrow \infty$ , as we would have to when using the Augmented Lagrangian method. This is an advantage of the ADMM-method since if we let  $\gamma \rightarrow \infty$  we might encounter numerical instabilities.

As  $\gamma$  remains constant we are also able to choose it such that the subproblems are easiest to solve, i.e. well conditioned or in a way that yields good convergence results. Alternatively one could think of an adaptive choice of the penalty parameter. We will discuss this in the next section.

Another advantage of the ADMM method is that it shows quick convergence when applied to specific objective functionals such as those which include a  $L^1$ -regularization term, c.f. [5].

### 3.4 Choosing the Penalty Parameter and different Error Measures

The penalty parameter  $\gamma$  has a huge influence on the convergence properties of the ADMM method. As a penalty parameter the natural choice for  $\gamma$  would be a value of large magnitude to enforce the penalty in the best way. However depending on the initial choice for  $u$ ,  $d$  and  $\lambda$  a large value for  $\gamma$  might lead to poor convergence results, especially in the case where the penalty is small for the initial values, however this choice is far off from the optimal solution. Thus  $\gamma$  has to be chosen carefully to provide good convergence results.

Besides choosing an appropriate penalty parameter it is also important to have meaningful error measures which indicate when to stop the iterations. Other than the distance from the last two iterates for  $u$  one can also use different indicators such as primal and dual residuals  $r^{\text{primal}}$  and  $r^{\text{dual}}$  showing the error related to primal and dual feasibility of the iterates. The primal residual  $r^{\text{primal}}$  clearly is the violation of the constraint  $d = Bu$ , i.e.

$$r_k^{\text{primal}} = d_k - Bu_k. \quad (3.14)$$

For the dual residual we consider the Lagrangian

$$\mathcal{L}(u, d, \lambda) = F(u) + E(d) + \langle j_Y(\lambda), d - Bu \rangle_{Y^* \times Y} \quad (3.15)$$

The optimal solution  $(u^*, d^*, \lambda^*)$  satisfies the optimality conditions

$$0 \in \partial_u \mathcal{L}(u^*, d^*, \lambda^*) = \partial F(u^*) - B^* j_Y(\lambda) \quad (3.16)$$

$$0 \in \partial_d \mathcal{L}(u^*, d^*, \lambda^*) = \partial E(d^*) + j_Y(\lambda) \quad (3.17)$$

which are the dual feasibility conditions.

Since  $d_{k+1}$  minimizes the augmented Lagrangian  $\mathcal{L}_\gamma(u_{k+1}, \cdot, \lambda_k)$  in each iteration  $k$  we have that

$$\begin{aligned} 0 \in \partial_d \mathcal{L}_\gamma(u_{k+1}, d_{k+1}, \lambda_k) &= \partial E(d_{k+1}) + j_Y(\lambda_k) + \gamma j_Y(d_{k+1} - Bu_{k+1}) \\ &= \partial E(d_{k+1}) + j_Y(\lambda_{k+1}) \end{aligned}$$

by using the update formula for  $\lambda$ . Thus we see that the iterates  $d_{k+1}$  and  $\lambda_{k+1}$  always satisfy the second dual feasibility condition (3.17).

Similarly, since  $u_{k+1}$  minimizes the augmented Lagrangian  $\mathcal{L}_\gamma(\cdot, d_k, \lambda_k)$ , we

get

$$\begin{aligned}
0 \in \partial_u \mathcal{L}_\gamma(u_{k+1}, d_k, \lambda_k) &= \partial F(u_{k+1}) - B^* j_Y(\lambda_k) + \gamma B^* j_Y(d_k - Bu_{k+1}) \\
&= \partial F(u_{k+1}) - B^* j_Y(\lambda_k) \\
&\quad + \gamma B^* j_Y(d_k - Bu_{k+1} + d_{k+1} - d_{k+1}) \\
&= \partial F(u_{k+1}) - B^* j_Y(\lambda_k) + \gamma B^* j_Y(d_{k+1} - Bu_{k+1}) \\
&\quad + \gamma B^* j_Y(d_k - d_{k+1}) \\
&= \partial F(u_{k+1}) - B^* j_Y(\lambda_{k+1}) + \gamma B^* j_Y(d_k - d_{k+1})
\end{aligned} \tag{3.18}$$

We see that we obtain dual feasibility condition (3.16) if the last term in (3.18) is equal to zero. Thus we set the dual residual to

$$r_k^{\text{dual}} = \gamma B^* j_Y(d_k - d_{k+1}). \tag{3.19}$$

Based on the primal and dual residual we have defined, we can now adapt the penalty parameter  $\gamma$ . When considering the Lagrangian (3.15) we see that if  $\gamma$  is large then the primal residual will be small. However the dual residual will be large as we can determine from (3.19). In the case where  $\gamma$  is small the situation is the opposite, i.e. the primal residual will be substantial whereas the dual residual will be small. That means we have to choose  $\gamma$  in a way such that the primal and dual residuals are balanced. As proposed in [14, section 4], this can be done by the choice

$$\gamma_{k+1} = \begin{cases} \tau \gamma_k & \text{if } \|r_k^{\text{primal}}\| > \zeta \rho \|r_k^{\text{dual}}\| \\ \tau^{-1} \gamma_k & \text{if } \|r_k^{\text{dual}}\| > \zeta^{-1} \rho \|r_k^{\text{primal}}\| \\ \gamma_k & \text{otherwise} \end{cases} \tag{3.20}$$

for parameters  $\mu, \tau, \zeta > 0$ . Often the choices  $\rho = 10$ ,  $\tau = 2$  and  $\zeta = 1$  yield far better results than a standardly chosen constant penalty parameter.

### 3.5 ADMM Approach to solve the TV denoising Problem

Minimizing functionals such as  $E$  defined by (2.5) or solving the ROF model (2.4) is hard because of the coupling of the  $L^2(\Omega)$ -norm and the TV-norm. Hence we decouple this problem by introducing  $d := \nabla u$ . Doing this and setting  $W = H_0^1(\Omega) \times Y$  we end up with the constrained optimization problem

$$\begin{aligned}
\min_{(u,d) \in W} \quad & \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \beta \int_{\Omega} |d|_2 \, dx \\
\text{subject to} \quad & d = \nabla u.
\end{aligned} \tag{3.21}$$

If we apply the ADMM method we have to iteratively minimize

$$\arg \min_{(u,d) \in W} \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \frac{\gamma}{2} \|\nabla u - d - b\|_Y^2 + \beta \int_{\Omega} |d|_2 \, dx \quad (3.22)$$

where  $b = \frac{\lambda}{\gamma}$ ,  $\gamma > 0$  is a penalty parameter and  $\lambda \in Y$  is a Lagrange-multiplier.

The most obvious choice for  $Y$  would be  $Y = L^2(\Omega)^d$ , but alternatively we could also use a convenient subspace of  $L^2(\Omega)^d$ .

We denote the curl of a three-dimensional differentiable vector field  $d$  by

$$\nabla \times d = \text{curl}(d) = \begin{pmatrix} \partial_2 d_3 - \partial_3 d_2 \\ \partial_3 d_1 - \partial_1 d_3 \\ \partial_1 d_2 - \partial_2 d_1 \end{pmatrix}$$

where  $\partial_i d_j = \frac{\partial d_j}{\partial x_i}$  and  $i, j = 1, 2, 3$ .

In  $2d$  we have

$$\text{curl}(d) = \partial_1 d_2 - \partial_2 d_1.$$

This leads to the definition of the subspace of  $L^2(\Omega)^d$  functions which have a curl in  $L^2(\Omega)^d$

$$H_{\text{curl}}(\Omega) := \left\{ v \in L^2(\Omega)^d : \text{curl}(d) \in L^2(\Omega)^d \right\}$$

equipped with the norm

$$\|u\|_{H_{\text{curl}}(\Omega)}^2 = \|u\|_{L^2(\Omega)}^2 + \|\text{curl}(u)\|_{L^2(\Omega)}^2.$$

This space occurs naturally for problems in the field of electromagnetism due to Maxwell's equation.

So far we have only assumed that  $\nabla u \in L^2(\Omega)^d$ , but we actually know more, namely that  $\nabla u \in H_{\text{curl}}(\Omega)$  and that the curl of  $\nabla u$  is zero almost everywhere in  $\Omega$ , i.e.  $\nabla u \in H_{\text{curl}}^0(\Omega)$ , where

$$H_{\text{curl}}^0(\Omega) := \left\{ v \in H_{\text{curl}}(\Omega) : \text{curl}(d) = 0 \right\}.$$

Thus instead of using the  $L^2(\Omega)$ -norm as a penalty as in (3.22) we could also consider penalizing with the stronger  $H_{\text{curl}}(\Omega)$ -norm, i.e. consider the following problem.

Let  $W = H_0^1(\Omega) \times H_{\text{curl}}(\Omega)$  find  $(\tilde{u}, \tilde{d}) \in \tilde{W}$  such that

$$(\tilde{u}, \tilde{d}) = \arg \min_{(u,d) \in \tilde{W}} \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \frac{\gamma}{2} \|\nabla u - d - b\|_{H_{\text{curl}}(\Omega)}^2 + \beta \int_{\Omega} |d|_2 \, dx \quad (3.23)$$

Note that both problems (3.22) and (3.23) will converge to the same solution as we let  $\gamma \rightarrow \infty$  since they both arise from the same original problem, which is minimizing (2.5). We just choose different penalty functionals. This means that the sequence of iterates will of course differ.



**The Choice**  $Y = L^2(\Omega)$ 

First we want to address the version with the  $L^2(\Omega)$ -penalty (3.22), i.e.  $Y = L^2(\Omega)$ .

The  $u$ -subproblem for the ADMM method then reads

$$u_{k+1} = \arg \min_{u \in H_0^1(\Omega)} \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \frac{\gamma}{2} \|\nabla u - d_k - b_k\|_{L^2(\Omega)}^2. \quad (3.24)$$

We set

$$G(u) = \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \frac{\gamma}{2} \|\nabla u - d_k - b_k\|_{L^2(\Omega)}^2$$

which is a convex and Fréchet-differentiable functional.

Hence the necessary and sufficient optimality condition for (3.24) reads  $G'(u)v = 0$  for all  $v \in H_0^1(\Omega)$ . This means solving the following linear variational problem.

For given  $d_k, b_k$  find  $u \in H_0^1(\Omega)$  such that for all  $v \in H_0^1(\Omega)$

$$(\mu + \gamma) \int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx + \gamma \int_{\Omega} (d_k + b_k) \cdot \nabla v \, dx \quad (3.25)$$

This is a weak formulation of a Dirichlet problem for the Poisson equation.

The associated strong formulatio reads

Find  $u \in C^2(\Omega) \cap C(\bar{\Omega})$  such that

$$\begin{aligned} -(\mu + \gamma)\Delta u &= f - \gamma \operatorname{div}(d_k + b_k) \quad \text{in } \Omega \\ u &= 0 \quad \text{on } \partial\Omega. \end{aligned}$$

**Lemma 3.5.1.** *Let  $\mu > 0$  and  $\gamma \geq 0$ . The weak-formulation (3.25) admits a unique solution  $u \in H_0^1(\Omega)$  for each  $f \in L^2(\Omega)$ ,  $d \in L^2(\Omega)$  and  $b \in L^2(\Omega)$ .*

*Proof.* This follows from the direct method of variational calculus and Theorem 2.1.1 or from the Lemma of Lax Milgram.  $\square$

The  $d$ -problem

$$d_{k+1} = \arg \min_{d \in L^2(\Omega)^d} \frac{\gamma}{2} \|\nabla u_{k+1} - d - b_k\|_{L^2(\Omega)}^2 + \beta \int_{\Omega} |d|_2 \, dx \quad (3.26)$$

also has a unique solution in  $L^2(\Omega)$  since the associated functional is strictly convex, continuous and coercive.

The optimality condition for (3.26) reads

$$\gamma(\nabla u - d - b) \in \beta \partial \|d\|_{L^1(\Omega)}$$

where

$$\|d\|_{L^1(\Omega)} = \int_{\Omega} |d|_2 \, dx.$$

The subdifferential of the  $L^1(\Omega)$  norm is

$$\partial\|d\|_{L^1(\Omega)} = \left\{ s \in L^\infty(\Omega) : \int_{\Omega} d \cdot s \, dx = \|d\|_{L^1(\Omega)}, \|s\|_{L^\infty(\Omega)} \leq 1 \right\}.$$

We can explicitly calculate this subdifferential and obtain

$$\nabla u - d - b = \begin{cases} \frac{\beta}{\gamma} \frac{d}{|d|_2} & : |d|_2 \neq 0 \\ B_{\frac{\beta}{\gamma}}(\mathbf{0}) & : |d|_2 = 0 \end{cases} \quad (3.27)$$

where  $B_r(\mathbf{0})$  denotes the ball of radius  $r > 0$  centered in  $\mathbf{0} \in \mathbb{R}^d$ . We see that the subdifferential for each  $x \in \Omega$  is just the projection onto the unit  $l^2$ -ball almost everywhere.

From (3.27) we conclude that if  $|\nabla u - b|_2 \leq \frac{\beta}{\gamma}$  we get  $d = 0$ . In the other case when  $|\nabla u - b|_2 > \frac{\beta}{\gamma}$  we obtain

$$d = \left( |\nabla u - b|_2 - \frac{\beta}{\gamma} \right) \frac{\nabla u - b}{|\nabla u - b|_2}$$

Shortly we can express this with shrinkage operators

$$d = \text{shrink} \left( \nabla u - b, \frac{\beta}{\gamma} \right) := \max \left( 0, |\nabla u - b|_2 - \frac{\beta}{\gamma} \right) \frac{\nabla u - b}{|\nabla u - b|_2} \quad (3.28)$$

We see that in the case of  $Y = L^2(\Omega)$  we can solve the  $d$ -problem (3.26) very efficiently without having to solve linear systems which is a great advantage of the ADMM method and makes it very efficient.

### The Choice $Y = H_{\text{curl}}(\Omega)$

Now we want to address solving the second version of our total variation denoising problem, i.e.

$$\arg \min_{(u,d) \in \tilde{W}} \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \frac{\gamma}{2} \|\nabla u - d - b\|_{H_{\text{curl}}(\Omega)}^2 + \beta \int_{\Omega} |d|_2 \, dx$$

Solving the  $u$ -subproblem remains exactly the same as for  $Y = L^2(\Omega)$  as

$$\|\nabla u - d - b\|_{H_{\text{curl}}(\Omega)}^2 = \|\nabla u - d - b\|_{L^2(\Omega)}^2 + \|\text{curl}(d + b)\|_{L^2(\Omega)}^2$$

meaning the penalty term is independant of  $u$ .

The  $d$ -subproblem now however becomes more difficult and more expensive to solve due to the  $H_{\text{curl}}$ -norm. We consider the regularized total variation term

$$\beta \int_{\Omega} |d + \varepsilon|_2 \, dx$$

then the  $d$ -problem is to solve find  $d \in H_{\text{curl}}(\Omega)$  such that for all  $w \in H_{\text{curl}}(\Omega)$

$$\begin{aligned} F(d, w) := & \int_{\Omega} (d - \nabla u + b) \cdot w \, dx + \int_{\Omega} \text{curl}(d + b) \cdot \text{curl}(w) \, dx \\ & + \frac{\beta}{\gamma} \int_{\Omega} \frac{d \cdot w}{\sqrt{|d|_2^2 + \varepsilon}} \, dx = 0 \end{aligned} \quad (3.29)$$

This is a nonlinear variational problem. We have regularized the problem with a regularization parameter  $\varepsilon > 0$  such that the associated functional in (3.29) is Fréchet-differentiable with respect to  $d$ . This allows us to solve (3.29) with Newton's method which means we have to solve a linear variational problem in each Newton step.

## Chapter 4

# Chambolle's and Pock's Algorithm

### 4.1 The Algorithm

Chambolle's and Pock's Algorithm is another first-order primal dual method for optimization problems of the form (3.4), i.e.

$$\inf_u E(Bu) + F(u)$$

We employ the same conventions for the functionals as in section 3.3. To formulate the corresponding dual problem we need the concept of the convex-conjugate of a function also referred to as the Fenchel-dual.

**Definition 4.1.1.** (*Fenchel-dual*)

Let  $f : X \rightarrow \mathbb{R}$  be a convex functional then the convex conjugate  $f^* : X^* \rightarrow \mathbb{R} \cup \{\infty\}$  of  $f$  is defined as

$$f(x^*) = \sup_{x \in X} \{ \langle x, x^* \rangle_{X \times X^*} - f(x) \}.$$

The associated dual problem of (3.4) is then given by

$$\sup_y -F^*(-B^*y) - E^*(y) \tag{4.1}$$

We assume that there exists a solution  $(u^\dagger, y^\dagger)$  of the primal and the dual problem respectively.

To derive the primal and dual first order optimality conditions for problem (3.4) it is useful to consider again the constrained problem

$$\begin{aligned} \inf_u \quad & E(d) + F(u) \\ \text{subject to} \quad & d = Bu. \end{aligned} \tag{4.2}$$

The corresponding Lagrangian can be written as

$$\begin{aligned} & \inf_u \sup_y E(d) + F(u) + \langle Bu - d, y \rangle_{Y \times Y^*} \\ &= \inf_u \sup_y F(u) - E^*(y) + \langle Bu, y \rangle_{Y \times Y^*} \end{aligned}$$

and thus the solution  $(u^\dagger, y^\dagger)$  has to satisfy the Karush-Kuhn-Tucker (KKT) conditions

$$Bu^\dagger \in \partial E^*(y^\dagger), \quad -B^*y^\dagger \in \partial F(u).$$

Since we are dealing with convex problems these conditions are also sufficient.

The algorithm is based on the following duality result. For a proof we refer to [12, section 31]

**Theorem 4.1.1.** (*Fenchel's Duality Theorem*)

Let  $F : X \rightarrow \mathbb{R} \cup \{\infty\}$ ,  $E : Y \rightarrow \mathbb{R} \cup \{\infty\}$  be proper, convex and lower semi-continuous functionals on Banach spaces  $X$  and  $Y$  and let  $B \in \mathcal{L}(X, Y)$ . Then the problems

$$\begin{aligned} P &= \inf_u E(Bu) + F(u) \\ P^* &= \sup_y -F^*(-B^*y) - E^*(y) \end{aligned}$$

satisfy weak-duality, i.e.  $P \geq P^*$ .

If there exists an  $u_0 \in X$  such that  $F(u_0) < \infty$ ,  $E(Bu_0) < \infty$  and  $E$  is continuous in  $Bu_0$  then we have strong duality, i.e.  $P = P^*$ .

Chambolle's and Pock's idea is to build a sequence of iterates  $u_k$  and  $y_k$  and add the residual of the current iterate in each step, i.e.

$$y_{k+1} = y_k + \sigma(B\bar{u}_k - \partial E^*(y_{k+1})) \quad (4.3)$$

and

$$u_{k+1} = u_k - \tau(B^*y_{k+1} + \partial F(u_{k+1})). \quad (4.4)$$

Since for the  $y$ -update we do not know  $u_{k+1}$  yet, we use an extrapolation to predict  $u_{k+1}$ , i.e.

$$\bar{u}_{k+1} = u_{k+1} + \theta(u_{k+1} - u_k). \quad (4.5)$$

This leads to Chambolle's and Pock's algorithm, c.f. [4].

**Algorithm 3:** Chambolle & Pock**Input** :  $u_0 = 0, y_0 = 0, \bar{u}_0 = 0, \sigma, \tau, \theta$ **Output:**  $u_k, y_k, \bar{u}_k$ 

- 1 Set  $k \leftarrow 0$
- 2  $y_{k+1} = (I + \sigma \partial E^*)^{-1}(y_k + \sigma B \bar{u}_k)$
- 3  $u_{k+1} = (I + \tau \partial F)^{-1}(u_k - \tau B^* y_{k+1})$
- 4  $\bar{u}_{k+1} = u_{k+1} + \theta(u_{k+1} - u_k)$
- 5 Convergence test, e.g. check if  $\|u_{k+1} - u_k\|_X$  is less than a given tolerance.
- 6 Set  $k \leftarrow k + 1$  and go back to 2.

Similarly as for the ADMM method, the choice of the parameters  $\tau, \sigma$  and  $\theta$  is essential for obtaining good convergence results.

## 4.2 Application to the Mosolov Model

To apply Chambolle's and Pock's algorithm to our problem we need the convex conjugate  $E^*(y)$  of the regularization functional  $E(d)$ , i.e.

$$\begin{aligned}
 E^*(y) &= \sup_{d \in L^2(\Omega)} \left\{ \langle d, y \rangle_{L^2(\Omega) \times L^2(\Omega)^*} - \beta \int_{\Omega} |d|_2 \, dx \right\} \\
 &= \sup_{d \in L^2(\Omega)} \left\{ \int_{\Omega} d \cdot y - \beta |d|_2 \, dx \right\}
 \end{aligned}$$

Since the second term of the integrand only depends on the Euclidean norm and the first term is maximized if  $d$  and  $y$  are linearly dependent we may assume  $d = cy$  for a constant  $c \in \mathbb{R}$ . Then we obtain

$$\begin{aligned}
 E^*(y) &= \sup_{c \in \mathbb{R}} \left\{ \int_{\Omega} c |y|_2 (|y|_2 - \text{sign}(c)\beta) \, dx \right\} \\
 &= \sup_{c \geq 0} \left\{ \int_{\Omega} c |y|_2 (|y|_2 - \beta) \, dx \right\} \\
 &= \begin{cases} \infty & |y|_2 > \beta \\ 0 & |y|_2 \leq \beta \end{cases} \\
 &= \delta_{\{|y|_2 \leq \beta\}}(y)
 \end{aligned} \tag{4.6}$$

where  $\delta_A : X \rightarrow \{0, \infty\}$  denotes the indicator function of the set  $A \subset \Omega$ .

Thus in the first step we have to solve

$$y_{k+1} = \arg \min_y \sigma \delta_{\{|y|_2 \leq \beta\}}(y) + \frac{1}{2} \|y - y_k - \sigma \nabla \bar{u}_k\|_{L^2(\Omega)}^2 \tag{4.7}$$

i.e. we need

$$0 \in \sigma \partial \delta_{\{|y|_2 \leq \beta\}}(y_{k+1}) + y_{k+1} - y_k - \sigma \nabla \bar{u}_k. \quad (4.8)$$

This problem can be solved very efficiently by using projections of the form

$$\text{proj}_{\{\|\cdot\|_2 \leq \beta\}}(x) = \begin{cases} x & |x|_2 \leq \beta \\ \beta \frac{x}{|x|_2} & \|x\|_2 > \beta \end{cases}$$

as we will see shortly.

The subdifferential of the indicator function  $\delta_A$  of a vector  $x \in \mathbb{R}^n$  is empty if  $x \notin A$ . If  $x \in A$  then it is the normal cone to the set  $A$ , i.e.

$$\partial \delta_A(x_0) = \{s \in \mathbb{R}^n : \langle s, x_0 - x \rangle \leq 0, x \in A\}.$$

In our case we obtain the pointwise expression

$$\partial \delta_{\{|y|_2 \leq \beta\}}(y) = \begin{cases} \{0\} & |y(x)|_2 < \beta \\ \mathbb{R}_{\geq 0} y & |y(x)|_2 = \beta \\ \emptyset & |y(x)|_2 > \beta \end{cases} \quad (4.9)$$

Combining (4.8) and (4.9) we find that the solution  $y_{k+1}$  to (4.7) is the projection onto the set  $\{|d|_2 \leq \beta\}$

$$y_{k+1} = \text{proj}_{\{\|\cdot\|_2 \leq \beta\}}(d_k + \sigma \nabla \bar{u}_k). \quad (4.10)$$

Now we tend to the solution of the  $u$ -problem. Since  $B^* = -\text{div}$  the problem reads

$$u_{k+1} = \arg \min_u \tau F(u) + \frac{1}{2} \|u - u_k - \tau j_X^{-1}(\text{div } y_{k+1})\|_{H_0^1(\Omega)}^2 \quad (4.11)$$

for

$$F(u) = \frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx.$$

Here  $j_X : H_0^1(\Omega) \rightarrow H_0^1(\Omega)^*$  denotes the Riesz-isomorphism on  $H_0^1(\Omega)$ . The optimality condition for (4.11) reads

$$0 \in \tau \partial F(u) + j_X(u - u_k - \tau j_X^{-1}(\text{div } y_{k+1}))$$

By

$$\langle v, u^* \rangle_{H_0^1(\Omega) \times H_0^1(\Omega)^*} = (v, u)_{H_0^1(\Omega)} = (\nabla v, \nabla u)_{L^2(\Omega)}$$

we see that  $u^* = j_X(u) = -\Delta u$ .

So the solution of (4.11) is a linear variation problem which reads.

Find  $u_{k+1} \in H_0^1(\Omega)$  such that for all  $v \in H_0^1(\Omega)$

$$(\tau\mu + 1) \int_{\Omega} \nabla u_{k+1} \cdot \nabla v \, dx = \int_{\Omega} \tau f v + (\nabla u_k - \tau y_{k+1}) \cdot \nabla v \, dx. \quad (4.12)$$

### 4.3 Connections between the Algorithms

Yet another way to look at the ADMM method is to consider the dual problems associated to the  $u$ -subproblem (3.11) and  $d$ -subproblem (3.12). For better readability we omit using the Riesz-isomorphism here.

We can use Moreau's identity

$$x = \text{prox}_{f,\gamma}(x) + \gamma \text{prox}_{f^*,\gamma^{-1}}\left(\frac{x}{\gamma}\right) \quad (4.13)$$

to easily derive the according optimality conditions.

Using  $u^*$  as the dual variable for  $Bu$  and  $d^*$  as the dual variable for  $d$  we get

$$\begin{aligned} u_{k+1}^* &= \arg \min_{u^*} F^*(-B^*u^*) + \frac{1}{2\gamma} \|u^* - \gamma(b_k + d_k)\| \\ d_{k+1}^* &= \arg \min_{d^*} E^*(d^*) + \frac{1}{2\gamma} \|d^* - \gamma(Bu_{k+1} - b_k)\|. \end{aligned}$$

Then (4.13) tells us

$$\gamma(b_k + d_k) = u_{k+1}^* + \gamma Bu_{k+1} \quad (4.14)$$

$$\gamma(Bu_{k+1} - b_k) = d_{k+1}^* + \gamma d_{k+1} \quad (4.15)$$

From (3.10) and (4.15) we further know  $b_{k+1} = -\gamma^{-1}d_{k+1}^*$ . Thus we have

$$u_{k+1}^* = \arg \min_{u^*} F^*(-B^*u^*) + \frac{1}{2\gamma} \|u^* - \gamma d_k + d_k^*\|$$

By (4.14) and summing (4.14) and (4.15) up we get

$$\begin{aligned} d_{k+1}^* &= \arg \min_{d^*} E^*(d^*) + \frac{1}{2\gamma} \|d^* - \gamma d_k + u_{k+1}^*\| \\ \gamma d_{k+1} &= \gamma d_k - u_{k+1}^* - d_{k+1}^*. \end{aligned} \quad (4.16)$$

If we introduce  $y_{k+1} = d_{k+1}^* + \gamma d_{k+1}$  this leads to

$$\begin{aligned} u_{k+1}^* &= \arg \min_{u^*} F^*(-B^*u^*) + \frac{1}{2\gamma} \|u^* + 2d_k^* - y_k\| \\ d_{k+1}^* &= \arg \min_{d^*} E^*(d^*) + \frac{1}{2\gamma} \|d^* + d_k^* - y_k + u_{k+1}^*\| \\ y_{k+1} &= y_k - u_{k+1}^* - d_k^*. \end{aligned} \quad (4.17)$$

Finally by substitution and interchanging the order we arrive at

$$y_{k+1} = y_k - \text{prox}_{F^*(-B^*\cdot),\gamma}(y_k - 2d_k^*) - d_k^* \quad (4.18)$$

$$d_{k+1}^* = \text{prox}_{E^*,\gamma}(y_{k+1}). \quad (4.19)$$



From (4.18)-(4.19) we see that this is the Douglas-Rachford splitting algorithm (2.2) we defined earlier applied to the dual problem

$$\sup_y -F^*(-B^*y) - E^*(y).$$

Hence the ADMM method coincides with the Douglas-Rachford method when applied to the dual problem.

Furthermore Chambolle and Pock's algorithm can be seen as a preconditioned variant of the ADMM method applied to the dual problem (4.1). For details on this we refer to [4, section 4.3].

## Chapter 5

# Discretization

In order to solve the problems discussed in the previous sections we will employ the Finite Element Method (FEM). We will use different kinds of finite elements such as the linear and continuous Lagrange elements ( $\mathcal{P}_1$  - elements) to discretize the  $u$ -problem. In order to solve the  $d$ -problem we will use different elements according to the choice of the space  $Y$ .

For  $Y = L^2(\Omega)$  we will use discontinuous Lagrange elements of lowest order ( $\mathcal{DG}_0$  - elements) and for the choice  $Y = H_{\text{curl}}(\Omega)$  we will use Nédélec's edge elements the latter of which we will describe further in the following section.

### 5.1 Nédélec's Edge Elements

Nédélec's edge elements most commonly appear in the context of electromagnetism. We will use them for the discretization of functions in  $H_{\text{curl}}(\Omega)$  such as  $\nabla u$ . We shortly repeat the definition of a finite element.

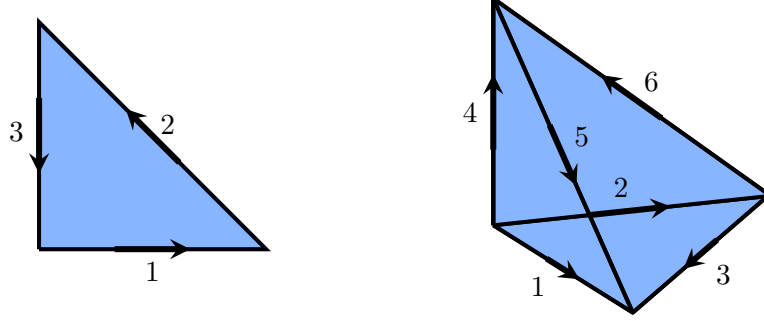
**Definition 5.1.1.** (*Ciarlet*)

A finite element is a triple  $(K, P, \Sigma)$  if the following conditions are met

- (a)  $K \subset \mathbb{R}^d$  is a compact polyeder with a non-empty interior.
- (b)  $P$  is a finite-dimensional vector space of dimension  $s \geq 1$  consisting of functions  $p : K \rightarrow \mathbb{R}^d$  defined on  $K$
- (c) There is a normed vector space  $V(K)$  of functions  $v : K \rightarrow \mathbb{R}$  such that  $P \subset V(K)$  and  $\Sigma = (\sigma_1, \dots, \sigma_s)$  is an ordered subset of the dual  $V(K)^*$  which is unisolvent, i.e. given that for some  $p \in P$  we have that  $\sigma_i(p) = 0$  for every  $i = 1, \dots, s$  then the implication  $p = 0$  in  $P$  is true.

The set  $\Sigma$  is also called the set of local degrees of freedom (dofs). Those functions  $p_j \in P$ ,  $j = 1, \dots, s$  which are uniquely determined by the relation  $\sigma_i(p_j) = \delta_{ij}$  are called local shape functions or basis functions.

As compact polyeders  $K$  we will use d-simplices, i.e. triangles for  $d = 2$  and tetrahedrons for  $d = 3$ .



**Figure 5.1:** Nédélec-element of first kind and lowest order with associated dofs.

**Definition and Proposition 5.1.1.** Let  $K \subset \mathbb{R}^d$  be a simplex and  $d \in \{2, 3\}$ . For  $d = 2$  we set

$$\mathcal{N}_0 := \left\{ [\mathbb{P}_0]^2 \oplus \mathbb{P}_0 \begin{pmatrix} x_2 \\ -x_1 \end{pmatrix} \right\}$$

and for  $d = 3$  we set

$$\mathcal{N}_0 := \left\{ [\mathbb{P}_0]^3 \oplus \left( x \times [\mathbb{P}_0]^3 \right) \right\}.$$

Then  $\mathcal{N}_0$  is a polynomial space of dimension  $s = \frac{1}{2}d(d+1)$ . Let  $e_i$  be the edges of the simplex  $K$ ,  $i = 1, \dots, s$  and let  $t_{e_i}$  be the unit vector that represents edge  $e_i$  oriented as in figure 5.1.

For  $p \in \mathcal{N}_0$  the local degrees of freedom are defined as the line integrals of the tangential part of  $p$  alongside the edges of  $K$ , i.e. for  $i = 1, \dots, s$

$$\sigma_i(p) = \int_{e_i} p \cdot t_{e_i} \, ds.$$

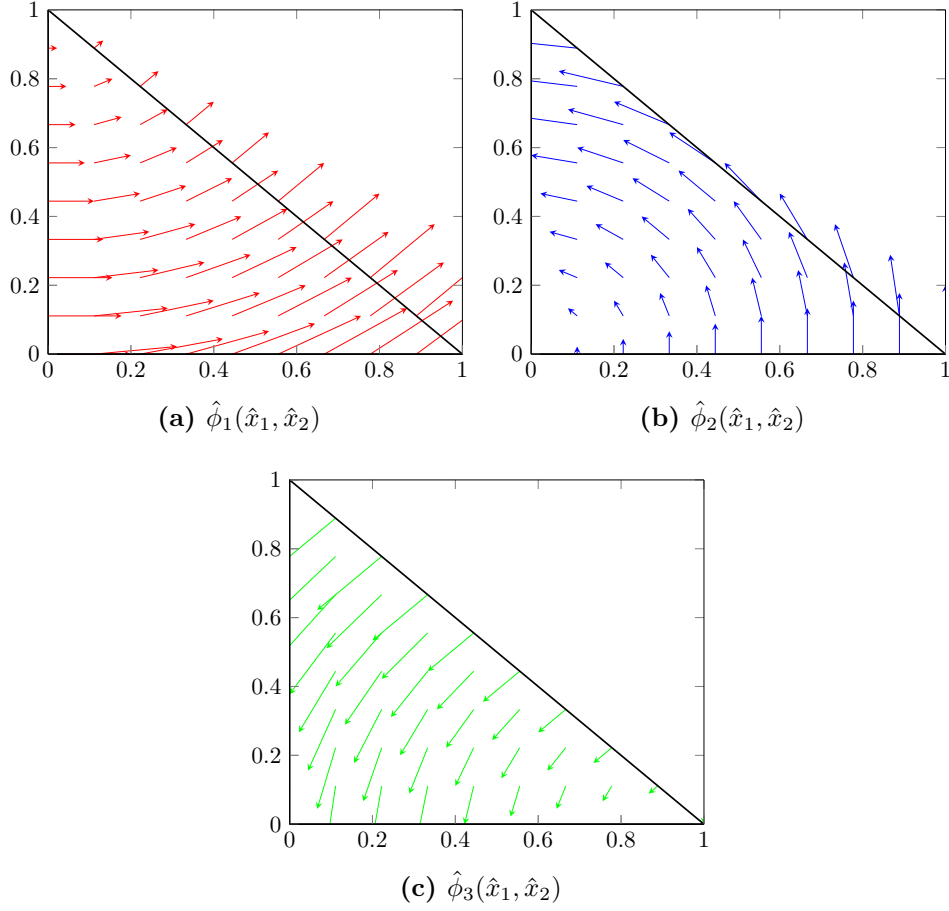
Then the triple  $(K, \mathcal{N}_0, \Sigma = \{\sigma_1, \dots, \sigma_s\})$  is a finite element which is called the Nédélec-element of first kind and lowest order ( $\mathcal{N}_0$ -element).

A proof for the unisolvence of Nédélec's elements can be found in [10, section 3]. As reference elements  $\hat{K}$  we use the unit triangle in 2D and the unit tetrahedron in 3D respectively. Together with the property  $\hat{\sigma}_i(\hat{p}_j) = \delta_{ij}$  for  $i, j = 1, 2, 3$  we obtain the local shape functions for the reference Nédélec-element  $(\hat{K}, \hat{\mathcal{N}}_0, \hat{\Sigma})$ .

In 2D these are

$$\hat{\phi}_1(\hat{x}_1, \hat{x}_2) = \begin{pmatrix} 1 - \hat{x}_2 \\ \hat{x}_1 \end{pmatrix}, \quad \hat{\phi}_2(\hat{x}_1, \hat{x}_2) = \begin{pmatrix} -\hat{x}_2 \\ \hat{x}_1 \end{pmatrix} \quad \text{and}$$

$$\hat{\phi}_3(\hat{x}_1, \hat{x}_2) = \begin{pmatrix} -\hat{x}_2 \\ \hat{x}_1 - 1 \end{pmatrix}.$$



**Figure 5.2:** Local shape functions for the reference  $\mathcal{N}_0$ -element in 2D.

In 3D for  $\hat{x} = (\hat{x}_1, \hat{x}_2, \hat{x}_3)^\top$  the basis functions are

$$\hat{\phi}_1(\hat{x}) = \begin{pmatrix} 1 - \hat{x}_2 - \hat{x}_3 \\ \hat{x}_1 \\ \hat{x}_1 \end{pmatrix}, \quad \hat{\phi}_2(\hat{x}) = \begin{pmatrix} \hat{x}_2 \\ 1 - \hat{x}_1 - \hat{x}_3 \\ \hat{x}_2 \end{pmatrix}, \quad \hat{\phi}_3(\hat{x}) = \begin{pmatrix} \hat{x}_2 \\ -\hat{x}_1 \\ 0 \end{pmatrix},$$

$$\hat{\phi}_4(\hat{x}) = \begin{pmatrix} \hat{x}_3 \\ \hat{x}_3 \\ 1 - \hat{x}_1 - \hat{x}_2 \end{pmatrix}, \quad \hat{\phi}_5(\hat{x}) = \begin{pmatrix} \hat{x}_3 \\ 0 \\ -\hat{x}_1 \end{pmatrix} \quad \text{and} \quad \hat{\phi}_6(\hat{x}) = \begin{pmatrix} 0 \\ -\hat{x}_3 \\ \hat{x}_2 \end{pmatrix}.$$

In order to obtain an arbitrary element in a triangulation  $\mathcal{T}$  we use the affine mapping  $T_K(\hat{x}) := B_K \hat{x} + b_K$  which maps the reference element  $\hat{K}$  to an arbitrary element  $K$  of the mesh of the domain  $\Omega$ . As the basis functions are vector fields we have to use a modified transformation to map the basis functions  $\hat{\phi}$  to a cell on the mesh in order to maintain the well-definedness of the local shape functions of each element. This transformation is given by

$$\phi(x) = B_K^{-\top} \hat{\phi} \circ T_K^{-1}(x) \quad (5.1)$$

and is called covariant Piola transform. It is the natural pullback map of a first order differential form.

The curl transforms differently according to the dimension as

$$2D : \quad \text{curl } \phi(x) = \frac{1}{\det(B_K)} \text{curl } \hat{\phi} \circ T_K^{-1}(x) \quad (5.2)$$

$$3D : \quad \text{curl } \phi(x) = \frac{1}{\det(B_K)} B_K \text{curl } \hat{\phi} \circ T_K^{-1}(x). \quad (5.3)$$

## 5.2 Global $H_{\text{curl}}$ Conformity, Orientation of the Edges and FEM Matrices

If we consider just a single element  $K$  of the mesh then a function which is in  $\mathcal{N}_0(K)$  clearly defines a function belonging to  $H_{\text{curl}}(K)$ . But how can we guarantee that this still remains true if we merge the elements together? For this purpose we consider a conforming triangulation  $\mathcal{T}_h$  of the domain  $\Omega$  and we define the space

$$\mathcal{R}_h(\Omega) = \{w \in L^1(\Omega) : w|_K \in \mathcal{N}_0, \forall K \in \mathcal{T}_h, \llbracket w \times n \rrbracket_f = 0, \forall \text{ facets } f\}$$

where  $\llbracket w \times n \rrbracket_f = 0$  denotes the jump of the tangential component of  $w$  along the inner facet  $f$ .

Then it holds that  $\mathcal{R}_h(\Omega) \subset H_{\text{curl}}(\Omega)$  and further the following theorem shows that we can achieve tangential continuity alongside facets by identifying the degrees of freedom on joint facets.

**Theorem 5.2.1.** *Let  $K_1$  and  $K_2$  be two elements such that  $f = K_1 \cap K_2$  is a common facet which both elements share. Then the condition  $\llbracket w \times n \rrbracket_f = 0$  is equivalent to*

$$\int_e w|_{K_1} \cdot t_e \, ds = \int_e w|_{K_2} \cdot t_e \, ds$$

for all edges  $e$  related to  $f$ .

As the degrees of freedom are integrals along edges we further need to ensure that an edge which belongs to multiple elements has the same orientation in a global setting, otherwise the tangential continuity will not be preserved.

Therefore we agree on the convention that an edge  $e_{ij}$  is oriented from the vertex with a smaller index to the vertex with a bigger index and we define the sign for an edge as

$$\text{sign}_K^{e_{ij}} = \begin{cases} +1 & \text{if } i < j \\ -1 & \text{otherwise} \end{cases}$$

on each element  $K$  which gives us a unique orientation for each edge. With this convention the global basis function can be obtained by just multiplying with the correct sign and we can compute the finite element matrices such as the mass matrix  $M_{ned}$  and the stiffness matrix  $A_{ned}$  as

$$\begin{aligned} 2D/3D : M_{ned}^K(m, n) &= |\det(B_K)| \int_{\hat{K}} \text{sign}_K^{e_m} B_K^{-\top} \hat{\phi}_m \cdot \text{sign}_K^{e_n} B_K^{-\top} \hat{\phi}_n \, d\hat{x} \\ 2D : A_{ned}^K(m, n) &= \frac{1}{|\det(B_K)|} \int_{\hat{K}} \text{sign}_K^{e_m} \text{curl } \hat{\phi}_m \cdot \text{sign}_K^{e_n} \text{curl } \hat{\phi}_n \, d\hat{x} \\ 3D : A_{ned}^K(m, n) &= \frac{1}{|\det(B_K)|} \int_{\hat{K}} \text{sign}_K^{e_m} B_K \text{curl } \hat{\phi}_m \cdot \text{sign}_K^{e_n} B_K \text{curl } \hat{\phi}_n \, d\hat{x} \end{aligned}$$

for  $K \in \mathcal{T}_h$ ,  $m, n = 1, 2, 3$  in 2D and  $m, n = 1, \dots, 6$  in 3D.

These entries can be computed exactly using second order gaussian quadrature formulas for the mass matrix and first order quadrature formulas for the stiffness matrix on the unit triangle or unit tetrahedron respectively. The entries of the mass matrix in 2D for example can thus be computed as

$$M_{ned}^K(m, n) = |\det(B_K)| \sum_{i=1}^3 w_i \text{sign}_K^{e_m} B_K^{-\top} \hat{\phi}_m(x_i) \cdot \text{sign}_K^{e_n} B_K^{-\top} \hat{\phi}_n(x_i)$$

with the weights  $w = \left(\frac{1}{6}, \frac{1}{6}, \frac{1}{6}\right)$  and the integration points

$$x_1 = \left(\frac{1}{6}, \frac{1}{6}\right), \quad x_2 = \left(\frac{2}{3}, \frac{1}{6}\right) \quad \text{and} \quad x_3 = \left(\frac{1}{6}, \frac{2}{3}\right).$$

We have implemented the edge elements in MATLAB following the ideas suggested in [13].

The connections between the mentioned finite elements, more precisely between the finite dimensional approximation spaces and the infinite dimensional counterparts, can be diagrammed using de Rham sequences as depicted in the following.

$$\begin{array}{ccccccc}
H^1 & \xrightarrow{\text{grad}} & H_{\text{curl}} & \xrightarrow{\text{curl}} & H_{\text{div}} & \xrightarrow{\text{div}} & L^2 \\
I_{\mathcal{T}}^{\mathcal{P}_1} \downarrow & & I_{\mathcal{T}}^{\mathcal{N}_0} \downarrow & & I_{\mathcal{T}}^{RT_0} \downarrow & & I_{\mathcal{T}}^{DG_0} \downarrow \\
\mathcal{P}_1 & \xrightarrow{\text{grad}} & \mathcal{N}_0 & \xrightarrow{\text{curl}} & RT_0 & \xrightarrow{\text{div}} & \mathcal{DG}_0
\end{array}$$

**Figure 5.3:** de Rham sequence

Here the space  $H_{\text{div}}$  is analogously defined as the space  $H_{\text{curl}}$ . Further  $RT_0$  denotes the Raviart-Thomas finite elements of lowest order and  $I_{\mathcal{T}}$  denotes the global interpolation operator of the corresponding finite element space

$$I_{\mathcal{T}} = \sum_i \sigma_i \phi_i.$$

### 5.3 The discrete Problems

#### The $u$ -problem

As mentioned before the  $u$ -problem boils down to solving a Poisson equation which, when employing the FEM method, can be done by solving a linear variational problem of the following form.

Find  $u \in V$  such that for all  $v \in V$

$$a[u, v] = F(v) \quad (5.4)$$

where  $V$  is some Hilberspace,  $F \in \mathcal{L}(V, \mathbb{R})$  and  $a : V \times V \rightarrow \mathbb{R}$  is a continuous and coercive bilinear form.

Instead of solving (5.4) we consider the discretized problem below.

Find  $u_h \in V_h$  such that for all  $v_h \in V_h$

$$a[u_h, v_h] = F(v_h) \quad (5.5)$$

where  $V_h$  is a finite-dimensional approximation space, usually a space of piecewise-defined polynomials. Since the space  $V_h$  is finite dimensional (5.5) holds if and only if (5.4) holds for a set of basis functions in  $V_h$ . This can be done by solving a linear system  $Au = F$ .

In our case we have  $V = H_0^1(\Omega)$  and we use  $\mathcal{P}_1$ -elements to approximate the solution  $u$ , i.e.  $V_h = I_{\mathcal{T}}^{\mathcal{P}_1}(C(\Omega))$  the image of the global interpolation operator  $I_{\mathcal{T}}^{\mathcal{P}_1}$ .

$$I_{\mathcal{T}}^{\mathcal{P}_1}(u) = \sum_{i=1}^{n_v} u(v_i) p_i(x) =: \sum_{i=1}^{n_v} u_i p_i(x) =: u_h(x)$$

where  $p_i(x)$ ,  $i = 1, \dots, n_v$  are the global shape functions of the  $\mathcal{P}_1$ -elements and  $n_v$  is the number of vertices of the triangulated domain. Furthermore we denote by  $u_c = [u_1, \dots, u_{n_v}]^\top$  the coefficient vector of  $u_h$ . Analogously  $d_c$  and  $b_c$  are the coefficient vectors of  $d_h$  and  $b_h$  which are defined by

$$d_h(x) = \sum_{k=1}^{n_W} d_k \phi_k(x), \quad b_h(x) = \sum_{k=1}^{n_W} b_k \phi_k(x) \quad (5.6)$$

where  $\phi_k(x)$ ,  $k = 1, \dots, n_W$  are the basis functions of the discrete approximation space  $Y_h$  for  $d$  depending on the choice of the corresponding finite element. For the choice of  $\mathcal{DG}_0$ -elements we have  $n_W = D n_{el}$ ,  $D = 2, 3$  with  $n_{el}$  the number of elements in the mesh and in the case of using the edge elements we have  $n_W = n_{edges}$ .

Further we denote by  $I_D \subset I = \{1, \dots, n_v\}$  the set of boundary nodes (Dirichlet nodes).

The linear variational problem (3.25) for the ADMM method then becomes

$$\begin{aligned} & (\mu + \gamma) \sum_{j=1}^{n_v} u_j \int_{\Omega} (\nabla p_i, \nabla p_j) dx \\ &= \int_{\Omega} f p_i dx + \gamma \sum_{k=1}^{n_W} (d_k + b_k) \int_{\Omega} (\nabla p_i, \phi_k) dx, \quad i = 1, \dots, n_v \end{aligned} \quad (5.7)$$

The equations (5.7) only have to hold for inner vertices of the mesh, i.e. only for non-Dirichlet dofs  $(i, j) \in (I \setminus I_D) \times (I \setminus I_D)$ .

Let us denote by  $A_{\mathcal{P}_1}$  the stiffness matrix for the  $\mathcal{P}_1$ -elements, the entries of which are defined by  $A_{\mathcal{P}_1}(i, j) = \int_{\Omega} (\nabla p_i, \nabla p_j) dx$ , further let

$$L_i = \int_{\Omega} f p_i dx, \quad i = 1, \dots, n_v$$

be the load vector. Finally let the rectangular matrix  $B$  be defined by

$$B = (b_{ij})_{i,j=1}^{n_v, n_W}, \quad b_{ij} = \int_{\Omega} (\nabla p_i, \phi_j) dx$$

then (5.7) simplifies to

$$(\mu + \gamma) A_{\mathcal{P}_1} u_c = L + \gamma B (d_c + b_c). \quad (5.8)$$

Again we stress the fact that the equations in (5.8) only have to be satisfied for non-boundary dofs. Since we infer homogenous Dirichlet boundary conditions we can simply set the remaining boundary dofs to zero  $u_j = 0$ ,  $j \in I_D$ . Analogously the  $u$ -problem for Chambolle and Pock's algorithm can be solved.



**The  $d$ -problem for the Choice  $Y = L^2(\Omega)$**

We first discuss the choice of the space  $Y = L^2(\Omega)$ . For this choice we will use discontinuous Lagrange elements of lowest order for the discretization. Since we have no global continuity conditions or similar requirements, there is no coupling between the elements or their local coefficient vectors. Thus it is convenient to represent  $d_h$  with respect to the local shape functions

$$d_h(x) = \sum_{K \in \mathcal{T}} \sum_{i=1}^2 d_{i,K} \phi_{i,K}(x) \quad (5.9)$$

where  $d_{i,K} = [d(x_{T,K})]_i$ ,  $x_{T,K}$  is the midpoint of the element  $K$  and  $[\cdot]_i$  denotes the  $i$ -th component of a vector. Further let  $d_{c,K} = (d_{i,K})_{i=1}^D$ .

Since  $u$  is a piecewise linear function we know that  $\nabla u_h$  is a piecewise-constant function and so is  $z_h := \nabla u_h - b_h$ . For  $\nabla u_h$  we obtain on each element  $K$

$$\nabla u_h|_K = B_K^{-\top} (u_{K,1} \nabla \hat{p}_{1,K} + u_{K,2} \nabla \hat{p}_{2,K} + u_{K,3} \nabla \hat{p}_{3,K}) \quad (5.10)$$

$$= B_K^{-\top} \left( u_{K,1} \begin{pmatrix} -1 \\ -1 \end{pmatrix} + u_{K,2} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + u_{K,3} \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right) = B_K^{-\top} \begin{pmatrix} u_{K,2} - u_{K,1} \\ u_{K,3} - u_{K,1} \end{pmatrix}. \quad (5.11)$$

with the local coordinates  $u_{i,K}$  on the element  $K$ .

The discretized problem associated with (3.26) reads

$$q_h(x) \in \partial d_h(x) = \begin{cases} \frac{d_h(x)}{|d_h(x)|_2} & : |d_h(x)|_2 > 0 \\ \{y \in \mathbb{R}^D : |y|_2 \leq 1\} & : \text{otherwise} \end{cases} \quad \text{a.e.}$$

$$\sum_{K \in \mathcal{T}} \int_K \beta q_h(x) \cdot w_h(x) dx + \gamma (d_h(x) - z_h(x)) \cdot w_h(x) dx = 0 \quad (5.12)$$

for all  $w_h \in Y_h$ . We define the semi-norm

$$\sqrt{d_{c,K}^\top M_K(x) d_{c,K}} = |d_{c,K}|_{M_K(x)}$$

for the local gramian matrix

$$M_K(x) = \left( \phi_{i,K}(x)^\top \phi_{j,K}(x) \right)_{i,j=1}^D, \quad x \in K$$

with the local shape functions  $\phi_{i,K}$  of the  $\mathcal{DG}_0$ -elements. Using (5.9) the problem (5.12) then reads

$$q_K(x) \in \partial d_h(x) = \begin{cases} \frac{d_{c,K}}{|d_{c,K}|_{M_K(x)}} & : |d_{c,K}|_{M_K(x)} > 0 \\ \{y \in \mathbb{R}^D : |y|_2 \leq 1\} & : \text{otherwise} \end{cases} \quad \text{a.e.}$$

$$\sum_{K \in \mathcal{T}} \int_K \beta M_K(x) q_K(x) + \gamma M_K(x) (d_{c,K} - z_{c,K}) dx = 0. \quad (5.13)$$

Since we use  $\mathcal{DG}_0$ -elements the corresponding mass matrix is a diagonal matrix with the diagonal entries the areas of the elements and the local gramian matrix  $M_K(x)$  is constant and just a  $D \times D$  - identity matrix,  $D = 2, 3$ . Thus on each element  $K$  we have the condition

$$\beta|K| \begin{cases} \frac{d_{c,K}}{|d_{c,K}|^2} & : |d_{c,K}|_2 > 0 \\ \{y \in \mathbb{R}^D : |y|_2 \leq 1\} & : \text{otherwise} \end{cases} + \gamma|K|(d_{c,K} - z_{c,K}) = 0$$

This leads to the shrinkage operator

$$d_{c,K} = \text{shrink} \left( z_{c,K}, \frac{\beta}{\gamma} \right).$$

We further need to determine the distance from the current iterate to the old iterate in each step. We can calculate this quantity using the mass matrix for the  $\mathcal{P}_1$ -elements  $M_{\mathcal{P}_1}$ , with the entries  $M_{\mathcal{P}_1}(i, j) = \int_{\Omega} p_i p_j dx$  and obtain

$$\text{dist}_{k,h} := \|u_{h,k+1} - u_{h,k}\|_{L^2(\Omega)} = \left[ (u_c^{k+1} - u_c^k)^\top M_{\mathcal{P}_1} (u_c^{k+1} - u_c^k) \right]^{\frac{1}{2}}. \quad (5.14)$$

We summarize the procedure in the following algorithm.

---

**Algorithm 4:** ADMM method for (3.22)

---

**Input** :  $\Omega, \beta, \mu, \gamma, f$

**Output:** Coefficient vectors  $u_c^k, d_c^k, b_c^k$

1 Set  $u_c^0 = \mathbf{0}, d_c^0 = \mathbf{0}, b_c^0 = \mathbf{0}$  and  $k = 0$

2 ***u*-subproblem**

Solve the system

$$(\mu + \gamma)A_{\mathcal{P}_1}u_c^{k+1} = L + \gamma B(d_c^k + b_c^k)$$

for  $u_c^{k+1}$ .

3 ***d*-subproblem**

For each element  $K$  set  $d_{c,K}^{k+1} = \text{shrink} \left( \nabla u_h|_K^{k+1} - b_{c,K}^k, \frac{\beta}{\gamma} \right)$ .

4 ***b*-update**

For each element  $K$  set  $b_{c,K}^{k+1} = b_{c,K}^k - \nabla u_h|_K^{k+1} + d_{c,K}^{k+1}$ .

5 Convergence test, e.g. check if  $\text{dist}_k$  is larger than a given tolerance.

6 Set  $k \leftarrow k + 1$  and go back to 2.

---

We have to solve a system with the symmetric and positive definite stiffness matrix  $A_{\mathcal{P}_1}$  in each step. Hence for moderate system sizes, with approximately up to one million unknowns, an efficient way to solve the  $u$ -subproblem is to compute a Cholesky factorization of the stiffness matrix

first. Then the solution of the  $u$ -problem is obtained by a forward and a backward substitution in each ADMM iteration. For larger system sizes iterative Krylov subspace methods such as GMRES with an incomplete LU factorization as preconditioners should be used. Since we are here in the setting of a symmetric and positive definite matrix we use a CG method preconditioned by an incomplete Cholesky factorization.

### The $d$ -problem for choosing Nédélec's Edge Elements

Instead of using  $\mathcal{DG}_0$ -elements to solve the  $d$ -problem we can also use Nédélec's edge-elements to discretize  $d$ . Ideally we would like to have  $d_h^* = \nabla u_h^*$ , with  $u_h^*$  the discrete solution of the TV-denoising problem. Thus to obtain convergence we have to verify that  $\nabla u_h$  can be represented as a Nédélec function. We know that on a common edge  $e = a + t(b - a)$  of two elements the function  $u_h \in V_h$  is continuously differentiable for all  $t \in (0, 1)$ .

Hence we know that  $\nabla u_h$  has a continuous tangential component along all the edges which means that the jump of the tangential component of  $\nabla u_h$  between internal edges is zero. Thus  $\nabla u_h \in \mathcal{R}_h(\Omega)$  and in particular  $\nabla u_h$  has a representation as a Nédélec function.

If we take a look at a function  $n_h \in \mathcal{N}_0$  we get

$$n_h|_K = \sum_{i=1}^3 n_{i,K} B_K^{-\top} \hat{\phi}_{i,K}(\hat{x}) = B_K^{-\top} \begin{pmatrix} n_{1,K} - \hat{x}_2(n_{1,K} + n_{2,K} + n_{3,K}) \\ -n_{3,K} + \hat{x}_1(n_{1,K} + n_{2,K} + n_{3,K}) \end{pmatrix},$$

with  $\hat{x} = T_K^{-1}x$ .

This means we need to require that the coefficients locally sum up to zero in order to get a constant function and by comparing coefficients with (5.11) we have that  $n_{1,K} = u_{2,K} - u_{1,K}$ ,  $n_{3,K} = u_{1,K} - u_{3,K}$  and  $n_{2,K} = u_{3,K} - u_{2,K}$  are the local coefficients of  $\nabla u$  in the Nédélec basis representation. We can then map these local coordinates to the global coordinates using the local-to-global dof mapping.

By (3.28) we know that in the continuous case the optimal solution to the  $d$ -problem is a shrinkage operator. Thus we suggest to use the best  $L^2$ -approximation of the shrinkage operator within the approximation space of Nédélec functions of lowest order, i.e. we need to solve

$$\arg \min_{d_h \in \mathcal{N}_0} \frac{1}{2} \left\| d_h - \text{shrink} \left( \nabla u_h - b_h, \frac{\beta}{\gamma} \right) \right\|_{L^2(\Omega)}^2 \quad (5.15)$$

where  $d_h$  has the representation (5.6).

This leads to solving the linear variational problem

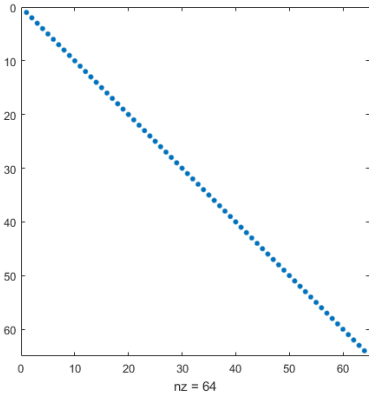
$$\begin{aligned} & \sum_{j=1}^{n_{\text{edges}}} d_j \int_{\Omega} (\phi_i, \phi_j) \, dx \\ &= \underbrace{\int_{\Omega} \left( \text{shrink} \left( \nabla u_h - b_h, \frac{\beta}{\gamma} \right), \phi_i \right) \, dx}_{:= [L_d]_i}, \quad i = 1, \dots, n_{\text{edges}} \end{aligned} \quad (5.16)$$

which is equivalent to solving

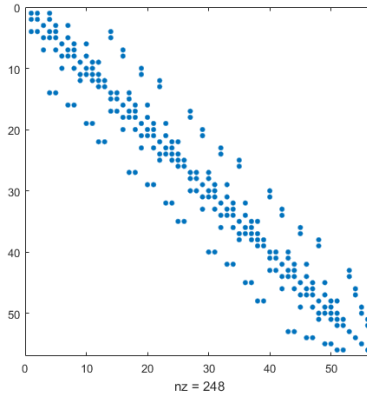
$$M_{ned} d_c = L_d \quad (5.17)$$

with  $d_c = (d_1, \dots, d_{n_{\text{edges}}})$ .

If we replace  $M_{ned}$  with  $M_{\mathcal{DG}_0}$  in (5.17) this approach leads to the same solution as for the choice  $Y_h = \mathcal{DG}_0$ . However in this case the mass matrix is no diagonal matrix and we have to solve the linear system (5.17), which is again best solved by employing a Cholesky decomposition.



(a)  $M_{\mathcal{DG}_0}$



(b)  $M_{ned}$

**Figure 5.4:** Mass matrices for the  $\mathcal{DG}_0$  and Nédélec elements of lowest order for conforming and regular mesh of the unit square with 32 elements. The entries of  $M_{\mathcal{DG}_0}$  are the areas of the elements. For the Nédélec elements each row contains three or five non-zero entries depending on whether the corresponding edge of the global shape function is a boundary edge or an interior edge.

**The choice**  $Y = H_{curl}(\Omega)$

Now we discuss the choice  $Y = H_{curl}(\Omega)$  which we will also handle using Nédélec's edge-elements. For this approach we further apply Newton's method to the discrete analogue of problem (3.29)

$$F(d_c) = 0$$

with  $F$  defined as

$$\begin{aligned} F(d_1, \dots, d_{n_{\text{edges}}}) &= \sum_{j=1}^{n_{\text{edges}}} (d_j - z_j) \int_{\Omega} (\phi_i, \phi_j) dx \\ &+ \sum_{j=1}^{n_{\text{edges}}} (d_j + b_j) \int_{\Omega} (\text{curl } \phi_i, \text{curl } \phi_j) dx \\ &+ \frac{\beta}{\gamma} \sum_{j=1}^{n_{\text{edges}}} d_j \int_{\Omega} \frac{(\phi_i, \phi_j)}{\sqrt{|\sum_{j=1}^{n_{\text{edges}}} d_j \phi_j|_2^2 + \varepsilon}} dx, \quad i = 1, \dots, n_{\text{edges}} \end{aligned} \quad (5.18)$$

Because of the regularization  $F$  is differentiable with respect to the  $d_i$ . We define the matrix

$$G(d_c^k) = \left( g_{ij}(d_c^k) \right)_{i,j=1}^{n_{\text{edges}}}$$

which is the Jacobian matrix of the total variation term in (5.18) and depends on the coefficient vector  $d_c^k$  of the  $k$ -th approximative iterate  $d_h^k$ . In 2D we have

$$d_h^k(x) = \sum_{i=1}^{n_{\text{edges}}} d_i^k \phi_i(x) =: \begin{pmatrix} d_{h,1}^k(x) \\ d_{h,2}^k(x) \end{pmatrix}$$

and the entry  $g_{ij}(d_c^k)$  is given by

$$\begin{aligned} g_{ij}(d_c^k) &= \int_{\Omega} \frac{1}{\sqrt{d_{h,1}^k(x)^2 + d_{h,2}^k(x)^2 + \varepsilon}^3} \cdot \\ &\phi_i(x)^\top \begin{pmatrix} d_{h,2}^k(x)^2 + \varepsilon & -d_{h,1}^k(x) \cdot d_{h,2}^k(x) \\ -d_{h,1}^k(x) \cdot d_{h,2}^k(x) & d_{h,1}^k(x)^2 + \varepsilon \end{pmatrix} \phi_j(x) dx \end{aligned} \quad (5.19)$$

We approximate these entries using a second order gaussian cubature rule as for the mass matrices.

Once we have these entries we can write the Jacobian  $J_F$  of  $F$  as

$$J_F(d_c^k) = \frac{\beta}{\gamma} G(d_c^k) + M_{ned} + A_{ned}$$

with the mass matrix  $M_{ned}$  and the stiffness matrix  $A_{ned}$ . Then a single Newton step is defined by

$$\begin{aligned} \text{Solve the system :} \quad & J_F(d_c^k) \Delta d = -F(d_c^k) \\ \text{Update :} \quad & d_c^{k+1} = d_c^k + \eta \Delta d. \end{aligned}$$

For our purpose we will mainly use full Newton steps, i.e.  $\eta = 1$ . Furthermore it is known, c.f. [2, section 3.4.4], that in order for the ADMM method to converge it is not necessary to solve the subproblems to full precision. Hence we will perform only one Newton step in each iteration, means we only have to solve one additional system per each iteration.

As the quality of the regularization model depends on the magnitude of the parameter  $\varepsilon$  we further tighten this parameter if the solution is sufficiently well approximated, e.g.

$$\varepsilon_{k+1} = \begin{cases} \tau \varepsilon_k & \text{if } \tilde{\varepsilon}_1 < \|u_{h,k+1} - u_{h,k}\|_{L^2(\Omega)} < \tilde{\varepsilon}_2 \\ \varepsilon_k & \text{otherwise} \end{cases} \quad (5.20)$$

with  $\tau \in (0, 1)$ .

## Chapter 6

# Numerical Results

### 6.1 Two-dimensional Flow Problem

We consider minimizing the functional  $E$ , given by

$$\frac{\mu}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \int_{\Omega} f u \, dx + \beta \int_{\Omega} |\nabla u|_2 \, dx \quad (6.1)$$

in 2d on the unit square, i.e.  $\Omega = (0, 1)^2$  using the methods from the previous sections. For the source term  $f$  we choose

$$f(x_1, x_2) = 10 \sin(3\pi x_1) \cos(2\pi x_2) (2 - 0.5x_1x_2).$$

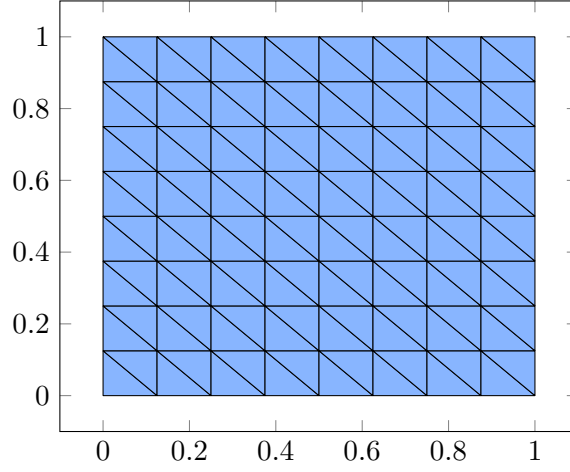
Further we set  $\mu = 1$  and  $\beta = 0.8$ .

For the domain  $\Omega$  we set up a geometrically conforming triangulation, i.e. no hanging vertices, by bisecting  $n_{cells}$  squares of length  $n_{cells}^{-1}$  in  $x$  and  $y$  direction. We continue this procedure until the unit square is triangulated.

The resulting mesh then has a total of  $(n_{cells} + 1)^2$  vertices and a total of  $2n_{cells}^2$  elements. To calculate the number of edges we can use the fact that the Euler characteristic of any two triangulations of a surface with the same boundary is equal, and in the case of the unit square equal to one. And thus by Euler's formula

$$V - E + F = 1$$

where  $V$  denotes the number of vertices,  $E$  the number of edges and  $F$  the number of facets. This results in a total of  $3n_{cells}^2 + 2n_{cells}$  edges.



**Figure 6.1:** Geometrically conforming triangulation of the unit square for  $n_{cells} = 8$ , i.e. 128 elements.

We solve 6.1 using the ADMM method for  $Y = L^2(\Omega)$  and  $Y = H_{\text{curl}}(\Omega)$  with the methods derived in section 3.5, 3.5 and 5.3 respectively.

To measure the quality of the approximation for the ADMM method we will use the distance between the current and last iterate  $\text{dist}_k$  and as another approach the error between the current iterate and a sufficiently well approximated solution, i.e.

$$e_{k,h} := \|u_{k,h} - u_{N,h}\|_{L^2(\Omega)} \quad (6.2)$$

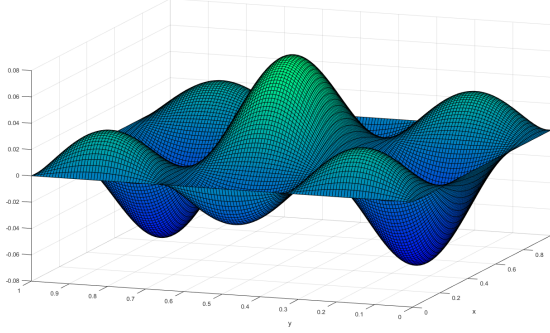
where  $u_{k,h}$  denotes the discretized solution of the  $u$ -subproblem after iteration  $k$ . That means  $e_{k,h}$  is a fairly good measure for the error of the ADMM method as a part of the total error which consists of the discretization error and the ADMM error, i.e.

$$\|u - u_k\|_{L^2(\Omega)} \leq \|u - u_h\|_{L^2(\Omega)} + \|u_h - u_{k,h}\|_{L^2(\Omega)} \quad (6.3)$$

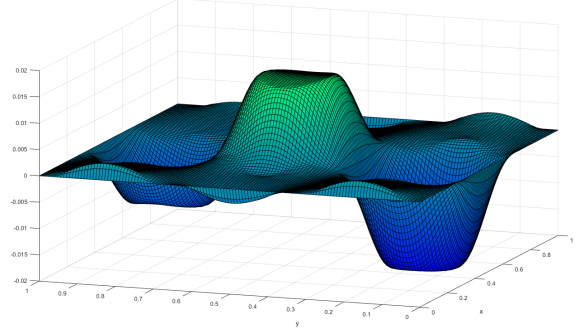
where  $u_h$  denotes the true solution of the discretized problem.

We will use  $N = 5000$  and compute iterates from  $k = 1, \dots, 200$ .

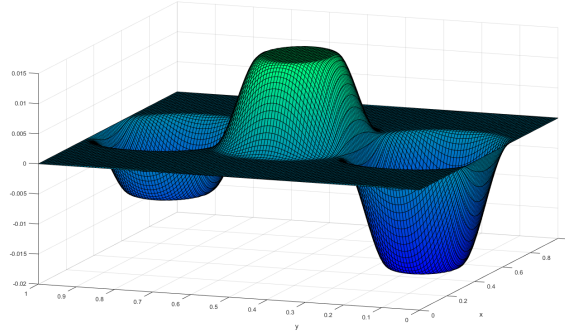




(a) solution of the poisson problem



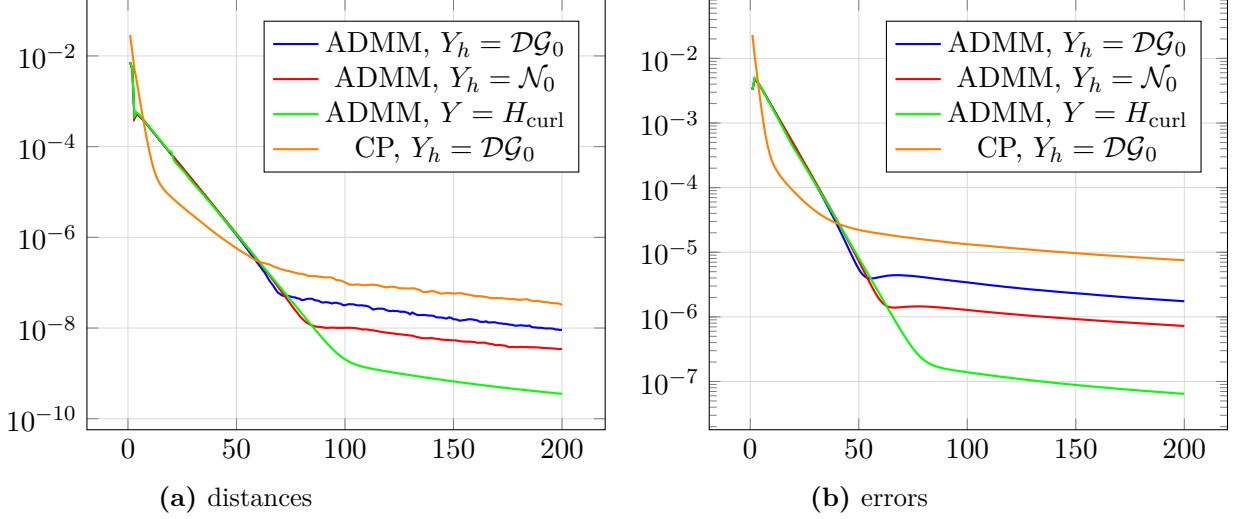
(b) solution after 8 iteration



(c) solution after 50 iterations

**Figure 6.2:** Solution of problem (6.1) after some iterations of the ADMM method ( $\gamma = 1$ ) and  $n_{cells} = 256$ .

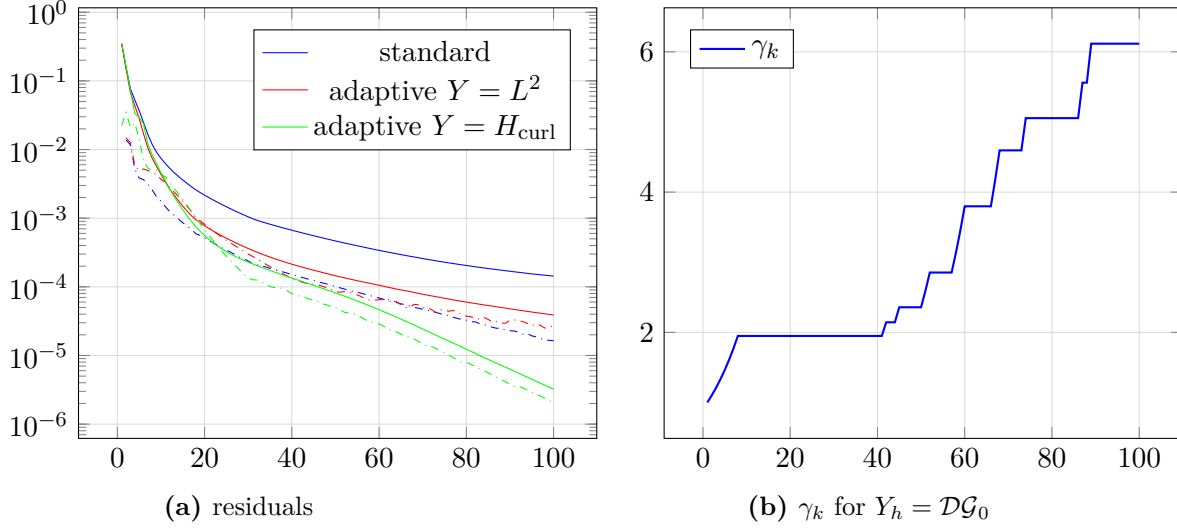
In Figure 6.2 we see that we only need very few ADMM iterations to arrive at a fairly denoised solution. Even after just a few iteration we can see a clear denoising effect. This is a general advantage of the ADMM method when applied to  $L^1$ -regularization problems as we stated in chapter 3.



**Figure 6.3:** Comparison of convergence speed measured by distance  $\text{dist}_{k,h}$  defined by (5.14) and error  $\text{err}_{k,h}$  defined by (6.2) for the different methods with constant penalty parameter  $\gamma = 7$ .

As Figure 6.3 shows, Chambolle and Pock’s method converges faster than the ADMM method at the start but flattens out very similar to the ADMM method afterwards. We further see that choosing  $Y_h = \mathcal{N}_0$  as the discretization space for the  $d$ -problem yields better convergence results by almost one order. This effect is even clearer when choosing the space  $Y = H_{\text{curl}}$  and penalizing the functional with the stronger  $H_{\text{curl}}$  norm.

However, to obtain convergence to the true solution we also need to refine the mesh, i.e. let the cell diameter  $h \rightarrow 0$ . Thus in this context if the error contribution of the discretization error in (6.3) is relatively large, the advantages of our method depicted above may not be so clear and the methods might show very similar convergence properties.



**Figure 6.4:** Primal and dual residuals for the iterates of the ADMM method when using the standard choice  $\gamma = 1$  for the penalty parameter or using an adaptive penalty parameter update scheme.

In Figure 6.4a we see that for the standard choice  $\gamma = 1$  for the penalty parameter we have a large gap between the primal and dual residual. Using an adaptive penalty parameter with the update rule (3.20) leads to a smaller primal residual as well as a highly reduced gap between the residuals. Using the  $H_{\text{curl}}$  norm for the penalty yields even better results. Further we see in Figure 6.4b that the residuals are quite sensitive to the choice of the penalty parameter, thus the step size  $\tau$  for updating the penalty parameter should not be too large.

## 6.2 Image denoising Problem

We now consider an image denoising problem and use the ROF model (2.4). This can be done analogously as described in section 5.3, i.e. we approximate the picture  $u$  as a  $\mathcal{P}_1$ -function. Alternatively one can also use a discontinuous approach and assume that  $u \in \mathcal{DG}_r$ ,  $r \geq 0$ , which leads to nonsymmetric interior penalty Galerkin (NIPG) method for the  $u$ -problem, as described in [6, section 5.1].

As a reference picture we use the famous picture of Lena, which has become a standard test object in image denoising.

To measure the quality of the reconstructed image  $A_{\text{rec}} \in \mathbb{R}^{m,n}$  from the noisy image we use the peak signal to noise ratio (PSNR), which is a standard tool for measuring the quality of compressed images in relation to the original image  $A \in \mathbb{R}^{m,n}$ .

PSNR is defined as

$$\text{PSNR} = 10 \log_{10} \left( \frac{mn \cdot \text{MAX}_P^2}{\|A - A_{\text{rec}}\|_F^2} \right) \quad (6.4)$$

where  $\|A\|_F$  denotes the Frobenius norm of  $A$  and  $\text{MAX}_P$  denotes the maximal pixel value of the image, e.g. 255 for a color depth of 8 bits. A higher PSNR value usually indicates a better reconstructed picture.



(a) true image  
512 × 512 pixels



(b) noisy image  
PSNR = 21.07



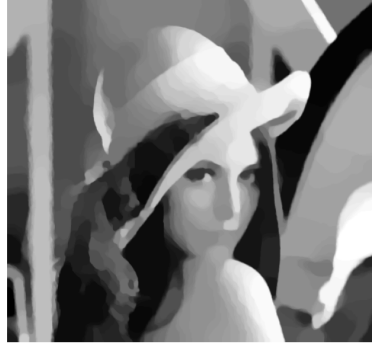
(c) 10 iterations  
PSNR = 28.34



(d) 500 iterations  
PSNR = 30.2

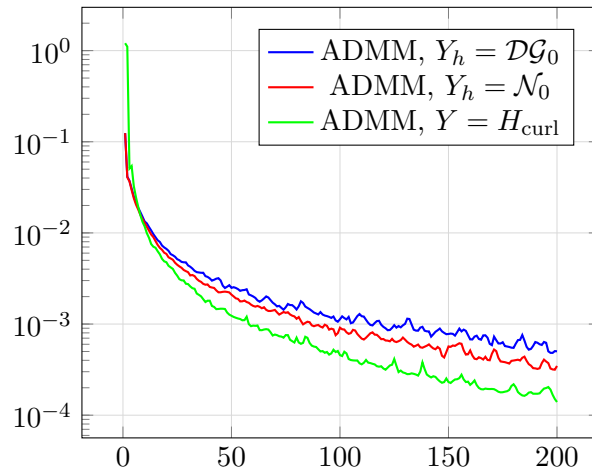
**Figure 6.5:** Comparison of original, noisy and denoised pictures after different numbers of iterations, denoising parameter  $\beta = 1.25 \cdot 10^{-4}$ .

Figure 6.5 confirms the observation that one only needs a small number of iterations to achieve reasonable denoising results.



**Figure 6.6:** Resulting image for  $\beta$  chosen too large.

The choice of the denoising parameter  $\beta$  is of significant importance for the resulting picture. For small values of  $\beta$  the resulting image might still be too noisy as sudden jumps in the pixels are not penalized harshly enough, whereas for large values of  $\beta$  sharp corners of the original will be lost and the image becomes cartoon-like, as can be seen in Figure 6.6.



**Figure 6.7:** Primal residuals (3.14) for the ROF model on a grid with  $256 \times 256$  cells,  $\beta = 1.25 \cdot 10^{-4}$  and  $\gamma = 10^{-3}$ .

In Figure 6.7 we can observe that we achieve faster convergence results in terms of the primal residuals when employing Nédélec elements. The other convergence indicators behave very similarly among the methods for this example.

## Chapter 7

# Summary and Future Work

### 7.1 Summary

In this work we have presented new strategies for discretizing subproblems in operator splitting methods, such as the ADMM method, that possess superior convergence properties when compared to the standard choice for discretization. Since they do however also require more computational cost, it depends on further factors whether it is necessary to spend this additional effort. For example if one does not need high precision results the usual discretization approach should be totally sufficient and efficient for this matter. However, as the the ADMM method is known to converge slowly in the usual setting, the proposed methods can be a remedy to this shortcoming, if one needs more precise results. In this regard though, it has to be mentioned that the regularization approach we use to solve the problem when penalizing with the  $H_{\text{curl}}$ -norm also incorporates an error contribution since it artificially smoothes the solution.

### 7.2 Future Work

Further work is necessary to better understand the convergence behaviour of the named fixed point iterative methods. Thus for example, it is generally not clear which implications one can make on the error with respect to the exact solution, when a stopping criteria has been triggered because a certain tolerance has been reached.

Thus the advantages of our methods, comparatively, might not be so clear when we incorporate the discretization error.

Many results for the aforementioned methods also are only available in a finite dimensional setting, thus a great achievement would be to find similar results for infinite dimensional problems.

# References

- [1] Viorel Barbu. *Nonlinear differential equations of monotone types in Banach spaces*. Springer Monographs in Mathematics. Springer, New York, 2010, pp. x+272. URL: <https://doi.org/10.1007/978-1-4419-5542-5> (cit. on p. 3).
- [2] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers”. In: *Found. Trends Mach. Learn.* 3.1 (Jan. 2011), pp. 1–122. URL: <http://dx.doi.org/10.1561/22000000016> (cit. on p. 43).
- [3] K. Bredies and D. Lorenz. *Mathematische Bildverarbeitung: Einführung in Grundlagen und moderne Theorie*. Vieweg+Teubner Verlag, 2010. URL: <http://books.google.de/books?id=0FVBkXCK3q8C> (cit. on pp. 5, 9).
- [4] Antonin Chambolle and Thomas 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. URL: <http://dblp.uni-trier.de/db/journals/jmiv/jmiv40.html#ChambolleP11> (cit. on pp. 26, 30).
- [5] Tom Goldstein and Stanley Osher. “The Split Bregman Method for L1-Regularized Problems”. In: *SIAM J. Imaging Sciences* 2.2 (2009), pp. 323–343. URL: <https://doi.org/10.1137/080725891> (cit. on pp. 14, 18).
- [6] M. Herrmann, R. Herzog, S. Schmidt, J. Vidal-Núñez, and G. Wachsmuth. “Discrete Total Variation with Finite Elements and Applications to Imaging”. In: *arXiv e-prints* (Apr. 2018). arXiv: [1804.07477](https://arxiv.org/abs/1804.07477) [math.NA] (cit. on p. 48).
- [7] Fumiaki Kohsaka and Wataru Takahashi. “Existence and Approximation of Fixed Points of Firmly Nonexpansive-Type Mappings in Banach Spaces”. In: *SIAM J. on Optimization* 19.2 (Aug. 2008), pp. 824–835. URL: <http://dx.doi.org/10.1137/070688717> (cit. on p. 6).

- [8] P.P. Mosolov and V.P. Miashikov. “On stagnant flow regions of a viscous-plastic medium in pipes”. In: *Journal of Applied Mathematics and Mechanics* 30.4 (1966), pp. 841–854. URL: <http://www.sciencedirect.com/science/article/pii/0021892866900359> (cit. on p. 8).
- [9] P.P. Mosolov and V.P. Myasnikov. “Variational methods in the theory of the fluidity of a viscous-plastic medium”. In: *PMM, Journal of Applied Mathematics and Mechanics* 29 (Dec. 1965) (cit. on p. 8).
- [10] J. C. Nédélec. “A new family of mixed finite elements in  $\mathbb{R}^3$ ”. In: *Numerische Mathematik* 50.1 (1986), pp. 57–81. URL: <https://doi.org/10.1007/BF01389668> (cit. on p. 32).
- [11] Stanley Osher, Martin Burger, Donald Goldfarb, Jinjun Xu, and Wotao Yin. “An iterative regularization method for total variation-based image restoration”. In: *MULTISCALE MODEL. SIMUL.* 4.2 (2005), pp. 460–489 (cit. on p. 15).
- [12] R. Tyrrell Rockafellar. *Convex analysis*. Princeton Mathematical Series. Princeton, N. J.: Princeton University Press, 1970 (cit. on p. 26).
- [13] Jan Valdman and Immanuel Anjam. “Fast MATLAB assembly of FEM matrices in 2D and 3D: Edge elements”. In: *CoRR* abs/1409.4618 (2014). arXiv: 1409.4618. URL: <http://arxiv.org/abs/1409.4618> (cit. on p. 35).
- [14] Brendt Wohlberg. “ADMM Penalty Parameter Selection by Residual Balancing”. In: *CoRR* abs/1704.06209 (2017) (cit. on p. 20).