

# A Kačanov Type Iteration for the $p$ -Poisson Problem

Dissertation  
zur Erlangung des Doktorgrades (Dr.rer.nat.)  
des Fachbereichs Mathematik/Informatik  
der Universität Osnabrück

vorgelegt  
von  
Maximilian Wank

aus  
München

Osnabrück, 2016/2017

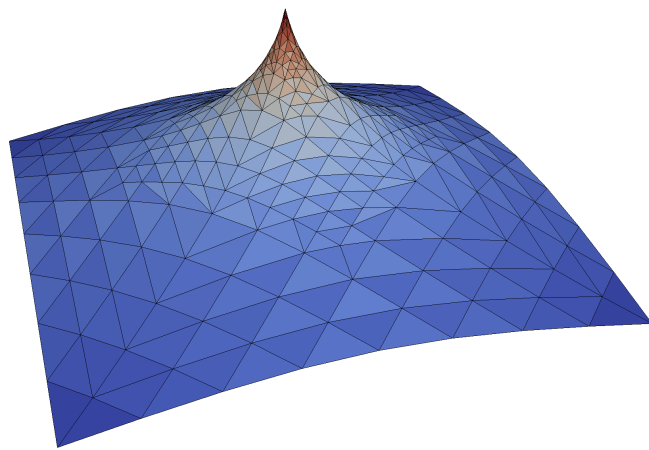


---

# A Kačanov Type Iteration for the $p$ -Poisson Problem

Maximilian Wank

---



Osnabrück 2016



---

# **A Kačanov Type Iteration for the $p$ -Poisson Problem**

**Maximilian Wank**

---

Dissertation  
am Institut für Mathematik  
der Universität Osnabrück

vorgelegt von  
Maximilian Wank  
aus München

Osnabrück, den 27. Oktober

Betreuer: Prof. Dr. Lars Dienen

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Outline . . . . .	1
1.2	Notation . . . . .	2
<b>2</b>	<b>Setting</b>	<b>3</b>
2.1	The $p$ -Poisson Equation and its Solution Space . . . . .	3
2.2	Further Function Spaces . . . . .	9
2.3	Presentation of the Algorithm . . . . .	18
<b>3</b>	<b>Convergence in the Relaxation Parameter</b>	<b>39</b>
3.1	$\Gamma$ -Convergence . . . . .	42
3.2	Error Bounds . . . . .	46
<b>4</b>	<b>The Kačanov Iteration</b>	<b>51</b>
4.1	Exponential Convergence . . . . .	51
4.2	An Example . . . . .	54
<b>5</b>	<b>Overall Convergence Analysis</b>	<b>61</b>
5.1	An Algebraic Rate . . . . .	61
5.2	Outlook On Adaptive Strategies for the Relaxation Parameter . . . . .	66
5.3	Numerical Examples . . . . .	74





# 1 Introduction

## 1.1 Outline

We start Section 2 by recapitulating the linear Poisson equation. However when describing non-Newtonian fluids, the viscosity may depend on the shear rate. In the model of those fluids, a generalization of the Laplace operator is needed. This is done by the non-linear  $p$ -Laplacian, which will be introduced in the further discussion of Section 2.1. After that we present the solution spaces for the simplest form of an equation containing the  $p$ -Laplace operator. Following, we use the structure of the problem to deduce our Kačanov based algorithm for the approximation of solution of the  $p$ -Poisson equation. Note that this approach as already been discussed for example in [Wei95] and [HJS97]. However, these results are not applicable to the  $p$ -Poisson problem.

In Chapter 3 we discuss a necessary relaxation that is introduced in Chapter 2. This relaxation is based a relaxation interval  $\varepsilon = (\varepsilon_-, \varepsilon_+)$  for two parameters  $0 < \varepsilon_- \leq 1 \leq \varepsilon_+ < \infty$ . We are mostly interested in the behaviour of  $\varepsilon_- \rightarrow 0$  and  $\varepsilon_+ \rightarrow \infty$  and show convergence results of the solutions of the relaxed problems to the solution of the original problem. Mainly we discuss two different approaches. At first we use the notion of  $\Gamma$ -convergence in Section 3.1 to deduce the required results. This techniques advantage is that it is able to be used for the most general case of the  $p$ -Poisson equation. However, no quantitative results will be available in this generality. In Section 3.2 we present estimates for the relaxation error introduced by  $\varepsilon$ . To do so, we will need to restrict ourselves to some classes of special cases of the equation that are known to provide certain regularity of the solutions of the original problem.

After discussing the error introduced by the relaxation we fix the relaxation parameters in Chapter 4 and discuss the error decay of the iterative part of our algorithm introduced in Chapter 2, where in each iteration step a linear equation needs to be solved. We will give an estimate yielding convergence at the cost of rates that depend badly on the relaxation parameters  $\varepsilon = (\varepsilon_-, \varepsilon_+)$ . Subsequently we will

present a fully computable example pointing out some details of the iteration for fixed relaxation parameters.

Finally, in Chapter 5 we combine the results of the previous chapters. In Section 5.1 we will present a strategy for the relaxation parameter combined with the iterative Kačanov approach that will lead to an algebraic error decay. Comparing this decay with the rates obtained in Chapter 4 we see that adaptive strategies might be much more suitable than a fixed strategy. We will present error estimators for a fully adaptive scheme in Section 5.2. In the last section we discuss numerical experiments for three different model problems.

## 1.2 Notation

The next table gives a small overview over the used notation. However, in the most cases we try to use widely accepted standard notation.

Expression	Explanation
$\Omega$	an open and bounded subset of $\mathbb{R}^d$ with Lipschitz continuous boundary
$ A $	the Lebesgue measure of $A \subset \mathbb{R}^d$
$f \lesssim g$	there is a constant $c > 0$ such that $f \leq cg$ pointwise
$f \gtrsim g$	there is a constant $c > 0$ such that $f \leq cg$ pointwise
$f \approx g$	$f \lesssim g$ and $f \gtrsim g$
$C_0^\infty$	space of test functions
$L_{\text{loc}}^1$	the space of locally integrable functions
$L^p$	the space of $p$ -integrable functions
$W^{1,p}$	the Sobolev space with $p$ -integrable weak derivatives
$W_0^{1,p}$	the Sobolev space with $p$ -integrable weak derivatives and zero boundary values

## 2 Setting

In this chapter, we will introduce the  $p$ -Poisson equation. The theory of solutions of the classical formulation of this partial differential equation is “too narrow”, even for the homogeneous version (see [Lin06, Chapter 2]). Hence, it is much more convenient to study its weak formulation which will be introduced in the first section. The  $p$ -Laplace operator, the  $p$ -Poisson equation and its solution spaces – the Sobolev spaces – are known and described very well. Therefore, we will not prove most of the statements in Section 2.1 but refer to the respective literature.

After introducing the original problem, we will make a little detour through the field of Orlicz spaces. Although they are very well described, too, they are not as “famous” as their special cases  $L^p$ . Combining the integration properties of Orlicz functions and the concept of weak derivatives we will end up with the notion of the Orlicz Sobolev Spaces  $W_0^{1,\varphi}(\Omega)$ .

### 2.1 The $p$ -Poisson Equation and its Solution Space

First of all, we want to recall the Laplace operator very briefly:

$$\begin{aligned}\Delta : C^2(\Omega) &\rightarrow C^0(\Omega) \\ w &\mapsto \Delta w := \sum_{i=1}^n \frac{\partial^2}{\partial x_i^2} w.\end{aligned}$$

This linear operator arises in very many applications. Hence, it is natural to be interested in its inverse. This is mostly stated as a partial differential equation in the following way: For a given open domain  $\Omega \subset \mathbb{R}^d$  and  $f \in C^0(\Omega)$  one searches for a function  $u : \overline{\Omega} \rightarrow \mathbb{R}$  satisfying the Poisson equation

$$\begin{cases} -\Delta u = f & \text{in } \Omega \text{ and} \\ u = 0 & \text{on } \partial\Omega. \end{cases} \quad (2.1)$$

With Gauss's Divergence Theorem it is easy to see that if additionally  $f \in L^2(\Omega)$  and  $u$  being a solution to (2.1) we get for any  $\xi \in C^1(\bar{\Omega})$  the identity

$$\int_{\Omega} f\xi \, dx = \int_{\Omega} (-\Delta u)\xi \, dx = \int_{\Omega} \nabla u \nabla \xi \, dx. \quad (2.2)$$

Only looking at the beginning and the end of the above equation one sees that much lower requirements on  $f$  and  $u$  are necessary. Most interesting in this formulation is that one only needs first order derivatives of  $u$ . Moreover, these derivatives do not need to exist in every point of  $\Omega$ . The suitable function spaces to solve the weak Poisson equation will be defined later in this section.

As already pointed out, the Laplacian is a very important differential operator and the Poisson equation is the standard example for a linear elliptic partial differential equation. However, it has a very natural generalization, namely the strong  $p$ -Laplacian  $\Delta_p^s$  which is defined for any  $p \in [1, \infty)$  by

$$\Delta_p^s w := \operatorname{div}(|\nabla w|^{p-2} \nabla w).$$

Just as in that very case one can use Gauss's Divergence Theorem to show that every solution  $u$  of  $-\Delta_p^s u = f \in L^{p'}(\Omega)$  satisfies the weak  $p$ -Poisson equation

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi \, dx = \int_{\Omega} f\xi \, dx, \quad (2.3)$$

provided  $f$  and  $\xi$  are "nice enough". This is exactly the equation we are going to study in this thesis. As in (2.2) we see that also in (2.5) it is enough to have first order derivatives. For the integral on the left hand side to exist we need certain integrability requirements. All of them are combined in the next definition.

**Definition 2.1.** A function  $w \in L^1_{\text{loc}}(\Omega)$  is called *weakly differentiable* with weak derivative in  $i$ -th direction  $\tilde{w}_i$  iff if for all functions  $\xi \in C_0^\infty(\Omega)$  and every index  $i \in \{1, \dots, d\}$  the identity

$$\int_{\Omega} w \left( \frac{\partial}{\partial x_i} \xi \right) \, dx = \int_{\Omega} \tilde{w}_i \xi \, dx$$

holds. We denote the weak derivative in  $i$ -th direction by  $\frac{\partial}{\partial x_i}$ . For  $p \in [1, \infty]$ , the *Sobolev space* is defined by

$$W^{1,p}(\Omega) := \{w \in L^p(\Omega) : w \text{ is w. d. and } \forall i \in \{1, \dots, d\} : \frac{\partial}{\partial x_i} w \in L^p(\Omega)\}.$$

Note that it is easy to see that differentiable functions in the classical sense are also weakly differentiable.

Sobolev spaces, which also exist for higher order derivatives, are known very well and there is very much literature about them (see for example [AF03]). Hence, we just recapitulate the most important results in the next theorem. The cases  $p = 1$  and  $p = \infty$  – which are also the crucial integrability indices for the  $L^p$  spaces – are not relevant for this work, so we do not take them into account.

**Theorem 2.2.** *For  $p \notin \{1, \infty\}$  the Sobolev spaces  $W^{1,p}(\Omega)$  have the following properties:*

1.  $W^{1,p}(\Omega)$  is normed with  $\|w\|_{W^{1,p}(\Omega)} := \sqrt[p]{\|w\|_{L^p(\Omega)}^p + \sum_{i=1}^d \|\frac{\partial}{\partial x_i} w\|_{L^p(\Omega)}^p}$ .
2.  $W^{1,p}(\Omega)$  is complete.
3.  $W^{1,p}(\Omega)$  is uniformly convex.
4.  $W^{1,p}(\Omega)$  is separable.
5.  $W^{1,p}(\Omega)$  is reflexive.

With choosing  $u, \xi \in W^{1,p}(\Omega)$  we get well-definedness of the term

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi \, dx.$$

If  $p'$  is the Hölder conjugate exponent of  $p$  meaning  $\frac{1}{p} + \frac{1}{p'} = 1$  and choosing the function  $f \in L^{p'}(\Omega)$  we also have well-definedness of the right hand side of (2.3), but there are some right hand sides that can not be covered by this integral representation. We will discuss the generalization of the right hand side later in this section.

An other part of (2.4) is not yet covered by (2.3), too: The boundary values in the second line. Even for nice examples of  $\Omega$  – by means of a very regular boundary – it is not clear by definition how to deal with “ $u = 0$  on  $\partial\Omega$ ” if  $u$  is just member of  $W^{1,p}(\Omega)$ : The Sobolev space is just a subspace of  $L^p(\Omega)$ . Hence, its elements are only defined almost everywhere. In particular, there is no point evaluation of  $u$ . However, we can come around that with the following definition that “hides” the boundary values in the function space.

**Definition 2.3.** We define the *Sobolev space with zero boundary values* as the  $W^{1,p}$ -closure of the test functions:

$$W_0^{1,p}(\Omega) := \overline{C_0^\infty(\Omega)}^{\|\cdot\|_{W^{1,p}(\Omega)}}.$$

Indeed, the boundary values of Sobolev functions are understood better as they exist in a  $\partial\Omega$  almost everywhere sense with a suitable measure on  $\partial\Omega$  – at least when we may assume the boundary of  $\Omega$  to be nice enough.

**Theorem 2.4** (Trace Theorem). *Let  $\Omega \subset \mathbb{R}^d$  with  $d \geq 2$  and Lipschitz boundary and  $p \in [1, \infty)$ . Then, there exists a unique linear and bounded trace operator  $T : W^{1,p}(\Omega) \rightarrow L^p(\partial\Omega)$  with*

$$w \in C^0(\overline{\Omega}) \cap W^{1,p}(\Omega) \implies Tw = w|_{\partial\Omega}.$$

The proof of this statement is not trivial. Since it is a classical result for Sobolev spaces we will not state it here but refer to [AF03]. Indeed, there is a completely different characterization of  $W_0^{1,p}(\Omega)$ .

**Theorem 2.5.** *For  $p \in [1, \infty)$  the space  $W_0^{1,p}(\Omega)$  coincides with the kernel of the trace operator:*

$$W_0^{1,p}(\Omega) = \ker(T) = \{w \in W^{1,p}(\Omega) : Tw = 0\}.$$

In particular,  $W_0^{1,p}(\Omega)$  is a closed subspace of  $W^{1,p}(\Omega)$ . Historically, at that time not including boundary values, this result goes back to the seminal paper [MS64] titled “ $H = W$ ”. A proof of the result stated in this form can be found in [AF03, Theorem 5.37].

Now that we can handle boundary values we can use the trace theorem to see that any function satisfying the strong  $p$ -Poisson equation

$$\begin{cases} -\Delta_p^s u = f & \text{in } \Omega \text{ and} \\ u = 0 & \text{on } \partial\Omega. \end{cases} \quad (2.4)$$

is actually not only a function in  $W^{1,p}(\Omega)$  but even in  $W_0^{1,p}(\Omega)$ .

To have the most general formulation of the  $p$ -Poisson equation we need to understand the  $p$ -Laplacian not only as a differential operator but as a map defined on the suitable Sobolev space. We recall the expression

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi \, dx$$

to see that it is well-defined for  $u, \xi \in W_0^{1,p}(\Omega)$  and that the boundary values are preserved in the trace sense. Furthermore, it is linear in  $\xi$  due to the linearity of integration. Hence, the above term can also be seen as map

$$-\Delta_p : W_0^{1,p}(\Omega) \rightarrow (W_0^{1,p}(\Omega))^*$$

$$u \mapsto \int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \cdot \, dx,$$

where  $(W_0^{1,p}(\Omega))^*$  denotes the dual of  $W_0^{1,p}(\Omega)$ . This shows clearly that when one studies the equation  $-\Delta_p u = f$  the restriction to  $f \in L^{p'}(\Omega)$  is not necessary – the most general case is  $f \in (W_0^{1,p}(\Omega))^*$ .

We combine everything up to now to formulate *the  $p$ -Poisson problem*: Given  $\Omega \subset \mathbb{R}^d$  with  $d \geq 2$  and Lipschitz boundary and  $f \in (W_0^{1,p}(\Omega))^*$  find  $u \in W_0^{1,p}(\Omega)$  such that

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi = \langle f, \xi \rangle. \quad (2.5)$$

holds for all  $\xi \in W_0^{1,p}(\Omega)$ .

Note that the  $p$ -Laplacian and its  $p$ -Poisson equation has the same model character for nonlinear elliptic partial differential equations as its linear special case one gets by setting  $p = 2$ .

For the existence theory of a solution of the  $p$ -Poisson problem one can use the so called direct method in the calculus of variations. Therefore, we define the energy functional (or just energy)

$$\mathcal{J} : W_0^{1,p}(\Omega) \rightarrow \mathbb{R}$$

$$w \mapsto \frac{1}{p} \int_{\Omega} |\nabla w|^p \, dx - \langle f, w \rangle.$$

Following this direct method of calculus of variations one can obtain the following theorem (compare [Arn07, Section 5.2] and subtract the linear term  $\langle f, \cdot \rangle$ ).

**Theorem 2.6.** *For each  $f \in (W_0^{1,p}(\Omega))^*$  there is a unique minimizer of  $\mathcal{J}$ .*

One should mention that the uniqueness of the minimizer is due to the strict convexity of  $\mathcal{J}$ : Assume that there are  $u \neq \tilde{u}$  minimizers of  $\mathcal{J}$ . Then,

$$\mathcal{J}\left(\frac{1}{2}(u + \tilde{u})\right) < \frac{1}{2}\mathcal{J}(u) + \frac{1}{2}\mathcal{J}(\tilde{u}) = \mathcal{J}(u)$$

which obviously contradicts the minimization property of  $u$ .

Interestingly, there is a one-to-one relation of minimizers of this energy and solutions to (2.5) as stated in the theorem below.

**Theorem 2.7.** *We have*

$$u = \arg \min_{w \in W_0^{1,p}(\Omega)} \mathcal{J}(w) \iff \forall \xi \in W_0^{1,p}(\Omega) : \int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi \, dx = \langle f, \xi \rangle.$$

The implication “ $\implies$ ” can be obtained by calculating the Euler-Lagrange equation of  $\mathcal{J}$  by means of  $\frac{d}{dt} \mathcal{J}(u+t\xi)|_{t=0} \stackrel{!}{=} 0$ . On the other hand it is shown in [Růž06, Theorem 1.30 and following remark] that the equation admits a unique solution. Since  $\mathcal{J}$  also has a unique minimizer, the minimizer and the solution have to coincide.

Note that Theorem 2.7 also carries the linear Poisson equation as a special case for  $p = 2$ .

An other interesting property of  $\mathcal{J}$  is that it can be described by a norm on  $W_0^{1,p}(\Omega)$  and a linear term. To see that, we need the following inequality (see for example [Dzi10, Theorem 3.23]).

**Theorem 2.8** (Poincaré’s Inequality). *For  $p \in [1, \infty)$ , there is a constant  $c > 0$  such that for all  $w \in W_0^{1,p}(\Omega)$  the inequality*

$$\int_{\Omega} |w|^p \, dx \leq c \int_{\Omega} |\nabla w|^p \, dx$$

*holds.*

It is clear that this inequality can not hold for all  $w \in W^{1,p}(\Omega)$ . One can construct a contradiction by simply adding a large number to  $w$ . One important feature of this inequality is that for all  $w \in W_0^{1,p}(\Omega)$  we can deduce

$$\|\nabla w\|_{L^p(\Omega)} \lesssim \|w\|_{W^{1,p}(\Omega)} \lesssim \|\nabla w\|_{L^p(\Omega)},$$

so  $\|\nabla \cdot\|_{L^p(\Omega)}$  is a norm on  $W_0^{1,p}(\Omega)$  that is equivalent to the norm of Definition 2.3.

Furthermore,  $\mathcal{J}$  carries norm structure since it can be rewritten as

$$\mathcal{J}(w) = \frac{1}{p} \|\nabla w\|_{L^p(\Omega)}^p - \langle f, w \rangle.$$

We will use that in Chapter 3.



## 2.2 Further Function Spaces

For our discussions in Section 2.3 the Sobolev spaces as defined in Section 2.1 are not enough. To introduce Orlicz spaces we need the notion of N-functions. A very detailed introduction to N-functions and Orlicz spaces is [KR61]. This section bases very much on that book.

### 2.2.1 N-functions

This and the next subsection mainly deal as base for the more interesting Subsection 2.2.3 and for referring to some statements later. We will directly start with the definition of N-functions as in [KR61, Chapter I, §1, 1.].

**Definition 2.9.** A function  $\varphi : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  is said to be an *N-function* iff there is a right-continuous, for  $t > 0$  positive, non-decreasing function  $\varphi' : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  with  $\varphi'(0) = 0$  and  $\lim_{t \rightarrow \infty} \varphi'(t) = \infty$  such that

$$\varphi(t) = \int_0^t \varphi'(\tau) d\tau.$$

It is clear by definition that in particular every differentiable function  $\varphi$  is an N-function, if its derivative has the above mentioned properties. Furthermore, N-functions are continuous,  $\varphi(0) = 0$ , are strictly increasing away from zero and convex (see [KR61, Chapter I, §1, 4.]).

Later it will turn out that the following property of N-functions is fundamental.

**Definition 2.10.** An N-function  $\varphi$  satisfies the  $\Delta_2$ -condition iff there is a constant  $c > 0$  such that for all  $t \geq 0$  the estimate  $\varphi(2t) \leq c\varphi(t)$  holds.

When dealing with N-functions, one is oftentimes interested in the relation  $\varphi(t) \approx \varphi'(t)t$ . This can be characterized by the Simonenko indices as introduced in [Sim64].

**Definition 2.11.** For an N-function  $\varphi$  we define the *Simonenko indices* via

$$p^- := \inf_{t>0} \frac{t\varphi'(t)}{\varphi(t)} \leq \sup_{t>0} \frac{t\varphi'(t)}{\varphi(t)} =: p^+.$$

A simple consequence is the following estimate. It shows the relation of growth of an N-function and yields a sufficient condition for the  $\Delta_2$ -condition.

**Lemma 2.12.** *Let  $\varphi$  be an N-function with Simonenko-indices  $1 < p^-$  and  $p^+ < \infty$ . Then for all  $s, t \geq 0$ , the inequalities*

$$\min\{s^{p^-}, s^{p^+}\}\varphi(t) \leq \varphi(st) \leq \max\{s^{p^-}, s^{p^+}\}\varphi(t) \quad (2.6)$$

*hold. In particular,  $\varphi$  satisfies the  $\Delta_2$ -condition if  $p^+ < \infty$ .*

*Proof.* We restate the sketch of the proof in [Wan13] and only prove the upper bound in the case  $s \geq 1$  since all other bounds can be proven similarly. The statement is clear for  $t = 0$ , so let  $t \neq 0$ . Directly by the definition we get

$$\ln(\varphi(t))' = \frac{\varphi'(t)}{\varphi(t)} \leq \frac{p^+}{t}.$$

Hence,

$$\begin{aligned} \ln(\varphi(st)) - \ln(\varphi(t)) &= \int_t^{st} \ln(\varphi(\tau))' d\tau \\ &\leq \int_t^{st} \frac{p^+}{\tau} d\tau \\ &= p^+(\ln(st) - \ln(t)) \\ &= \ln(s^{p^+}). \end{aligned}$$

Applying the exponential function, the statement follows directly.  $\square$

To formulate the concept of duality of Orlicz spaces nicely, we use the next definition.

**Definition 2.13.** For an N-function  $\varphi$ , we define the *right inverse of  $\varphi'$*  by

$$\begin{aligned} (\varphi^*)' : \mathbb{R}_{\geq 0} &\rightarrow \mathbb{R}_{\geq 0} \\ t &\mapsto \sup_{\varphi'(s) \leq t} s \end{aligned}$$

and the *complementary N-function* by

$$\begin{aligned} \varphi^* : \mathbb{R}_{\geq 0} &\rightarrow \mathbb{R}_{\geq 0} \\ t &\mapsto \int_0^t (\varphi^*)'(\tau) d\tau. \end{aligned}$$

Note that  $(\varphi^*)'$  is the inverse of  $\varphi'$  if  $\varphi'$  is invertible. Additionally one can show that  $(\varphi^*)'$  satisfies all the requirements of  $\varphi'$  in Definition 2.9. Hence,  $\varphi^*$  is an N-function, too. Even the Simonenko indices of  $\varphi^*$  can be calculated by the Simonenko indices of  $\varphi$ , as shown in [FK97, Proposition 2.1].

**Lemma 2.14.** *Let  $\varphi$  be an N-function with Simonenko indices  $p^-, p^+ \in (1, \infty)$ . Then,*

$$(p^-)^* := \inf_{t>0} \frac{(\varphi^*)'(t)t}{\varphi^*(t)} \leq \sup_{t>0} \frac{(\varphi^*)'(t)t}{\varphi^*(t)} := (p^+)^*.$$

*In particular,  $\varphi^*$  satisfies the  $\Delta_2$  condition.*

An important tool where one can see how an N-function and its complementary N-function go hand in hand is Young's Inequality as for example stated in [KR61, I, §2, 2.]

**Lemma 2.15** (Young's Inequality). *Let  $\varphi$  be an N-function,  $\varphi^*$  its complementary N-function and  $s, t \geq 0$ . Then,*

$$st \leq \varphi(s) + \varphi^*(t).$$

More often, we will use a properly weighted version:

**Lemma 2.16** (Young's Inequality – Weighted Version). *Let  $\varphi$  be an N-function,  $\varphi^*$  its complementary N-function, both satisfying the  $\Delta_2$ -condition and  $s, t \geq 0$ . Then, for each  $\delta > 0$  there is  $c_\delta > 0$  such that*

$$st \leq \delta\varphi(s) + c_\delta\varphi^*(t).$$

*Proof.* The statement is a consequence of the unweighted version when  $\delta \geq 1$ , so let  $\delta \in (0, 1)$ . As already stated on page 9,  $\varphi$  is convex and  $\varphi(0) = 0$  and therefore  $\varphi(\delta s) \leq \delta\varphi(s)$ . Now, choose  $N_\delta \in \mathbb{N}$ , such that  $\delta 2^{N_\delta} \geq 1$ . Then,  $\varphi^*(\frac{t}{\delta}) \leq c^{N_\delta} \varphi^*(\frac{t}{2^{N_\delta} \delta}) \leq c^{N_\delta} \varphi^*(t)$  by induction where  $c$  is the constant of the  $\Delta_2$ -condition of  $\varphi^*$ . Hence, by the unscaled version of Young's Inequality we get

$$\begin{aligned} st &= \delta s \frac{t}{\delta} \\ &\leq \varphi(\delta s) + \varphi^*\left(\frac{t}{\delta}\right) \\ &\leq \delta\varphi(s) + c^{N_\delta} \varphi^*(t). \end{aligned} \quad \square$$

Of course one can refine the last estimate when one knows the Simonenko indices of  $\varphi$  and  $\varphi^*$ , respectively.

### 2.2.2 Shifted N-functions

The concept of shifted N-function goes back to [DE08] (respectively a preprint from 2005) and [RD07]. Certainly, the definition in these papers is slightly different – where we use  $a \vee t$ , the shifted N-functions in these papers had the term  $a + t$ . Once more we will see later, namely in Chapter 4, why our definition is more suitable for the use in our case, although one loses differentiability of  $\varphi_a$ .

**Definition 2.17.** For a given N-function  $\varphi$  and a shift  $a \geq 0$  we set

$$\varphi'_a(t) := \frac{\varphi'(a \vee t)}{a \vee t} t$$

and define the *shifted N-function* of  $\varphi$  via

$$\begin{aligned} \varphi_a : \mathbb{R}_{\geq 0} &\rightarrow \mathbb{R}_{\geq 0} \\ t &\mapsto \int_0^t \varphi'_a(\tau) d\tau. \end{aligned}$$

Note that  $\varphi'_a(t)$  satisfies all the properties of  $\varphi'$  as in Definition 2.9. Hence,  $\varphi_a$  is an N-function, too. It has quadratic growth for small arguments and coincides with  $\varphi$  for large arguments. The main strength of those shifted N-function lies in its relation with certain other terms appearing naturally when minimizing functionals like the one we defined in Section 2.1. We will define these quantities now.

**Definition 2.18.** For  $P \in \mathbb{R}^d$  we define  $A_\varphi, V_\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^d$  by

$$A_\varphi(P) := \begin{cases} \frac{\varphi'(|P|)}{|P|} P & \text{if } P \neq 0 \text{ and} \\ 0 & \text{if } P = 0 \end{cases} \quad \text{and} \quad V_\varphi(P) := \begin{cases} \sqrt{\frac{\varphi'(|P|)}{|P|}} P & \text{if } P \neq 0 \text{ and} \\ 0 & \text{if } P = 0. \end{cases}$$

The most important relation between  $\varphi_a$ ,  $A_\varphi$  and  $V_\varphi$  is stated as the consequence in the following lemma.

**Lemma 2.19.** *Let  $\varphi$  be an N-function with  $p^-, p^+ \in (1, \infty)$ . Furthermore, let there exist constants  $c_1, c_2 > 0$  such that for all  $P, Q \in \mathbb{R}^d$  the estimates*

$$\begin{aligned} (A_\varphi(P) - A_\varphi(Q)) \cdot (P - Q) &\geq c_1 \frac{\varphi'(|P| + |P - Q|)}{|P| + |P - Q|} |P - Q|^2 \text{ and} \\ |A_\varphi(P) - A_\varphi(Q)| &\leq c_2 \frac{\varphi'(|P| + |P - Q|)}{|P| + |P - Q|} |P - Q|^2 \end{aligned}$$

hold. Then,

$$\begin{aligned} (A_\varphi(P) - A_\varphi(Q))(P - Q) &\approx |V_\varphi(P) - V_\varphi(Q)|^2 \\ &\approx \varphi_{|P|}(|P - Q|) \\ &\approx \frac{\varphi'(|P| \vee |Q|)}{|P| \vee |Q|} |P - Q|^2. \end{aligned}$$

*Proof.* The lemma was firstly stated in the preprint of [DE08, Lemma 3] with the additional assumption the  $\varphi \in C^2((0, \infty))$ . However, this was only a technique issue. A proof of the statement without that assumption can be found in [RD07, Lemma 6.16]. Both versions only state the equivalences for the shifted N-function defined as in these papers. Note that for  $t, a \geq 0$  we have  $a + t \approx a \vee t$ . Hence by the Definition of the Simonenko indices and Lemma 2.12 we get

$$\frac{\varphi'(a + t)}{a + t} \approx \frac{\varphi(a + t)}{(a + t)^2} \approx \frac{\varphi(a \vee t)}{(a \vee t)^2} \approx \frac{\varphi'(a \vee t)}{a \vee t}. \quad (2.7)$$

Therefore, the shifted N-function of [DE08, Lemma 3] and [RD07, Lemma 6.16] is equivalent to our definition, so the result can also be applied to our version of shifted N-functions.  $\square$

It can be useful to change the shift. How this can be done is explained by the preprint of [DK08, Corollary 26], respectively in [RD07, Lemma 5.15]. Since the shift introduced here is pointwise equivalent to the shift used in these papers (see (2.7)), we can use the result without any further proof.

**Lemma 2.20** (Change of Shift). *Let  $\varphi$  be an N-function such that both  $\varphi$  and its complementary N-function  $\varphi^*$  satisfy the  $\Delta_2$  condition. Then, for all  $\delta \in (0, 1)$  there is a constant  $c_\delta$  such that for all  $P, Q, \in \mathbb{R}^d$  and  $t \geq 0$  the estimate*

$$\varphi_{|P|}(t) \leq c_\delta \varphi_{|Q|}(t) + \delta \varphi_{|Q|}(|P - Q|) \quad (2.8)$$

holds.

### 2.2.3 Orlicz and Orlicz Sobolev Spaces

After all the results we learned on N-functions we are finally ready to define a generalization of the usual  $L^p$  spaces. As [KR61] is a very good and detailed book about that topic, we only restate the important results.

**Definition 2.21.** Let  $\varphi$  be an N-function. We define

$$\mathcal{L}^\varphi(\Omega) := \{w : \Omega \rightarrow \mathbb{R} : w \text{ measurable and } \int_{\Omega} \varphi(|w|) dx < \infty\}.$$

Furthermore we say

$$w \sim \tilde{w} : \iff w = \tilde{w} \text{ almost everywhere.}$$

Then, we define the *Orlicz class* as

$$L^\varphi := \tilde{\mathcal{L}}^\varphi / \sim.$$

Note that when choosing  $\varphi(t) := \frac{1}{p}t^p$  one ends up the usual  $L^p$  spaces. Note that for an arbitrary N-function it is not necessarily true that it is a vector space, the next theorem (see [KR61, Chapter II, §8, 3., Theorem 8.2]) gives a characterization for that.

**Theorem 2.22.**  $L^\varphi(\Omega)$  is a vector space if and only if  $\varphi$  satisfies the  $\Delta_2$ -condition.

We collect some properties of Orlicz spaces.

**Theorem 2.23.** Let  $\varphi$  be an N-function such that both  $\varphi$  and its complementary N-function  $\varphi^*$  satisfy the  $\Delta_2$  condition. Then, the following statements are true:

1.  $L^\varphi(\Omega)$  is normed with  $\|w\|_{L^\varphi(\Omega)} := \inf\{\lambda > 0 : \int_{\Omega} \varphi(\frac{|w|}{\lambda}) dx \leq 1\}$ , see [KR61, Chapter II, §9, 7.].
2. A sequence of functions  $(w_n)_n$  in  $L^\varphi(\Omega)$  converges to  $w \in L^\varphi(\Omega)$  with respect to  $\|\cdot\|_{L^\varphi(\Omega)}$  if and only if  $\int_{\Omega} \varphi(|w_n - w|) dx \xrightarrow{n \rightarrow \infty} 0$ , see [KR61, Chapter II, §9, 6., Theorem 9.4].
3.  $L^\varphi(\Omega)$  is complete, see [KR61, Chapter II, §9, 2., Theorem 9.2].
4.  $L^\varphi(\Omega)$  is separable, see [KR61, Chapter II, §10, 3. & 4.].
5.  $L^\varphi(\Omega)$  is reflexive, see [KR61, Chapter II, §9, 5. and Chapter II, §14].

For the  $L^p$  spaces there are good comparison criteria, provided  $|\Omega| < \infty$  which we are assuming in this work. This can be generalized for Orlicz spaces as well, for a proof see [KR61, Chapter II, §8, 3.].

**Theorem 2.24.** *The inclusion  $L^{\varphi_1}(\Omega) \subseteq L^{\varphi_2}(\Omega)$  holds if and only if there are  $c, t_0$  such that  $\varphi_2(t) \leq c\varphi_1(t)$  for all  $t \geq t_0$ .*

With this result and Lemma 2.12 it is easy to deduce the following relations as already shown in [Wan13, Corollary 1.7].

**Corollary 2.25.** *Let  $\varphi$  be an N-function with  $p^-, p^+ \in (1, \infty)$ . Then*

$$L^{p^+}(\Omega) \subseteq L^\varphi(\Omega) \subseteq L^{p^-}(\Omega).$$

At this point it is natural to combine the notion of weak differentiability and non-standard growth.

**Definition 2.26.** For an N-function satisfying the  $\Delta_2$  condition we define the Orlicz Sobolev space via

$$W^{1,\varphi}(\Omega) := \{w \in L^\varphi(\Omega) : \forall i \in \{1, \dots, d\} \text{ we get } \frac{\partial}{\partial x_i} w \in L^\varphi(\Omega)\}$$

where  $\frac{\partial}{\partial x_i} w$  denotes the weak derivative in  $i$ -th direction.

If one combines not only the definitions of weak differentiability and integrability but also the properties of the underlying function spaces, one ends up with the following theorem.

**Theorem 2.27.** *Let  $\varphi$  be an N-function with  $p^-, p^+ \in (1, \infty)$ . Then, the following statements are true:*

1.  $W^{1,\varphi}(\Omega)$  is normed with  $\|w\|_{W^{1,\varphi}(\Omega)} := \|w\|_{L^\varphi(\Omega)} + \sum_{i=1}^d \|\frac{\partial}{\partial x_i} w\|_{L^\varphi(\Omega)}$ .
2. A sequence of functions  $(w_n)_n$  in  $W^{1,\varphi}(\Omega)$  converges to  $w \in W^{1,\varphi}(\Omega)$  with respect to  $\|\cdot\|_{W^{1,\varphi}(\Omega)}$  if and only if  $\int_\Omega \varphi(|w_n - w|) dx \xrightarrow{n \rightarrow \infty} 0$  and for all  $i \in \{1, \dots, d\}$  also  $\int_\Omega \varphi(|\frac{\partial}{\partial x_i}(w_n - w)|) dx \xrightarrow{n \rightarrow \infty} 0$ .
3.  $W^{1,\varphi}(\Omega)$  is complete.
4.  $W^{1,\varphi}(\Omega)$  is separable.
5.  $W^{1,\varphi}(\Omega)$  is reflexive.

As for the Sobolev spaces we can equip  $W^{1,\varphi}(\Omega)$  with zero boundary values.

**Definition 2.28.** We define the *Orlicz Sobolev space with zero boundary values* as the  $W^{1,\varphi}(\Omega)$ -closure of the test functions:

$$W_0^{1,\varphi}(\Omega) := \overline{C_0^\infty(\Omega)}^{\|\cdot\|_{W^{1,\varphi}(\Omega)}}.$$

Note that Poincaré's Inequality also holds in the Orlicz Sobolev setting.

**Theorem 2.29** (Poincaré's Inequality – Orlicz version). *For any  $N$ -function  $\varphi$  with  $p^-, p^+ \in (1, \infty)$  there is a constant  $c > 0$  such that for all  $w \in W_0^{1,\varphi}(\Omega)$  the inequality*

$$\int_{\Omega} \varphi(|w|) dx \leq c \int_{\Omega} \varphi(|\nabla w|) dx$$

holds.

*Proof.* From [DRS10, Theorem 6.5] we know that for all  $w \in W^{1,\varphi}(\Omega)$  with mean zero, i. e.  $\int_{\Omega} w dx = 0$ , we have – additionally using that  $\varphi$  satisfies the  $\Delta_2$  condition – the estimate

$$\int_{\Omega} (|w|) dx \lesssim \int_{\Omega} \varphi(|\nabla w|) dx.$$

Without loss of generality we may assume that there is  $R > 1$  such that  $\Omega \subset [1, R]^d$  due to the translation invariance of the Lebesgue measure and since we always assume  $\Omega$  to be bounded. Now suppose  $w \in W_0^{1,p}(\Omega)$  arbitrary. Then, we define  $\tilde{\Omega} := \{x - (R+1)\mathbb{1} : x \in \Omega\}$  where  $\mathbb{1}$  is the  $\mathbb{R}^d$ -vector only containing ones. It is clear that  $\tilde{\Omega} \in [-R, -1]^d$ . Hence,  $\Omega \cup \tilde{\Omega} \subset [-R, R]^d$  and  $\Omega \cap \tilde{\Omega} = \emptyset$ . Since the boundary of  $\Omega$  is Lipschitz continuous and  $W^{1,\varphi}(\Omega) \subseteq W^{1,p^-}(\Omega)$  we know that the function

$$\tilde{w} : (-R, R)^d \rightarrow \mathbb{R}$$

$$x \mapsto \begin{cases} w(x) & \text{for } x \in \Omega, \\ -w(x + (R+1)\mathbb{1}) & \text{for } x \in \tilde{\Omega} \text{ and} \\ 0 & \text{for } x \in (-R, R)^d \setminus (\Omega \cup \tilde{\Omega}) \end{cases}$$

is weakly differentiable (see for example [AF03, Theorem 5.29]). Furthermore, it is clear that  $\tilde{w} \in W^{1,\varphi}((-R, R)^d)$ , since we have

$$\int_{(-R,R)^d} \varphi(|\tilde{w}|) dx = 2 \int_{\Omega} \varphi(|\tilde{w}|) dx$$

$$< \infty,$$



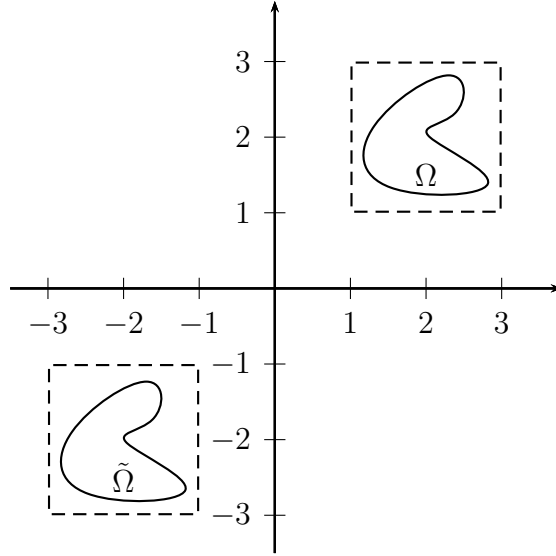


Figure 2.1: Exemplary picture for  $d = 2$  and  $R = 3$  showing  $\Omega$  and  $\tilde{\Omega}$ .

and the same argument applies to  $|\frac{\partial}{\partial x_i} \tilde{w}|$  instead of  $|\tilde{w}|$  for any  $i \in \{1, \dots, d\}$ . Additionally,

$$\int_{(-R,R)^d} \tilde{w} \, dx = 0.$$

This yields that we can use the first inequality of this proof on  $(-R, R)^d$  to deduce

$$\begin{aligned} \int_{\Omega} \varphi(|w|) \, dx &= \frac{1}{2} \int_{(-R,R)^d} \varphi(|\tilde{w}|) \, dx \\ &\lesssim \frac{1}{2} \int_{(-R,R)^d} \varphi(|\nabla \tilde{w}|) \, dx \\ &= \int_{\Omega} \varphi(|\nabla w|) \, dx. \end{aligned} \quad \square$$

With that, we can state a theorem which is well-known for  $\varphi(t) := \frac{1}{p} t^p$ , respectively  $W_0^{1,\varphi}(\Omega) = W_0^{1,p}(\Omega)$ .

**Theorem 2.30.** *For any N-function  $\varphi$  with  $p^-, p^+ \in (1, \infty)$  and all functions  $w \in W_0^{1,\varphi}(\Omega)$  we have*

$$\|w\|_{W^{1,\varphi}(\Omega)} \approx \|\nabla w\|_{L^\varphi(\Omega)}.$$

In particular,  $\|\nabla \cdot\|_{L^\varphi(\Omega)}$  is a to  $\|\cdot\|_{W^{1,\varphi}(\Omega)}$  equivalent norm on  $W_0^{1,\varphi}(\Omega)$ .

*Proof.* It is clear from the definition in Theorem 2.23 and from the monotonicity of  $\varphi$  that for  $w, \tilde{w} \in L^\varphi(\Omega)$  the estimate  $|w| \leq |\tilde{w}|$  almost everywhere implies that  $\|w\|_{L^\varphi(\Omega)} \leq \|\tilde{w}\|_{L^\varphi(\Omega)}$ . Hence, for all  $i \in \{1, \dots, d\}$  we get the estimate  $\|\frac{\partial}{\partial x_i} w\|_{L^\varphi(\Omega)} \leq \|\nabla w\|_{L^\varphi(\Omega)}$ , where Theorem 2.29 yields  $\|w\|_{L^\varphi(\Omega)} \lesssim \|\nabla w\|_{L^\varphi(\Omega)}$ . Therefore,

$$\begin{aligned} \|w\|_{W^{1,\varphi}(\Omega)} &= \|w\|_{L^\varphi(\Omega)} + \sum_{i=1}^d \left\| \frac{\partial}{\partial x_i} w \right\|_{L^\varphi(\Omega)} \\ &\lesssim \|\nabla w\|_{L^\varphi(\Omega)}. \end{aligned}$$

On the other hand, all norms on  $\mathbb{R}^d$  are equivalent, so  $|\nabla w| \lesssim \sum_{i=1}^d |\frac{\partial}{\partial x_i} w|$ . So

$$\begin{aligned} \|\nabla w\|_{L^\varphi(\Omega)} &\lesssim \sum_{i=1}^d \left\| \frac{\partial}{\partial x_i} w \right\|_{L^\varphi(\Omega)} \\ &\leq \|w\|_{W^{1,\varphi}(\Omega)}. \end{aligned} \quad \square$$

With this and since  $W_0^{1,\varphi}(\Omega)$  is a closed subspace of  $W^{1,\varphi}(\Omega)$  by definition, we can restate Theorem 2.27 for Orlicz Sobolev spaces with zero boundary values.

**Theorem 2.31.** *Let  $\varphi$  be an  $N$ -function with  $p^-, p^+ \in (1, \infty)$ . Then, the following statements are true:*

1.  $W_0^{1,\varphi}(\Omega)$  is normed with  $\|\nabla \cdot\|_{L^\varphi(\Omega)}$ .
2. A sequence of functions  $(w_n)_n$  in  $W_0^{1,\varphi}(\Omega)$  converges to  $w \in W_0^{1,\varphi}(\Omega)$  with respect to  $\|\nabla \cdot\|_{L^\varphi(\Omega)}$  if and only if  $\int_\Omega \varphi(|\nabla(w_n - w)|) dx \xrightarrow{n \rightarrow \infty} 0$ .
3.  $W_0^{1,\varphi}(\Omega)$  is complete.
4.  $W_0^{1,\varphi}(\Omega)$  is separable.
5.  $W_0^{1,\varphi}(\Omega)$  is reflexive.

## 2.3 Presentation of the Algorithm

As we have seen in Section 2.1 there are two legitimate views on weak partial differential equations: The equation an the energy point of view. We will discuss

the derivation of the algorithm from an equation point of view in the first subsection whereas we will find an energy point of view in the subsection after that.

From now on we restrict ourselves to the case where  $p \in (1, 2)$ . We will see in Section 4.2 that this is not only a problem of the technique of the proof but show an example where our algorithm fails in parts and totally for  $p > 2$  and  $p > 3$ , respectively.

### 2.3.1 Derivation of the Algorithm

We recall the  $p$ -Poisson problem from page 7: Given  $\Omega \subset \mathbb{R}^d$  with  $d \geq 2$  and Lipschitz boundary and  $f \in (W_0^{1,p}(\Omega))^*$  find  $u \in W_0^{1,p}(\Omega)$  such that

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi = \langle f, \xi \rangle. \quad (2.5)$$

holds for all  $\xi \in W_0^{1,p}(\Omega)$ .

If one pretends to know  $u$  already and one “hides” the term  $|\nabla u|^{p-2}$  which is responsible for the non-linearity of the equation in  $a := |\nabla u|$ , equation (2.5) becomes

$$\int_{\Omega} a^{p-2} \nabla u \nabla \xi = \langle f, \xi \rangle. \quad (2.9)$$

Note that this equation is linear in  $u$  and well known as a weighted Poisson problem. In particular, many numerical approaches are known for this problem such as the Finite Element Method. The representation in (2.9) now suggests to define a sequence  $(v_n)_n$  recursively as solutions to

$$\int_{\Omega} |\nabla v_n|^{p-2} \nabla v_{n+1} \nabla \xi = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega). \quad (2.10)$$

For the weighted Poisson equation to be well defined, the weight needs to be in  $L^\infty$  and away from zero in the sense that there is a constant  $c > 0$  such that  $a^{p-2} > c$  almost everywhere. As the solutions and therefore the iterated  $v_n$  will be elements of  $W_0^{1,2}(\Omega)$  it is not possible to assume even one of those two requirements: the weight  $|\nabla v_n|^{p-2}$  degenerates for  $|\nabla v_n| = 0$  and  $|\nabla v_n| = \infty$ . The fact that for our choices of  $p$  we have  $W_0^{1,2}(\Omega) \subseteq W_0^{1,p}(\Omega)$  ensures that  $f \in (W_0^{1,p}(\Omega))^* \subseteq (W_0^{1,2}(\Omega))^*$ .

Hence, the right hand side does not need to be modified to ensure (2.10) to be well-defined.

We will cope the task of degenerating weights by introducing a relation parameter. Therefore, let  $0 < \varepsilon_- \leq 1 \leq \varepsilon_+ < \infty$  and we define the relaxation interval  $\varepsilon := (\varepsilon_-, \varepsilon_+) \subset \mathbb{R}$  (for  $\varepsilon_- = \varepsilon_+ = 1$  we define  $\varepsilon := \{1\}$ ). We will truncate  $a$  to the closure of that interval which is nothing but the projection onto  $\varepsilon$ :

$$\Pi_\varepsilon(a) := \varepsilon_- \vee a \wedge \varepsilon_+ := \min\{\max\{\varepsilon_-, a\}, \varepsilon_+\}$$

which is illustrated in Figure 2.2 as well as the affect on the weight.

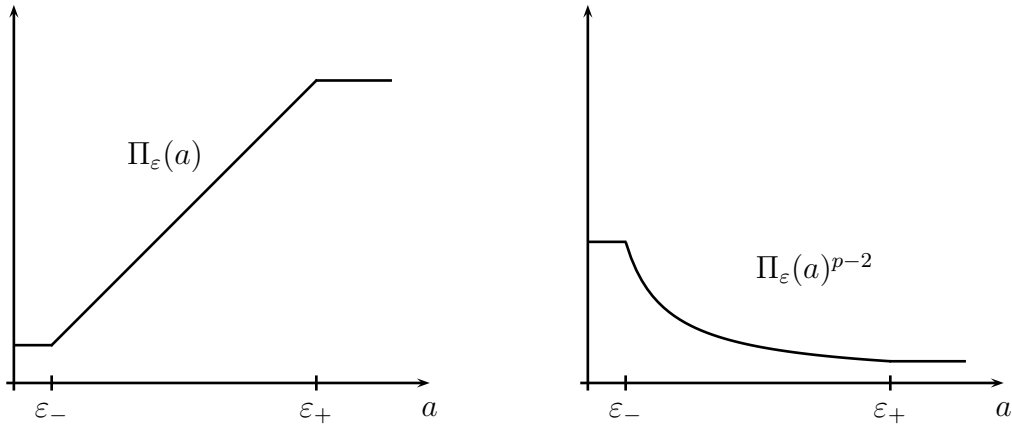


Figure 2.2: Plot of the constraint and the weight for  $p = 1.1$ .

Replacing  $a^{p-2}$  by  $\Pi_\varepsilon(a)^{p-2}$  in (2.10) we get an admissible inductive clause and

formulate the following algorithm.

**Algorithm:** The relaxed  $p$ -Kačanov algorithm

**Data:**  $f \in (W_0^{1,p}(\Omega))^*$

**Result:** Approximate solution of the  $p$ -Poisson problem (2.5).

$n := 0$ ;

$\varepsilon_{-1,-} := 1$ ;

$\varepsilon_{-1,+} := 1$ ;

**while** desired accuracy is not achieved **yet do**

    Calculate  $v_n$  by means of

$$\int_{\Omega} (a_{n-1})^{p-2} \nabla v_n \cdot \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega);$$

    Choose a new relaxation interval  $\varepsilon_n \supseteq \varepsilon_{n-1}$ ;

    Define  $a_n := \Pi_{\varepsilon_n}(|\nabla v_n|)$ ;

    Update  $n \rightsquigarrow n + 1$ ;

**end**

Note that it is not a problem not to initialize  $a_{-1}$  since we have chosen the initial values for  $\varepsilon$  such that  $\varepsilon_{-1,-} \vee a_{-1} \wedge \varepsilon_{-1,+} = 1 \vee a_{-1} \wedge 1 = 1$  anyway. In particular,  $v_0$  solves the linear Poisson equation

$$\int_{\Omega} \nabla v_0 \cdot \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega).$$

### 2.3.2 The Algorithm from an Energy Point of View

As we have already seen weak formulated partial differential equations can often-times be linked to energy minimizing problems. In this section we aim for such an energy formulation for the algorithm presented in Subsection 2.3.1. For the ease of readability we fix  $\varepsilon_n \equiv \varepsilon$ .

Similar to Theorem 2.7 it is known from the theory for the weighted Poisson equation that

$$\int_{\Omega} \Pi_{\varepsilon}(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega)$$

with  $v_{n+1} \in W_0^{1,2}(\Omega)$  is equivalent to

$$v_{n+1} = \arg \min_{w \in W_0^{1,2}} \frac{1}{2} \int_{\Omega} \Pi_{\varepsilon}(|\nabla v_n|)^{p-2} |\nabla w|^2 dx - \langle f, w \rangle. \quad (2.11)$$

The theory for this problem does not respect changes of the weight – it is mostly hidden in the scalar product. Therefore, it does not fit very well when the weight changes from iteration step to iteration step, i. e. when you want to compare

$$\int_{\Omega} \Pi_{\varepsilon}(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi dx \quad \text{to} \quad \int_{\Omega} \Pi_{\varepsilon}(|\nabla v_{n+1}|)^{p-2} \nabla v_{n+1} \nabla \xi dx.$$

Furthermore, the equation for  $v_n$  does not yield pointwise information, so it can not be plugged in the inductive clause iteratively.

To overcome this problem we want to have the energy to provide some minimization property for a change in the weight. Therefore, we add a function  $h : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  to the energy as defined in (2.11) such that

$$v_{n+1} = \arg \min_{w \in W_0^{1,2}} \int_{\Omega} \frac{1}{2} \Pi_{\varepsilon}(|\nabla w|)^{p-2} |\nabla v_{n+1}|^2 + h(|\nabla w|) dx - \langle f, w \rangle. \quad (2.12)$$

We will see that the function  $h : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  with  $h(t) := (\frac{1}{p} - \frac{1}{2}) \Pi_{\varepsilon}(t)^p$  is a proper choice and define

$$\begin{aligned} \mathcal{J}_{\varepsilon}^s : W_0^{1,2}(\Omega) \times L_{\text{loc}}^1(\Omega) &\rightarrow \mathbb{R} \\ (w, a) &\mapsto \int_{\Omega} \frac{1}{2} \Pi_{\varepsilon}(|a|)^{p-2} |\nabla w|^2 + (\frac{1}{p} - \frac{1}{2}) \Pi_{\varepsilon}(|a|)^p dx - \langle f, w \rangle. \end{aligned}$$

First, we need to prove a technical lemma.

**Lemma 2.32.** *For  $\beta : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$  measurable we consider the functional*

$$\begin{aligned} \mathcal{K} : \mathcal{M} &\rightarrow \overline{\mathbb{R}} \\ a &\mapsto \int_{\Omega} \beta(x, a(x)) dx \end{aligned}$$

where  $\mathcal{M}$  is a vector space of equivalence classes of functions mapping  $\Omega$  to  $\mathbb{R}$  with  $\sim$  as in Definition 2.21. Then,

$$m \in \arg \min_{a \in \mathcal{M}} \mathcal{J}(a) \iff \forall a \in \mathcal{M} : \beta(x, m(x)) \leq \beta(x, a(x)) \quad a.e.$$

*Proof.* The implication " $\Leftarrow$ " is easy to see by integration over  $\Omega$ . We prove " $\Rightarrow$ " by proving the contra position. Therefore, we assume there is  $r \in \mathcal{M}$  and  $\omega \subseteq \Omega$  with  $|\omega| \neq 0$  such that for almost every  $x \in \omega$  we have the inequality  $\beta(x, m(x)) > \beta(x, r(x))$ . Now, we define

$$s(x) := \begin{cases} r(x) & x \in \Omega' \\ m(x) & x \in \Omega \setminus \Omega' \end{cases}$$

and get the contradiction  $\mathcal{K}(s) < \mathcal{K}(m)$  by direct calculation.  $\square$

With the aid of the last lemma we are able to prove the required minimization property for the weight.

**Lemma 2.33.** *Let  $p \in (1, 2)$ . Then, for fixed  $w \in W_0^{1,2}(\Omega)$  there is a global minimum of  $\mathcal{J}_\varepsilon^s(w, \cdot)$  in  $L_{\text{loc}}^1(\Omega)$ . Under the additional constraint that the minimizer satisfies  $a = \Pi_\varepsilon(a)$ , the minimizer is unique and admits the representation*

$$\arg \min_{\substack{a \in L_{\text{loc}}^1(\Omega) \\ a = \Pi_\varepsilon(a)}} \mathcal{J}_\varepsilon^s(w, a) = \Pi_\varepsilon(|\nabla w|).$$

*Proof.* Using Lemma 2.32 we reduce the minimization problem to a minimization problem in one real variable. Therefore, we define

$$\beta(x, a) := \frac{1}{2} \Pi_\varepsilon(|a|)^{p-2} |\nabla w(x)|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) \Pi_\varepsilon(|a|)^p$$

with the aim of minimizing  $\beta$  in the variable  $a$  for a.e.  $x \in \Omega$ . This leads to an a.e. defined function  $a : \Omega \rightarrow \mathbb{R}$ . So let  $x \in \Omega$  be fixed. Thus, we can rewrite

$$\beta(x, a) = \begin{cases} \frac{1}{2} \varepsilon^{p-2} |\nabla w(x)|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) \varepsilon^p & \text{for } |a| \leq \varepsilon_- \\ \frac{1}{2} |a|^{p-2} |\nabla w(x)|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) |a|^p & \text{for } |a| \in \varepsilon \\ \frac{1}{2} \varepsilon^{2-p} |\nabla w(x)|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) \varepsilon^{-p} & \text{for } |a| \geq \varepsilon_+ \end{cases}$$

to see that  $\beta(x, \mathbb{R}) = \beta(x, \bar{\varepsilon})$ . Together with continuity of  $\beta(x, \cdot)$  this already yields the existence of a minimizer. Furthermore,  $\beta(x, \cdot)$  is differentiable a.e. (i.e. on  $\mathbb{R} \setminus \{\pm\varepsilon_-, \pm\varepsilon_+\}$ ) and we get

$$\frac{d\beta}{da}(x, a) = \begin{cases} 0 & \text{for } |a| < \varepsilon_- \\ \frac{2-p}{2} |a|^{p-4} a (|a|^2 - |\nabla w(x)|^2) & \text{for } |a| \in \varepsilon \\ 0 & \text{for } |a| > \varepsilon_+. \end{cases}$$

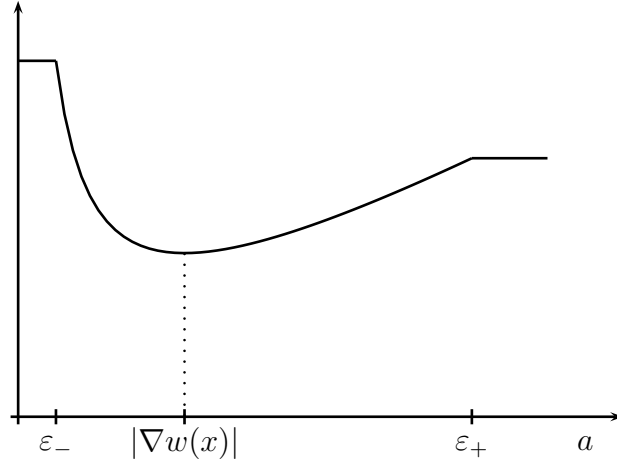


Figure 2.3: Plot of  $\beta(|\nabla w(x)|, a)$  for  $|\nabla w(x)| \in \varepsilon$ .

To get the desired representation of the minimizer under the assumption  $a = \Pi_\varepsilon(a)$  it is sufficient to search for minimizers in  $\bar{\varepsilon}$ . We have to distinguish three cases to determine the minimizer uniquely:

$|\nabla u(x)| \leq \varepsilon_-$ : Then,  $\frac{d\beta}{da}(x, a) > 0$  on  $\varepsilon$ . Due to continuity of  $\beta(x, \cdot)$  the unique minimum is attained in  $a = \varepsilon_-$ .

$|\nabla u(x)| \in \varepsilon$ : Then,  $\frac{d\beta}{da}(x, a) < 0$  on the interval  $(\varepsilon_-, |\nabla u(x)|)$  and  $\frac{d\beta}{da}(x, a) > 0$  on  $(|\nabla u(x)|, \varepsilon_+)$ . So the unique minimum is attained in  $a = |\nabla u(x)|$ .

$|\nabla u(x)| \geq \varepsilon_+$ : Then,  $\frac{d\beta}{da}(x, a) < 0$  on  $\varepsilon$ . Due to continuity of  $\beta(x, \cdot)$  the unique minimum is attained in  $a = \varepsilon_+$ .

In deed, this procedure leads to a minimizer since  $\beta(x, \mathbb{R}) = \beta(x, \bar{\varepsilon})$ .  $\square$

Note that the last lemma is not interesting for the case  $p = 2$  since  $\mathcal{J}_\varepsilon^s(w, \cdot)$  is constant then. For  $p > 2$  the statement is not true. This is due to the fact that  $\beta(x, \cdot)$  is increasing on  $-\varepsilon := (-\varepsilon_+, -\varepsilon_-)$  and decreasing on  $\varepsilon$  and hence  $\pm\varepsilon_+ = \arg \min_{\{a \in \mathbb{R}\}} \beta(x, a)$ .

As a consequence of Lemma 2.33 we can rewrite the algorithm presented in Subsection 2.3.1 on page 21 based on alternately minimizing

$$\begin{aligned} \mathcal{J}_\varepsilon^s : W_0^{1,2}(\Omega) \times L_{\text{loc}}^1(\Omega) &\rightarrow \mathbb{R} \\ (w, a) &\mapsto \int_{\Omega} \frac{1}{2} \Pi_\varepsilon(|a|)^{p-2} |\nabla w|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) \Pi_\varepsilon(|a|)^p dx - \langle f, w \rangle. \end{aligned}$$



in each argument.

**Algorithm:** The relaxed  $p$ -Kačanov algorithm

**Data:**  $f \in (W_0^{1,p}(\Omega))^*$

**Result:** Approximate solution of the  $p$ -Poisson problem (2.5).

$n := 0$ ;

$\varepsilon_{-1,-} := 1$ ;

$\varepsilon_{-1,+} := 1$ ;

**while** desired accuracy is not achieved yet **do**

    Calculate  $v_n$  by means of

$$v_n = \arg \min_{w \in W_0^{1,2}(\Omega)} \mathcal{J}_{\varepsilon_{n-1}}(w, a_{n-1});$$

    Choose a new relaxation interval  $\varepsilon_n \supseteq \varepsilon_{n-1}$ ;

    Calculate  $v_n$  by means of

$$a_n = \arg \min_{\substack{a \in L_{\text{loc}}^1(\Omega) \\ a = \Pi_\varepsilon(a)}} \mathcal{J}_\varepsilon^s(v_n, a);$$

    Update  $n \rightsquigarrow n + 1$ ;

**end**

### 2.3.3 The Constrained $p$ -Poisson Equation and Its Relaxed Energy

For the further discussion, we neglect the iteration of  $\varepsilon_n$  and look at a fixed  $\varepsilon$ . As we already observed, our algorithm is based on alternatingly minimizing  $\mathcal{J}_\varepsilon^s$ . Furthermore we have already seen that convexity yields good minimization properties.

**Lemma 2.34.**  $\mathcal{J}_\varepsilon^s$  is strictly convex on the set

$$W_0^{1,2}(\Omega) \times \{a \in L_{\text{loc}}^1(\Omega) : a = \Pi_\varepsilon(a)\}.$$

*Proof.* It suffices to show strict convexity of

$$\begin{aligned} f : \mathbb{R}_{>0} \times \mathbb{R} &\rightarrow \mathbb{R} \\ (x, y) &\mapsto x^{p-2}y^2. \end{aligned}$$

We can calculate the Hessian matrix

$$H_f(x, y) = \begin{pmatrix} (p-2)(p-3)x^{p-4}y^2 & 2(p-2)x^{p-3}y \\ 2(p-2)x^{p-3}y & 2x^{p-2} \end{pmatrix}$$

and will show that it is positive definite for all  $(x, y) \in \mathbb{R}_{>0} \times \mathbb{R}$ . For  $a, b \in \mathbb{R}$  we get

$$\langle \begin{pmatrix} a \\ b \end{pmatrix}, H_f(x, y) \begin{pmatrix} a \\ b \end{pmatrix} \rangle = x^{p-4}((p-2)(p-3)a^2y^2 - 4(2-p)abxy + 2b^2x^2)$$

which is positive if and only if

$$4(2-p)abxy < (p-2)(p-3)a^2y^2 + 2b^2x^2. \quad (2.13)$$

Applying Young's inequality we get

$$4(2-p)abxy \leq 2(2-p)^2a^2y^2 + 2b^2x^2$$

and since  $2(2-p)^2 < (p-2)(p-3)$  for  $p \in (1, 2)$  inequality (2.13) holds for all  $a, b \in \mathbb{R}$ . The application of convexity of  $\mathbb{R}_{>0} \times \mathbb{R}$  finishes the proof.  $\square$

Although one could now already argue with the direct method of variations to get a minimizing tuple of  $\mathcal{J}_\varepsilon$  we first restrict the problem to a smaller set of functions. As we already have seen in Lemma 2.33 minimizing with respect to the weight yields a special representation for the minimum, namely  $a = \Pi_\varepsilon(|\nabla w|)$ . We will see that this is enough to give reason to define

$$\begin{aligned} \mathcal{J}_\varepsilon : W_0^{1,2}(\Omega) &\rightarrow \mathbb{R} \\ w &\mapsto \mathcal{J}_\varepsilon^s(w, |\nabla w|) = \int_{\Omega} \kappa_\varepsilon(|\nabla w|) dx - \langle f, w \rangle \end{aligned} \quad (2.14)$$

where  $\kappa_\varepsilon : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  admits the representation

$$\kappa_\varepsilon(t) := \begin{cases} \frac{1}{2}\varepsilon_-^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_-^p & \text{for } t \leq \varepsilon_-, \\ \frac{1}{p}t^p & \text{for } t \in \varepsilon \text{ and} \\ \frac{1}{2}\varepsilon_+^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p & \text{for } t \geq \varepsilon_+. \end{cases} \quad (2.15)$$

The next theorem can be seen exactly like the existence of the unique minimizer of  $\mathcal{J}$  for Theorem 2.6.

**Theorem 2.35.** *For each  $f \in (W_0^{1,p}(\Omega))^*$  there is a unique minimizer of  $\mathcal{J}_\varepsilon$ .*

As we will see in the next theorem  $\mathcal{J}_\varepsilon$  and the minimizer of  $\mathcal{J}_\varepsilon$  will play an important role so we define

$$u_\varepsilon := \arg \min_{w \in W_0^{1,2}(\Omega)} \mathcal{J}_\varepsilon(w). \quad (2.16)$$

It is also the unique solution of the Euler Lagrange equation of  $\mathcal{J}_\varepsilon$ :

$$\int_{\Omega} \frac{\kappa'_\varepsilon(|\nabla u_\varepsilon|)}{|\nabla u_\varepsilon|} \nabla u_\varepsilon \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega). \quad (2.17)$$

We want to state one nice property of  $\mathcal{J}_\varepsilon(u_\varepsilon)$ .

**Lemma 2.36.** *The function*

$$\begin{aligned} \rho : (0, 1) \times (1, \infty) &\rightarrow \mathbb{R} \\ (\varepsilon_-, \varepsilon_+) &\mapsto \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) \end{aligned} \quad (2.18)$$

is convex.

*Proof.* We show this basically by proving that

$$\begin{aligned} \kappa : [0, \infty) \times (0, 1) \times (1, \infty) &\rightarrow \mathbb{R}_{\geq 0} \\ (t, \varepsilon_-, \varepsilon_+) &\mapsto \kappa_\varepsilon(t) \end{aligned}$$

is convex. Therefore, let  $\theta \in (0, 1)$  and choose two arbitrary triples

$$(t, \varepsilon_-, \varepsilon_+), (t', \varepsilon'_-, \varepsilon'_+) \in [0, \infty) \times (0, 1) \times (1, \infty).$$

For  $P, Q \in \mathbb{R}^k$  we use the notation

$$[P, Q]_\theta := \theta P + (1 - \theta)Q$$

and define  $\theta_-$  as solution to  $[\varepsilon_-, \varepsilon'_-]_{\theta_-} = [t, t']_{\theta_-}$  and  $\theta_+$  as solution to  $[\varepsilon_+, \varepsilon'_+]_{\theta_+} = [t, t']_{\theta_+}$  – in case only one or none of those values exists the proof is easier and similar to this one. Note that  $\theta_-$  and  $\theta_+$  are unique. Without loss of generality we may assume that  $\theta_- \leq \theta_+$ .

We already know from the proof of Lemma 2.34 that  $x^{p-2}y^2$  is convex. Hence, the terms

$$\frac{1}{2}\varepsilon_-^{p-2}t^2 + \left(\frac{1}{p} - \frac{1}{2}\right)\varepsilon_-^p, \quad \frac{1}{p}t^p \quad \text{and} \quad \frac{1}{2}\varepsilon_+^{p-2}t^2 + \left(\frac{1}{p} - \frac{1}{2}\right)\varepsilon_+^p$$

are convex with respect to  $(t, \varepsilon_-, \varepsilon_+)$  and in particular with respect to  $\theta$  when one plugs in  $([t, t']_\theta, [\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta)$ . This implies that  $\kappa([t, t']_\theta, [\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta)$  is convex with respect to  $\theta$  on  $[0, \theta_-]$ ,  $[\theta_-, \theta_+]$  and  $[\theta_+, 1]$ . Furthermore,  $\kappa$  is  $C^1$  in each component and hence  $C^1$  with respect to  $\theta$ . Since any piecewise convex and continuously differentiable function in one variable is convex on the whole domain we get

$$\kappa([t, t']_\theta, [\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta) \leq [\kappa(t, \varepsilon_-, \varepsilon_+), \kappa(t', \varepsilon'_-, \varepsilon'_+)]_\theta.$$

To deduce the statement it is enough to show that  $\mathcal{J}_\varepsilon(u_\varepsilon)$  is convex in  $(\varepsilon_-, \varepsilon_+)$ . This follows by

$$\begin{aligned} \mathcal{J}_{([\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta)}(u_{([\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta)}) &\leq \mathcal{J}_{([\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta)}([u_\varepsilon, u_{\varepsilon'}]_\theta) \\ &\leq \int_{\Omega} \kappa_{([\varepsilon_-, \varepsilon'_-]_\theta, [\varepsilon_+, \varepsilon'_+]_\theta)}(|\nabla u_\varepsilon|, |\nabla u_{\varepsilon'}|)_\theta dx - \langle f, [u_\varepsilon, u_{\varepsilon'}]_\theta \rangle \\ &\leq \int_{\Omega} [\kappa_\varepsilon(|\nabla u_\varepsilon|), \kappa_{\varepsilon'}(|\nabla u_{\varepsilon'}|)]_\theta dx - [\langle f, u_\varepsilon \rangle, \langle f, u_{\varepsilon'} \rangle]_\theta \\ &\leq [\mathcal{J}_\varepsilon(u_\varepsilon), \mathcal{J}_{\varepsilon'}(u_{\varepsilon'})]_\theta, \end{aligned}$$

where we used the triangular inequality and the monotonicity of  $\kappa_\varepsilon(t)$  in  $t$ .  $\square$

The next theorem shows the relation between minimizing  $\mathcal{J}_\varepsilon^s$  and  $\mathcal{J}_\varepsilon$ .

**Theorem 2.37.** *A tuple  $(w, a) \in W_0^{1,2}(\Omega) \times \{a \in L_{\text{loc}}^1(\Omega) : a = \Pi_\varepsilon(a)\}$  minimizes  $\mathcal{J}_\varepsilon^s$  if and only if  $(w, a) = (u_\varepsilon, \Pi_\varepsilon(|\nabla u_\varepsilon|))$ . In particular, there is a unique minimization tuple of  $\mathcal{J}_\varepsilon^s$ .*

*Proof.* Let  $(w_m, a_m) \in W_0^{1,2}(\Omega) \times \{a \in L_{\text{loc}}^1(\Omega) : a = \Pi_\varepsilon(a)\}$  be a minimizing tuple of  $\mathcal{J}_\varepsilon^s$ . We know by Lemma 2.33 that this already implies  $a_m = \Pi_\varepsilon(|\nabla w_m|)$ . Hence for any  $w \in W_0^{1,2}(\Omega)$

$$\begin{aligned} \mathcal{J}_\varepsilon(w) &= \mathcal{J}_\varepsilon^s(w, |\nabla w|) \\ &= \mathcal{J}_\varepsilon^s(w, \Pi_\varepsilon(|\nabla w|)) \\ &\geq \mathcal{J}_\varepsilon^s(w_m, a_m) \\ &= \mathcal{J}_\varepsilon^s(w_m, \Pi_\varepsilon(|\nabla w_m|)) \\ &= \mathcal{J}_\varepsilon(w_m), \end{aligned}$$

so  $w_m$  is minimizer of  $\mathcal{J}_\varepsilon$ . Due to the uniqueness of the minimizer  $\mathcal{J}_\varepsilon$  we get  $w_m = u_\varepsilon$  and therefore  $a_m = \Pi_\varepsilon(|\nabla w_m|) = \Pi_\varepsilon(|\nabla u_\varepsilon|)$ .

On the other hand let  $(w_m, a_m) = (u_\varepsilon, \Pi_\varepsilon(|\nabla u_\varepsilon|))$ . Then, for any  $w \in W_0^{1,2}(\Omega)$  and  $a \in \{\tilde{a} \in L_{\text{loc}}^1(\Omega) : \tilde{a} = \Pi_\varepsilon(\tilde{a})\}$  we have

$$\begin{aligned} \mathcal{J}_\varepsilon^s(w, a) &\geq \mathcal{J}_\varepsilon^s(w, \Pi_\varepsilon(|\nabla w|)) \\ &= \mathcal{J}_\varepsilon^s(w, |\nabla w|) \\ &= \mathcal{J}_\varepsilon(w) \\ &\geq \mathcal{J}_\varepsilon(u_\varepsilon) \\ &= \mathcal{J}_\varepsilon^s(u_\varepsilon, |\nabla u_\varepsilon|) \\ &= \mathcal{J}_\varepsilon^s(u_\varepsilon, \Pi_\varepsilon(|\nabla u_\varepsilon|)) \\ &= \mathcal{J}_\varepsilon^s(w_m, a_m), \end{aligned}$$

so  $(w_m, a_m)$  is the minimizer of  $\mathcal{J}_\varepsilon^s$  on  $W_0^{1,2}(\Omega) \times \{a \in L_{\text{loc}}^1(\Omega) : a = \Pi_\varepsilon(a)\}$ .  $\square$

For a fixed relaxation interval  $\varepsilon$  we may hope that alternatingly minimizing  $\mathcal{J}_\varepsilon^s$  leads to a minimizing tuple of  $\mathcal{J}_\varepsilon^s$  and therefore to a minimizer of  $\mathcal{J}_\varepsilon$ .

First of all we note that  $\kappa_\varepsilon$  not an N-function since  $\kappa_\varepsilon(0) = (\frac{1}{p} - \frac{1}{2})\varepsilon_-^p \neq 0$ . But  $\kappa_\varepsilon \in C^1(\mathbb{R}_{\geq 0})$  with

$$\kappa'_\varepsilon(t) = \begin{cases} \varepsilon_-^{p-2}t & \text{for } t \leq \varepsilon_-, \\ t^{p-1} & \text{for } t \in \varepsilon \text{ and} \\ \varepsilon_+^{p-2}t & \text{for } t \geq \varepsilon_+. \end{cases}$$

Note that the identity

$$\begin{aligned} \frac{\varphi'_\varepsilon(t)}{t} &= \begin{cases} \varepsilon_-^{p-2} & \text{for } t \leq \varepsilon_- \\ t^{p-2} & \text{for } t \in \bar{\varepsilon} \text{ and} \\ \varepsilon_+^{p-2} & \text{for } t \geq \varepsilon_+ \end{cases} \\ &= \Pi_\varepsilon(t)^{p-2} \end{aligned} \tag{2.19}$$

holds. Furthermore, the derivative  $\kappa'_\varepsilon$  satisfies all requirements on  $\varphi'$  as in Definition 2.9, hence

$$\varphi_\varepsilon(t) := \int_0^t \kappa'_\varepsilon(\tau) d\tau$$

is an N-function. Obviously,  $\varphi'_\varepsilon(t) = \kappa'_\varepsilon(t)$  so (2.17) reads as

$$\int_\Omega \frac{\varphi'_\varepsilon(|\nabla u_\varepsilon|)}{|\nabla u_\varepsilon|} \nabla u_\varepsilon \nabla \xi dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega). \tag{2.20}$$

We will use this formulation more often than (2.17) since it fits better to the general  $\varphi$ -Laplace setting.

Furthermore,  $\varphi_\varepsilon(t) = \kappa_\varepsilon(t) - \kappa_\varepsilon(0)$ . Hence,

$$\begin{aligned} \varphi_\varepsilon : \mathbb{R}_{\geq 0} &\rightarrow \mathbb{R} \\ t &\mapsto \begin{cases} \frac{1}{2}\varepsilon_-^{p-2}t^2 & \text{for } t \leq \varepsilon_-, \\ \frac{1}{p}t^p - \left(\frac{1}{p} - \frac{1}{2}\right)\varepsilon_-^p & \text{for } t \in \varepsilon \text{ and} \\ \frac{1}{2}\varepsilon_+^{p-2}t^2 + \left(\frac{1}{p} - \frac{1}{2}\right)(\varepsilon_+^p - \varepsilon_-^p) & \text{for } t \geq \varepsilon_+. \end{cases} \end{aligned} \quad (2.21)$$

As we already know about the importance of the Simonenko indices it is meaningful to calculate them.

**Lemma 2.38.** *For  $\varphi_\varepsilon$  as defined above we get*

$$p^- = \frac{p}{1 - \frac{1}{2}(2-p)\left(\frac{\varepsilon_-}{\varepsilon_+}\right)^p} \quad \text{and} \quad p^+ = 2.$$

In particular,  $\varphi_\varepsilon$  and  $\varphi_\varepsilon^*$  satisfy the  $\Delta_2$  condition.

*Proof.* We calculate the quantities

$$\inf_{0 < t \leq \varepsilon_-} \frac{t\varphi'_\varepsilon(t)}{\varphi_\varepsilon(t)} = \inf_{0 < t \leq \varepsilon_-} \frac{\varepsilon_-^{p-2}t^2}{\frac{1}{2}\varepsilon_-^{p-2}t^2} = 2$$

as well as

$$\begin{aligned} \inf_{t \in \varepsilon} \frac{t\varphi'_\varepsilon(t)}{\varphi_\varepsilon(t)} &= \inf_{t \in \varepsilon} \frac{t^p}{\frac{1}{p}t^p - \left(\frac{1}{p} - \frac{1}{2}\right)\varepsilon_-^p} \\ &= \inf_{t \in \varepsilon} \frac{1}{\frac{1}{p} - \left(\frac{1}{p} - \frac{1}{2}\right)\varepsilon_-^p t^{-p}} \\ &= \frac{p}{1 - \left(1 - \frac{p}{2}\right)\left(\frac{\varepsilon_-}{\varepsilon_+}\right)^p} \end{aligned}$$

and

$$\begin{aligned} \inf_{t \geq \varepsilon_+} \frac{t\varphi'_\varepsilon(t)}{\varphi_\varepsilon(t)} &= \inf_{t \geq \varepsilon_+} \frac{\varepsilon_+^{p-2}t^2}{\frac{1}{2}\varepsilon_+^{p-2}t^2 + \left(\frac{1}{p} - \frac{1}{2}\right)(\varepsilon_+^p - \varepsilon_-^p)} \\ &= \inf_{t \geq \varepsilon_+} \frac{\varepsilon_+^{p-2}}{\frac{1}{2}\varepsilon_+^{p-2} + \left(\frac{1}{p} - \frac{1}{2}\right)(\varepsilon_+^p - \varepsilon_-^p)t^{-2}} \\ &= \frac{2}{1 + \left(\frac{2}{p} - 1\right)\left(1 - \left(\frac{\varepsilon_-}{\varepsilon_+}\right)^p\right)}. \end{aligned}$$

so

$$p^- = \min \left\{ 2, \frac{p}{1 - (1 - \frac{p}{2})(\frac{\varepsilon_-}{\varepsilon_+})^p}, \frac{2}{1 + (\frac{2}{p} - 1)(1 - (\frac{\varepsilon_-}{\varepsilon_+})^p)} \right\}.$$

For our choices of  $p$  and  $\varepsilon$  it is clear that

$$2 \geq \frac{2}{1 + (\frac{2}{p} - 1)(1 - (\frac{\varepsilon_-}{\varepsilon_+})^p)}.$$

Furthermore

$$\begin{aligned} \frac{p}{1 - (1 - \frac{p}{2})(\frac{\varepsilon_-}{\varepsilon_+})^p} &= \frac{2}{1 + (\frac{2}{p} - 1)(1 - (\frac{\varepsilon_-}{\varepsilon_+})^p)} \\ &\iff p + (2 - p)(1 - (\frac{\varepsilon_-}{\varepsilon_+})^p) = 2 - (2 - p)(\frac{\varepsilon_-}{\varepsilon_+})^p \\ &\iff 1 - (\frac{\varepsilon_-}{\varepsilon_+})^p + (\frac{\varepsilon_-}{\varepsilon_+})^p = 1 \end{aligned}$$

which also holds obviously true. Analogue, we get

$$\sup_{0 < t \leq \varepsilon_-} \frac{t\varphi'_\varepsilon(t)}{\varphi_\varepsilon(t)} = \sup_{0 < t \leq \varepsilon_-} \frac{\varepsilon_-^{p-2}t^2}{\frac{1}{2}\varepsilon_-^{p-2}t^2} = 2,$$

as well as

$$\begin{aligned} \sup_{t \in \varepsilon} \frac{t\varphi'_\varepsilon(t)}{\varphi_\varepsilon(t)} &= \sup_{t \in \varepsilon} \frac{t^p}{\frac{1}{p}t^p - (\frac{1}{p} - \frac{1}{2})\varepsilon_-^p} \\ &= \sup_{t \in \varepsilon} \frac{1}{\frac{1}{p} - (\frac{1}{p} - \frac{1}{2})\varepsilon_-^p t^{-p}} \\ &= 2 \end{aligned}$$

and

$$\begin{aligned} \sup_{\frac{1}{\varepsilon} \leq t} \frac{t\varphi'_\varepsilon(t)}{\varphi_\varepsilon(t)} &= \sup_{\frac{1}{\varepsilon} \leq t} \frac{\varepsilon_+^{p-2}t^2}{\frac{1}{2}\varepsilon_+^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})(\varepsilon_+^p - \varepsilon_-^p)} \\ &= \sup_{t \geq \varepsilon_+} \frac{\varepsilon_+^{p-2}}{\frac{1}{2}\varepsilon_+^{p-2} + (\frac{1}{p} - \frac{1}{2})(\varepsilon_+^p - \varepsilon_-^p)t^{-2}} \\ &= 2. \end{aligned}$$

which clearly proofs  $p^+ = 2$ .

That  $\varphi_\varepsilon$  and  $\varphi_\varepsilon^*$  satisfy the  $\Delta_2$  condition is a simple consequence of the Lemmas 2.12 and 2.14.  $\square$

As  $p \leq \frac{p}{1 - \frac{1}{2}(2-p)(\frac{\varepsilon_-}{\varepsilon_+})^p}$  we can use Lemma 2.12 to deduce the following lemma.

**Lemma 2.39.** *For all  $\varepsilon \in \mathbb{R}_{\geq 0}$  and  $s, t \geq 0$ , the inequalities*

$$\min\{s^p, s^2\}\varphi_\varepsilon(t) \leq \varphi_\varepsilon(st) \leq \max\{s^p, s^2\}\varphi_\varepsilon(t) \quad (2.22)$$

*hold.*

Note that it is not surprising that the Simonenko indices of  $\varphi_\varepsilon$  are close to  $p$  for  $\varepsilon_-$  small and  $\varepsilon_+$  large and 2 since  $\varphi_\varepsilon$  has quadratic growth at 0 and  $\infty$  and  $p$ -growth on  $\varepsilon$ .

Up to now we always understood  $\mathcal{J}_\varepsilon$  as an energy defined on  $W_0^{1,2}(\Omega)$ . By the definition of  $\varphi_\varepsilon$  we have

$$\mathcal{J}_\varepsilon(w) = \int_{\Omega} \varphi_\varepsilon(|\nabla w|) dx - |\Omega|\kappa_\varepsilon(0) - \langle f, w \rangle.$$

Hence, it is much more natural to imagine  $\mathcal{J}_\varepsilon : W_0^{1,\varphi_\varepsilon}(\Omega) \rightarrow \mathbb{R}$ . From a set theoretical point of view we have  $W_0^{1,2}(\Omega) = W_0^{1,\varphi_\varepsilon}(\Omega)$  which is due to the quadratic growth of  $\varphi_\varepsilon(t)$  for  $t \geq \varepsilon_+$  and Theorem 2.24.

Now that we know some things about the “right” N-function  $\varphi_\varepsilon$  we want to use it to deduce on of our central equivalences. To do so, we need the following central statements.

Due to the importance of the expression  $A_{\varphi_\varepsilon}$  we write according to Definition 2.18

$$A_\varepsilon(P) := A_{\varphi_\varepsilon}(P) \quad \text{and} \quad V_\varepsilon(P) := V_{\varphi_\varepsilon}(P).$$

By  $\varphi_{\varepsilon,a} := (\varphi_\varepsilon)_a$  we denote the shifted N-function of  $\varphi_\varepsilon$ .

We want to show that  $\varphi_\varepsilon$  satisfies the requirements of Lemma 2.19 uniformly in the relaxation interval  $\varepsilon$  – meaning that the constants do not depend on  $\varepsilon$  – to use it in the later chapters.

This will also imply that all statements care  $\varepsilon = (0, \infty)$  as a special case.

**Lemma 2.40.** *Let  $\varphi_\varepsilon$  be as defined above. Then,*

$$\frac{\varphi'_\varepsilon(s \vee t)}{s \vee t} |t - s| \leq \frac{1}{p-1} |\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)| \leq \frac{1}{p-1} \frac{\varphi'_\varepsilon(s \vee t)}{s \vee t} |t - s|$$

*for all  $s, t \geq 0$ . In particular,  $\varphi'_\varepsilon$  is locally Lipschitz continuous.*



*Proof.* Without loss of generality we assume  $s \leq t$  and distinguish six cases.

$s \leq t \leq \varepsilon_-$ : It follows directly that also  $s \leq \varepsilon$ . Then,

$$\frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{\frac{\varepsilon_-^{p-2}t}{t}(t-s)}{\varepsilon_-^{p-2}(t-s)} = 1.$$

$s \leq \varepsilon_- \leq t \leq \varepsilon_+$ : Now  $\varepsilon_-^{p-2} \geq t^{p-2}$  implies  $t^{p-2}(t-s) \geq t^{p-1} - \varepsilon_-^{p-2}s$  and hence

$$\frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{t^{p-2}(t-s)}{t^{p-1} - \varepsilon_-^{p-2}s} \geq 1.$$

On the other hand,

$$\begin{aligned} \frac{t^{p-2}(t-s)}{t^{p-1} - \varepsilon_-^{p-2}s} \leq \frac{1}{p-1} &\iff t^{p-1} - st^{p-2} \leq \frac{1}{p-1}(t^{p-1} - \varepsilon_-^{p-2}s) \\ &\iff s \leq \frac{\frac{2-p}{p-1}t^{p-1}}{\frac{1}{p-1}\varepsilon_-^{p-2} - t^{p-2}} \\ &\iff s \leq \frac{(2-p)t^{p-1}}{\varepsilon_-^{p-2} - (p-1)t^{p-2}} =: h_\varepsilon(t). \end{aligned}$$

The function  $h_\varepsilon$  is increasing since  $\varepsilon_-^{p-2} \geq t^{p-2}$  implies

$$\begin{aligned} h'_\varepsilon(t) &= \frac{(2-p)(p-1)t^{p-2}(\varepsilon_-^{p-2} - (p-1)t^{p-2}) + (p-1)(p-2)t^{p-3}(2-p)t^{p-1}}{(\varepsilon_-^{p-2} - (p-1)t^{p-2})^2} \\ &= \frac{(2-p)(p-1)t^{p-2}(\varepsilon_-^{p-2} - (p-1)t^{p-2} + (p-2)t^{p-2})}{(\varepsilon_-^{p-2} - (p-1)t^{p-2})^2} \\ &= \frac{(2-p)(p-1)t^{p-2}(\varepsilon_-^{p-2} - t^{p-2})}{(\varepsilon_-^{p-2} - (p-1)t^{p-2})^2} \\ &\geq 0 \end{aligned}$$

In particular,  $h_\varepsilon(t) \geq h_\varepsilon(\varepsilon_-) = \varepsilon_-$ . But  $s \leq \varepsilon_-$  by assumption.

$s \leq \varepsilon_- \leq \varepsilon_+ \leq t$ : In this case we use  $-\varepsilon_+^{p-2} \geq -\varepsilon_-^{p-2}$  to get

$$\frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{\varepsilon_+^{p-2}(t-s)}{\varepsilon_+^{p-2}t - \varepsilon_-^{p-2}s} \geq 1.$$

On the other hand we write  $\varepsilon_- = \sigma\varepsilon_+$  and  $s = \theta t$ . The assumptions on this case imply  $\theta \leq \sigma \leq 1$ . Hence

$$\begin{aligned} \frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} &= \frac{\varepsilon_+^{p-2}(t-s)}{\varepsilon_+^{p-2}t - \varepsilon_-^{p-2}s} \\ &= \frac{\varepsilon_+^{p-2}t(1-\sigma)}{\varepsilon_+^{p-2}t(1-\sigma^{p-2}\theta)} \\ &\leq \frac{1-\theta}{1-\theta^{p-1}} \end{aligned}$$

For  $p < 2$  this is increasing in  $\theta$ . Therefore we get by L'Hôpital's rule

$$\begin{aligned} \frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} &\leq \lim_{\theta \rightarrow 1} \frac{1-\theta}{1-\theta^{p-1}} \\ &= \lim_{\theta \rightarrow 1} \frac{-1}{-(p-1)\theta^{p-2}} \\ &= \frac{1}{p-1}. \end{aligned}$$

$\varepsilon_- \leq s \leq t \leq \varepsilon_+$ : With  $s \leq t$  we get

$$\frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{t^{p-2}(t-s)}{t^{p-1} - s^{p-1}} \geq 1.$$

By concavity of  $t \mapsto t^{p-1}$  we get  $s^{p-1} \leq t^{p-1} - (p-1)t^{p-2}(t-s)$  and so

$$\frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{t^{p-2}(t-s)}{t^{p-1} - s^{p-1}} \leq \frac{1}{p-1}.$$

$\varepsilon_- \leq s \leq \varepsilon_+ \leq t$ : Wit  $s \leq \varepsilon_+$  we get  $s^{p-1} \geq \varepsilon_+^{p-2}s$  and so

$$\frac{\frac{\varphi'_\varepsilon(s\sqrt{t})}{s\sqrt{t}}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{\varepsilon_+^{p-2}(t-s)}{\varepsilon_+^{p-2}t - s^{p-1}} \geq 1.$$

$$\varepsilon^{2-p}(t-s) \sim \varepsilon^{2-p}t - s^{p-1}.$$

On the other hand we have

$$\begin{aligned} \frac{\varepsilon_+^{p-2}(t-s)}{\varepsilon_+^{p-2}t - s^{p-1}} \leq \frac{1}{p-1} &\iff (p-1)\varepsilon_+^{p-2}(t-s) \leq \varepsilon_+^{p-2}t - s^{p-1} \\ &\iff s^{p-1} - (p-1)\varepsilon_+^{p-2}s \leq (2-p)\varepsilon_+^{p-2}t. \end{aligned}$$

This holds true since  $\frac{\partial}{\partial s}(s^{p-1} - (p-1)\varepsilon_+^{p-2}s) = (p-1)(s^{p-2} - \varepsilon_+^{p-2}) \geq 0$  and so

$$s^{p-1} - (p-1)\varepsilon_+^{p-2}s \leq (2-p)\varepsilon_+^{p-1} = (2-p)\varepsilon_+^{p-2}\varepsilon_+ \leq (2-p)\varepsilon_+^{p-2}t.$$

$\varepsilon_+ \leq s \leq t$ : We get

$$\frac{\frac{\varphi'_\varepsilon(s \vee t)}{s \vee t}|t-s|}{|\varphi'_\varepsilon(t) - \varphi'_\varepsilon(s)|} = \frac{\varepsilon^{2-p}(t-s)}{\varepsilon^{2-p}t - \varepsilon^{2-p}s} = 1.$$

Building the minima and maxima of all constants we obtain the statement.  $\square$

With that we will first of all show that  $A_\varepsilon$  admits so called  $\varphi_\varepsilon$ -structure, which is stated in the next lemma.

**Lemma 2.41.** *The N-function  $\varphi_\varepsilon$  satisfies the requirements of Lemma 2.19 in the sense that*

$$(A_\varepsilon(P) - A_\varepsilon(Q))(P - Q) \approx \frac{\varphi'_\varepsilon(|P| + |P - Q|)}{|P| + |P - Q|} |P - Q|^2.$$

*In particular,*

$$\begin{aligned} (A_\varepsilon(P) - A_\varepsilon(Q))(P - Q) &\approx |V_\varepsilon(P) - V_\varepsilon(Q)|^2 \\ &\approx \varphi_{\varepsilon, |P|}(|P - Q|) \\ &\approx \frac{\varphi'_\varepsilon(|P| \vee |Q|)}{|P| \vee |Q|} |P - Q|^2. \end{aligned}$$

*Proof.* We plug  $t = a + \tilde{t}$  and  $s = a$  in Lemma 2.40 to deduce that

$$\frac{\varphi'_\varepsilon(a+t)t}{a+t} \approx \varphi'_\varepsilon(a+t) - \varphi'_\varepsilon(a).$$

As shown in [RD07, Lemma 6.14] this is a sufficient condition for the first statement. The second estimates are a direct consequence of Lemma 2.19.  $\square$

The next theorem is fundamental for the proofs in the next two chapters. It allows us to switch between energy differences and differences of equations.

**Theorem 2.42.** *Let  $u_\varepsilon$  be the minimizer of  $\mathcal{J}_\varepsilon$  and  $w \in W_0^{1, \varphi_\varepsilon}(\Omega)$ . Then,*

$$\mathcal{J}_\varepsilon(w) - \mathcal{J}_\varepsilon(u_\varepsilon) \leq \int_{\Omega} (A_\varepsilon(\nabla w) - A_\varepsilon(\nabla u_\varepsilon)) \nabla(w - u_\varepsilon) dx \lesssim \mathcal{J}_\varepsilon(w) - \mathcal{J}_\varepsilon(u_\varepsilon).$$

*Proof.* We define

$$\begin{aligned} f &: [0, 1] \rightarrow \mathbb{R} \\ t &\mapsto \mathcal{J}_\varepsilon(tw + (1-t)u_\varepsilon). \end{aligned}$$

Since  $\varphi_\varepsilon$  is in  $C^1(\mathbb{R}_{\geq 0})$ ,  $f$  has the same properties with

$$\begin{aligned} f'(t) &= \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla(tw + (1-t)u_\varepsilon)|)}{|\nabla(tw + (1-t)u_\varepsilon)|} \nabla(tw + (1-t)u_\varepsilon) \cdot \nabla(w - u_\varepsilon) dx - \langle f, w - u_\varepsilon \rangle \\ &= \int_{\Omega} A_\varepsilon(\nabla(tw + (1-t)u_\varepsilon)) \cdot \nabla(w - u_\varepsilon) dx - \langle f, w - u_\varepsilon \rangle. \end{aligned}$$

In particular,

$$f'(1) = \int_{\Omega} A_\varepsilon(\nabla w) \cdot \nabla(w - u) dx - \langle f, w - u \rangle.$$

The convexity of  $f$  implies  $f(0) \geq f(1) - f'(1)$ . With the equation for  $u_\varepsilon$  (see (2.20) on page 29) we get

$$\begin{aligned} \mathcal{J}_\varepsilon(w) - \mathcal{J}_\varepsilon(u_\varepsilon) &= f(1) - f(0) \\ &\leq f'(1) \\ &= \int_{\Omega} A_\varepsilon(\nabla w) \cdot \nabla(w - u) dx - \langle f, w - u \rangle \\ &= \int_{\Omega} A_\varepsilon(\nabla w) \cdot \nabla(w - u) dx - \int_{\Omega} A_\varepsilon(\nabla u_\varepsilon) \nabla(w - u_\varepsilon) dx \\ &= \int_{\Omega} (A_\varepsilon(\nabla w) - A_\varepsilon(\nabla u_\varepsilon)) \nabla(w - u_\varepsilon) dx. \end{aligned}$$

Due to the definition of the Simonenko indices and Lemma 2.39 we have for any  $s \in [0, 1]$

$$\varphi'_\varepsilon(st) \geq \frac{p\varphi_\varepsilon(st)}{st} \geq \frac{ps^2\varphi_\varepsilon(t)}{st} \geq \frac{p}{2}s\varphi'_\varepsilon(t).$$

With this, the first equivalence of Lemma 2.41 and the monotonicity of  $\frac{\varphi'_\varepsilon(t)}{t}$  we

can deduce

$$\begin{aligned}
& (A_\varepsilon(\nabla(sw + (1-s)u_\varepsilon)) - A_\varepsilon(u_\varepsilon))\nabla(w - u_\varepsilon) \\
&= (A_\varepsilon(u_\varepsilon) - A_\varepsilon(\nabla(sw + (1-s)u_\varepsilon)))\nabla(u_\varepsilon - w) \\
&\gtrsim \frac{\varphi'_\varepsilon(|\nabla u_\varepsilon| + |\nabla u_\varepsilon - s\nabla w - (1-s)\nabla u_\varepsilon|)}{|\nabla u_\varepsilon| + |\nabla u_\varepsilon - s\nabla w - (1-s)\nabla u_\varepsilon|} |\nabla u_\varepsilon - s\nabla w - (1-s)\nabla u_\varepsilon|^2 \\
&= \frac{\varphi'_\varepsilon(|\nabla u_\varepsilon| + s|\nabla(w - u_\varepsilon)|)}{|\nabla u_\varepsilon| + s|\nabla(w - u_\varepsilon)|} |\nabla(w - u_\varepsilon)|^2 s^2 \\
&\geq \frac{\varphi'_\varepsilon(s|\nabla u_\varepsilon| + s|\nabla(w - u_\varepsilon)|)}{s|\nabla u_\varepsilon| + s|\nabla(w - u_\varepsilon)|} |\nabla(w - u_\varepsilon)|^2 s^2 \\
&\gtrsim \frac{\varphi'_\varepsilon(|\nabla u_\varepsilon| + |\nabla(w - u_\varepsilon)|)}{|\nabla u_\varepsilon| + |\nabla(w - u_\varepsilon)|} |\nabla(w - u_\varepsilon)|^2 s^2 \\
&\gtrsim s^2 (A_\varepsilon(\nabla w) - A_\varepsilon(\nabla u_\varepsilon))\nabla(w - u_\varepsilon).
\end{aligned}$$

This implies

$$\begin{aligned}
\mathcal{J}_\varepsilon(w) - \mathcal{J}_\varepsilon(u_\varepsilon) &= f(1) - f(0) = \int_0^1 f'(s) ds \\
&= \int_0^1 \int_\Omega (A_\varepsilon(\nabla(sv + (1-s)u_\varepsilon)) - A_\varepsilon(u_\varepsilon)) \cdot \nabla(v - u_\varepsilon) dx ds \\
&\gtrsim \int_0^1 \int_\Omega s^2 (A_\varepsilon(\nabla w) - A_\varepsilon(\nabla u_\varepsilon))\nabla(w - u_\varepsilon) dx ds \\
&\gtrsim \int_\Omega (A_\varepsilon(\nabla w) - A_\varepsilon(\nabla u_\varepsilon))\nabla(w - u_\varepsilon) dx. \quad \square
\end{aligned}$$



### 3 Convergence in the Relaxation Parameter

As explained in the last chapter, we introduced a relaxation interval into the  $p$ -Poisson equation for our iteration to be well defined. As suggested by Theorem 2.37 we may hope that this leads to a minimizer of  $\mathcal{J}_\varepsilon$ . However, to solve the  $p$ -Poisson equation we want to minimize  $\mathcal{J}$ . It is clear that the minimizer  $u_\varepsilon$  of  $\mathcal{J}_\varepsilon$  and the minimizer  $u$  of  $\mathcal{J}$  may not coincide. We will describe the relation between  $u_\varepsilon$  and  $\varepsilon$  in this chapter.

The right space to measure the distance of  $u_\varepsilon$  and  $u$  is  $W_0^{1,p}(\Omega)$  and not  $W_0^{1,2}(\Omega)$  or  $W_0^{1,\varphi_\varepsilon}(\Omega)$  since  $u \notin W_0^{1,2}(\Omega) = W_0^{1,\varphi_\varepsilon}(\Omega)$  for certain  $f \in (W_0^{1,p}(\Omega))^*$ . A further disadvantage of  $W_0^{1,\varphi_\varepsilon}(\Omega)$  is that its norm depends on  $\varepsilon$ .

But first of all we will prove one important property of  $\mathcal{J}_\varepsilon$ , namely its monotonicity in  $\varepsilon$ .

**Lemma 3.1.** *Let  $\varepsilon \subseteq \delta \subseteq (0, \infty)$ . Then for all  $w \in W_0^{1,\varphi_\varepsilon}(\Omega)$  we have*

$$\mathcal{J}_\varepsilon(w) \geq \mathcal{J}_\delta(w) \geq \mathcal{J}(w).$$

*Proof.* Due to the definition of  $\mathcal{J}_\varepsilon$  it suffices to show  $\kappa_\varepsilon(t) \geq \kappa_\delta(t) \geq \varphi(t) := \frac{1}{p}t^p$ .

$t \leq \delta_- \leq \varepsilon_-$ : We easily see

$$\kappa'_\varepsilon(t) = \varepsilon_-^{p-2}t \leq \delta_-^{p-2}t = \kappa'_\delta(t) \leq t^{p-1} = \varphi'(t).$$

$\delta_- \leq t \leq \varepsilon_-$ : In this case

$$\kappa'_\varepsilon(t) = \varepsilon_-^{p-2}t \leq t^{p-1} = \kappa'_\delta(t) = \varphi'(t).$$

$\varepsilon_+ \leq t \leq \delta_+$ : Here

$$\kappa'_\varepsilon(t) = \varepsilon_+^{p-2}t \geq t^{p-1} = \kappa'_\delta(t) = \varphi'(t).$$

$\varepsilon_+ \leq \delta_+ \leq t$ : Under these assumptions we get

$$\kappa'_\varepsilon(t) = \varepsilon_+^{p-2}t \geq \delta_+^{p-2}t = \kappa'_\delta(t) \geq t^{p-1} = \varphi'(t).$$

Combining the first two cases imply the estimates  $\kappa'_\varepsilon(t) \leq \kappa'_\delta(t) \leq \varphi'(t)$  for  $t \in [0, \varepsilon_-]$ . Additionally we have  $\kappa_\varepsilon(\varepsilon_-) = \kappa_\delta(\varepsilon_-) = \varphi(\varepsilon_-)$ . Hence for those choices of  $t$

$$\begin{aligned} \kappa_\varepsilon(t) &= \kappa_\varepsilon(\varepsilon_-) - \int_t^{\varepsilon_-} \kappa'_\varepsilon(\tau) d\tau \\ &\geq \kappa_\delta(\varepsilon_-) - \int_t^{\varepsilon_-} \kappa'_\delta(\tau) d\tau \\ &= \kappa_\delta(t). \end{aligned}$$

For  $t \in [\varepsilon_-, \varepsilon_+]$  we have  $\kappa_\varepsilon(t) = \kappa_\delta(t) = \varphi(t)$  by definition.

If  $t \in [\varepsilon_+, \infty)$  we have  $\kappa'_\varepsilon(t) \geq \kappa'_\delta(t) \geq \varphi'(t)$  and  $\kappa_\varepsilon(\varepsilon_+) = \kappa_\delta(\varepsilon_+) = \varphi(\varepsilon_+)$  and therefore

$$\begin{aligned} \kappa_\varepsilon(t) &= \kappa_\varepsilon(\varepsilon_+) + \int_{\varepsilon_+}^t \kappa'_\varepsilon(\tau) d\tau \\ &\geq \kappa_\delta(\varepsilon_+) + \int_{\varepsilon_+}^t \kappa'_\delta(\tau) d\tau \\ &= \kappa_\delta(t). \end{aligned}$$

The proof of  $\kappa_\delta(t) \geq \varphi(t)$  is analogue to  $\kappa_\varepsilon(t) \geq \kappa_\delta(t)$ . □

To show the convergence results in this section we will need a tool which is based on another tool known from the harmonic analysis, see for example [SM93, I, §1].

**Definition 3.2.** For  $w \in L^1(\mathbb{R}^d)$  we define the *Hardy-Littlewood maximal function*

$$M(w)(x) := \sup_{r>0} \int_{B_r(x)} |w(y)| dy := \sup_{r>0} \frac{1}{|B_r(0)|} \int_{B_r(x)} |w(y)| dy.$$



For  $p \in (1, \infty]$  the Hardy-Littlewood maximal operator  $M$  is continuous from  $L^p(\mathbb{R}^d)$  to  $L^p(\mathbb{R}^d)$ , i.e. there is a constant  $c_p$  such that for all  $w \in L^p(\mathbb{R}^d)$

$$\|M(w)\|_{L^p(\mathbb{R}^d)} \leq c_p \|w\|_{L^p(\mathbb{R}^d)}. \quad (3.1)$$

The Hardy-Littlewood maximal function is of weak type (1,1), i.e. there exists a constant  $c > 0$  such that for all  $w \in L^1(\mathbb{R}^d)$

$$\sup_{\lambda > 0} \lambda |\{M(w) > \lambda\}| \leq c \|w\|_{L^1(\mathbb{R}^d)}. \quad (3.2)$$

Now, for  $w \in W_0^{1,p}(\Omega) \subset W_0^{1,1}(\Omega)$  we may assume  $w \in W_0^{1,1}(\mathbb{R}^d)$  by extending  $w$  by zero on  $\Omega^c$  (see for example [AF03, Theorem 5.29]) since we assumed  $\Omega$  to be bounded and having Lipschitz boundary.

**Definition 3.3.** For  $w \in W_0^{1,p}(\mathbb{R}^d)$  and  $\lambda \in (0, \infty)$  we define the *bad set*

$$\mathcal{O}_\lambda(w) := \{x \in \mathbb{R}^d : M(|\nabla w|)(x) > \lambda\}.$$

Note that we will write  $\mathcal{O}_\lambda$  instead of  $\mathcal{O}_\lambda(w)$  if the choice of  $w$  is clear.

We will use the Lipschitz truncation provided by [DKS13, Subsection 3.5].

**Theorem 3.4.** *There is a constant  $c > 1$  such that for any  $\lambda > 0$  and  $w \in W_0^{1,p}(\Omega)$  there is  $w_\lambda \in W_0^{1,\infty}(\Omega)$  with the following properties:*

- (i)  $\{w \neq w_\lambda\} \subset \mathcal{O}_\lambda(w) \cap \Omega$ ,
- (ii)  $\|w_\lambda\|_{L^p(\Omega)} \lesssim \|w\|_{L^p(\Omega)}$ ,
- (iii)  $\|\nabla w_\lambda\|_{L^p(\Omega)} \lesssim \|\nabla w\|_{L^p(\Omega)}$ ,
- (iv)  $|\nabla w_\lambda| \leq c\lambda \chi_{\mathcal{O}_\lambda(w) \cap \Omega} + |\nabla w| \chi_{\mathcal{O}_\lambda(w)^c \cap \Omega} \leq c\lambda$  almost everywhere on  $\Omega$  and
- (v)  $\|\nabla(w - w_\lambda)\|_{L^p(\Omega)} \lesssim \|\nabla w\|_{L^p(\mathcal{O}_\lambda(w))}$ .

The constant in (iv) is mentioned explicitly since it will be of importance later. Note that (v) was not proven in this work but is a consequence of (i) and (iii).

### 3.1 $\Gamma$ -Convergence

First of all, we will use the concept of  $\Gamma$ -convergence to show  $u_\varepsilon \xrightarrow{\varepsilon \rightarrow (0, \infty)} u$ , although we will prove the same result later again. We use a definition that is only equivalent to the general definition in the literature, if the underlying topological space satisfies the first axiom of countability (see [dM93, Proposition 8.1]).

**Definition 3.5.** Let  $X$  be a topological space,  $F_n, F : X \rightarrow \overline{\mathbb{R}}$ . We say  $F$  is the  $\Gamma$ -limit of  $(F_n)$  or  $F = \Gamma\text{-}\lim_{n \rightarrow \infty} F_n$  iff if the lim inf-condition (i) and the lim sup-condition (ii) are satisfied.

- (i) If  $x_n \rightarrow x$ , then  $F(x) \leq \liminf_{n \rightarrow \infty} F_n(x_n)$ .
- (ii) For any  $x \in X$  there is  $(x_n)$  with  $x_n \xrightarrow{n \rightarrow \infty} x$  and  $F(x) \geq \limsup_{n \rightarrow \infty} F_n(x_n)$ .

The sequence in (ii) is called *recovering sequence*.

Note that if (i) holds we have that (ii) is equivalent to  $F(x) = \lim_{n \rightarrow \infty} F_n(x_n)$  since

$$F(x) \geq \limsup_{n \rightarrow \infty} F_n(x_n) \geq \liminf_{n \rightarrow \infty} F_n(x_n) \geq F(x).$$

For the rest of this chapter, let  $(\varepsilon_n)$  be a monotone sequence of intervals converging to  $(0, \infty)$  in the sense that  $\varepsilon_{n+1} \supseteq \varepsilon_n$  for each  $n$  and  $\lim_{n \rightarrow \infty} \varepsilon_{-,n} = 0$  and  $\lim_{n \rightarrow \infty} \varepsilon_{+,n} = \infty$ .

Although we are aware that the weak topology on  $W_0^{1,p}(\Omega)$  is not first countable we will show that, according to our definition,  $\mathcal{J} = \Gamma\text{-}\lim_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}$  on  $W_0^{1,p}(\Omega)$  endowed with the strong as well as the weak topology. Indeed, there is no canonical statement that  $\Gamma$ -convergence with respect to one topology implies  $\Gamma$ -convergence with respect to the other topology. This is due to the fact that the lim inf-condition is easier to prove for the strong convergence, whereas the lim sup-condition is easier to prove for the weak convergence.

It remains to extend  $\mathcal{J}_\varepsilon$  via

$$\mathcal{J}_\varepsilon : W_0^{1,p}(\Omega) \rightarrow \overline{\mathbb{R}}$$

$$w \mapsto \begin{cases} \mathcal{J}_\varepsilon(w) & \text{for } w \in W_0^{1,2}(\Omega) \text{ and} \\ +\infty & \text{for } w \in W_0^{1,p}(\Omega) \setminus W_0^{1,2}(\Omega). \end{cases}$$

where we wrote  $W_0^{1,2}(\Omega)$  instead of  $W_0^{1,\varphi_\varepsilon}(\Omega)$  since this is more robust under the process  $\varepsilon_{+,n} \xrightarrow{n \rightarrow \infty} \infty$  and at that point only the set theoretical point of view is important – the topology is induced by  $W_0^{1,p}(\Omega)$  anyway.

**Theorem 3.6.**  $\mathcal{J} = \Gamma\text{-}\lim_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}$  with respect to both, the strong and the weak topology on  $W_0^{1,p}(\Omega)$ .

*Proof.* to prove the lim inf-condition let  $w, w_n \in W_0^{1,p}(\Omega)$  with  $w_n \rightharpoonup w$  (which is also sufficed if  $w_n \rightarrow w$ ). Then, using the monotonicity of  $\mathcal{J}_{\varepsilon_n}$  and the weak lower semi-continuity of the norm we can directly estimate

$$\liminf_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}(w_n) \geq \liminf_{n \rightarrow \infty} \mathcal{J}(w_n) \geq \mathcal{J}(w).$$

It remains to show the existence of a recovery sequence that satisfies the lim sup-condition with respect to the strong topology. We choose a sequence  $(\lambda_n)_n \subset (0, \infty)$  with  $\lambda_n \xrightarrow{n \rightarrow \infty} \infty$  satisfying  $c\lambda_n \leq \varepsilon_{+,n}$  where  $c$  is the constant of Theorem 3.4, (iv). We proof that the sequence of Lipschitz truncations  $w_n := w_{\lambda_n}$  is an admissible recovering sequence for each function  $w \in W_0^{1,p}(\Omega)$ , so we need to show

$$\mathcal{J}_{\varepsilon_n}(w_n) - \mathcal{J}(w) = \int_{\Omega} \kappa_{\varepsilon_n}(|\nabla w_n|) - \frac{1}{p}|\nabla w|^p dx - \langle f, w_n - w \rangle \xrightarrow{n \rightarrow \infty} 0. \quad (3.3)$$

We split the domain of integration into  $\mathcal{O}_{\lambda_n} \cap \Omega$  and  $\mathcal{O}_{\lambda_n}^c \cap \Omega$ . For the latter one we have  $\nabla w_n = \nabla w$ ,  $|\nabla w_n| \leq c\lambda_n \leq \varepsilon_{+,n}$  and for  $t \in [\varepsilon_{-,n}, \varepsilon_{+,n}]$  we have  $\kappa_{\varepsilon_n}(t) = \frac{1}{p}t^p$ , so

$$\begin{aligned} \left| \int_{\mathcal{O}_{\lambda_n}^c \cap \Omega} \kappa_{\varepsilon_n}(|\nabla w_n|) - \frac{1}{p}|\nabla w|^p dx \right| &= \left| \int_{\mathcal{O}_{\lambda_n}^c \cap \Omega} \kappa_{\varepsilon_n}(|\nabla w_n|) - \frac{1}{p}|\nabla w_n|^p dx \right| \\ &= \int_{\mathcal{O}_{\lambda_n}^c \cap \{|\nabla w_n| < \varepsilon_n\} \cap \Omega} \kappa_{\varepsilon_n}(|\nabla w_n|) - \frac{1}{p}|\nabla w_n|^p dx \\ &\leq \int_{\Omega} \kappa_{\varepsilon}(\varepsilon_{-,n}) dx \\ &\leq \frac{1}{p}|\Omega|\varepsilon_{-,n}^p \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

On the bad set we have  $|\nabla w_n| \lesssim \lambda < M(|\nabla w|)$ . From (3.1) we know that  $M(|\nabla w|) \in L^p(\Omega)$ . Beside that (3.2) ensures

$$|\mathcal{O}_{\lambda_n}| \lesssim \frac{\|\nabla w\|_{L^1(\Omega)}}{\lambda_n} \xrightarrow{n \rightarrow \infty} 0.$$

Therefore,

$$\begin{aligned}
\left| \int_{\mathcal{O}_{\lambda_n} \cap \Omega} \kappa_{\varepsilon_n}(|\nabla w_n|) - \frac{1}{p} |\nabla w|^p dx \right| &\leq \int_{\mathcal{O}_{\lambda_n} \cap \Omega} \kappa_{\varepsilon_n}(|\nabla w_n|) + \frac{1}{p} |\nabla w|^p dx \\
&\lesssim \int_{\mathcal{O}_{\lambda_n} \cap \Omega} \lambda^p + |\nabla w|^p dx \\
&\leq \int_{\mathcal{O}_{\lambda_n} \cap \Omega} \underbrace{(M(|\nabla w|))^p + |\nabla w|^p}_{\in L^1(\Omega)} dx \xrightarrow{n \rightarrow \infty} 0
\end{aligned}$$

by dominated convergence theorem. By Theorem 3.4, (v) we deduce that  $w_n \rightarrow w$  in  $W_0^{1,p}(\Omega)$ . Hence,  $w_n \rightharpoonup w$  and therefore,  $\langle f, w_n - w \rangle \xrightarrow{n \rightarrow \infty} 0$ . Altogether, this yields (3.3), so the lim sup-condition.  $\square$

With the  $\Gamma$ -convergence we can prove the following corollary.

**Corollary 3.7.**  $\lim_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}(u_{\varepsilon_n}) = \mathcal{J}(u)$ .

*Proof.* Let  $(w_n)$  be the recovering sequence of  $u$  according to Theorem 3.6. Then, with the aid of Lemma 3.1 we can directly estimate

$$\begin{aligned}
\mathcal{J}(u) &= \lim_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}(w_n) \\
&\geq \lim_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}(u_{\varepsilon_n}) \\
&\geq \liminf_{n \rightarrow \infty} \mathcal{J}_{\varepsilon_n}(u_{\varepsilon_n}) \\
&\geq \liminf_{n \rightarrow \infty} \mathcal{J}(u_{\varepsilon_n}) \\
&\geq \mathcal{J}(u).
\end{aligned}
\quad \square$$

The functional  $\mathcal{J}$  carries norm structure. Hence, we can get the first convergence result without using the whole theory of  $\Gamma$ -convergence.

**Corollary 3.8.**  $u_\varepsilon \xrightarrow{\varepsilon \rightarrow (0, \infty)} u$  in  $W_0^{1,p}(\Omega)$ .

*Proof.* It is easy to see that for any N-function  $\varphi$  that if it is shifted by the argument this yields  $\varphi_t(t) = \int_0^t \frac{\varphi'(t \vee s)}{t \vee s} s ds = \frac{1}{2} \varphi'(t) t \leq c \varphi(t)$ . Now, we use Lemma 2.20 for  $\varphi(t) := \frac{1}{p} t^p$  to get for any  $\delta > 0$  that

$$\begin{aligned}
\varphi(t) &\leq c_\delta \varphi_{|P|}(t) + \frac{\delta}{c} \varphi_{|P|}(|P|) \\
&\leq c_\delta \varphi_{|P|}(t) + \delta \varphi(|P|).
\end{aligned}$$

With  $\varphi(t) := \frac{1}{p}t^p$  we get by Lemma 2.41 and Lemma 3.1 that

$$\begin{aligned} \int_{\Omega} |\nabla(u_{\varepsilon} - u)|^p dx &\lesssim c_{\delta} \int_{\Omega} \varphi_{|\nabla u|}(|\nabla(u_{\varepsilon} - u)|) dx + \delta \int_{\Omega} |\nabla u|^p dx \\ &\lesssim \mathcal{J}(u_{\varepsilon}) - \mathcal{J}(u) + \delta \int_{\Omega} |\nabla u|^p dx \\ &\leq \mathcal{J}_{\varepsilon}(u_{\varepsilon}) - \mathcal{J}(u) + \delta \int_{\Omega} |\nabla u|^p dx. \end{aligned}$$

Now choosing  $\delta, \varepsilon_-$  sufficiently small and  $\varepsilon_+$  sufficiently large yields that the gradient norm  $\|\nabla(u_{\varepsilon} - u)\|_{L^p(\Omega)} \approx \|u - u_{\varepsilon}\|_{W_0^{1,p}(\Omega)}$  is arbitrary small.  $\square$

There are also two other approaches leading to the convergence result as stated in Corollary 3.8. We will just sketchy note them here. Note that we need  $\Gamma$ -convergence with respect to the weak topology for this argumentation since the family  $(\mathcal{J}_{\varepsilon})$  is equi-coercive (which we will need afterwards) if and only if there is a lower semi-continuous and coercive  $\Phi : W_0^{1,p}(\Omega) \rightarrow \overline{\mathbb{R}}$  with  $\Phi \leq \mathcal{J}_{\varepsilon}$ . This is satisfied for  $\Phi = \mathcal{J}$  but only when  $W_0^{1,p}(\Omega)$  is endowed with the weak topology as it is well known that bounded sets in  $L^p$  are only weakly compact.

The first approach is showing that  $\mathcal{J}_{\varepsilon}$  is decreasing and converges pointwisely to the functional  $\tilde{\mathcal{J}}$ , where  $\tilde{\mathcal{J}} = \mathcal{J}$  on  $W_0^{1,2}(\Omega)$  and  $\tilde{\mathcal{J}} \equiv \infty$  on  $W_0^{1,p}(\Omega) \setminus W_0^{1,2}(\Omega)$ . This implies by [dM93, Proposition 5.7] that  $\Gamma\text{-}\lim_{\varepsilon \rightarrow (0,\infty)} \mathcal{J}_{\varepsilon} = \text{sc}^- \tilde{\mathcal{J}}$ . Then one can show that with respect to the weak topology we have  $\text{sc}^- \tilde{\mathcal{J}} = \mathcal{J}$ .

The other approach deals with the sequential characterization as we stated it. But to be precise, this characterization is only equivalent to  $\Gamma$ -convergence if the underlying topological space satisfies the first axiom of countability – which  $W_0^{1,p}(\Omega)$  does not with respect to the weak topology. One can get around this problem by choosing a bounded  $X \subset W_0^{1,p}(\Omega)$  carrying all minimizers and all recovering sequences. Then, the induced weak topology on  $X$  is metrizable (see [dM93, Proposition 8.7]) which implies the first axiom of countability.

After performing one of those two approaches one uses the equi-coercivity and the uniqueness of the minimizer of  $\mathcal{J}$  to deduce that  $u_{\varepsilon} \rightarrow u$  in  $W_0^{1,p}(\Omega)$  directly (see [dM93, Proposition 7.24]). Now one can use  $\mathcal{J}_{\varepsilon}(u_{\varepsilon}) \rightarrow \mathcal{J}(u)$  to get the convergence  $\|\nabla u_{\varepsilon}\|_{L^p(\Omega)} \rightarrow \|\nabla u\|_{L^p(\Omega)}$  – at this point we used again that  $\mathcal{J}$  carries norm structure. But  $\|\nabla \cdot\|_{L^p(\Omega)}$  is a uniformly convex norm on  $W_0^{1,p}(\Omega)$ , so weak convergence and the convergence of the norm imply  $u_{\varepsilon} \rightarrow u$  in  $W_0^{1,p}(\Omega)$ .

## 3.2 Error Bounds

In this section we want to overcome that we only achieved qualitative results in the last section. As we already have used in the last section we have

$$\mathcal{J}(u_\varepsilon) - \mathcal{J}(u) \leq \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u)$$

so it obviously suffices to show bounds for the latter term.

We will immediately start with the following estimate.

**Lemma 3.9.** *For all  $\lambda \leq \frac{\varepsilon_+}{c}$  with  $c$  from Theorem 3.4 we have*

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) \lesssim \varepsilon_-^p + \int_{\mathcal{O}_\lambda(u)} |\nabla u|^p dx. \quad (3.4)$$

*Proof.* Let  $\lambda > 0$  be as required and let  $T_\lambda u$  be the Lipschitz truncation of  $u$ . Since  $u_\varepsilon$  minimizes  $\mathcal{J}_\varepsilon$  and due to the equation for  $u$  (2.5) we may estimate

$$\begin{aligned} \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) &\leq \mathcal{J}_\varepsilon(T_\lambda u) - \mathcal{J}(u) \\ &= \int_{\Omega} \kappa_\varepsilon(|\nabla T_\lambda u|) - \frac{1}{p} |\nabla u|^p dx - \langle f, T_\lambda u - u \rangle \\ &= \int_{\Omega} \kappa_\varepsilon(|\nabla T_\lambda u|) - \frac{1}{p} |\nabla u|^p dx - \int_{\Omega} |\nabla u|^{p-2} \nabla(T_\lambda u - u) dx \end{aligned}$$

Now we use  $|\nabla T_\lambda u| \leq c\lambda \leq \varepsilon_+$ ,  $\kappa_\varepsilon(t) = \frac{1}{p} t^p$  on  $\bar{\varepsilon}$  and  $\nabla T_\lambda u = \nabla u$  on  $\Omega \setminus \mathcal{O}_\lambda(u)$  to see

$$\kappa_\varepsilon(|\nabla T_\lambda u|) - \frac{1}{p} |\nabla u|^p \leq \begin{cases} \frac{1}{p} \varepsilon_-^p & \text{on } \{|\nabla T_\lambda u| \leq \varepsilon_-\}, \\ 0 & \text{on } (\Omega \setminus \mathcal{O}_\lambda(u)) \cap \{|\nabla T_\lambda u| > \varepsilon_-\} \text{ and} \\ \frac{1}{p} |\nabla T_\lambda u|^p & \text{on } \Omega \cap \mathcal{O}_\lambda(u) \cap \{|\nabla T_\lambda u| > \varepsilon_-\}. \end{cases}$$

Combining this with the last estimate, the Hölder Inequality and Theorem 3.4, (v) we deduce

$$\begin{aligned} \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) &\leq |\Omega| \frac{1}{p} \varepsilon_-^p + \int_{\mathcal{O}_\lambda(u)} \frac{1}{p} |\nabla T_\lambda u|^p dx + \int_{\mathcal{O}_\lambda(u)} |\nabla u|^{p-1} |\nabla(T_\lambda u - u)| dx \\ &\lesssim \varepsilon_-^p + \int_{\mathcal{O}_\lambda(u)} |\nabla T_\lambda u|^p dx + \int_{\mathcal{O}_\lambda(u)} |\nabla u|^p dx + \int_{\mathcal{O}_\lambda(u)} |\nabla(T_\lambda u - u)|^p dx \\ &\lesssim \varepsilon_-^p + \int_{\mathcal{O}_\lambda(u)} |\nabla u|^p dx. \quad \square \end{aligned}$$

Note that Lemma 3.9 also implies Corollary 3.7 almost directly with the weak type estimate for the Hardy-Littlewood maximal function (3.2) on page 41.

The estimate (3.4) clearly separates the impact of  $\varepsilon_-$  and  $\varepsilon_+$ . The  $\varepsilon_-$ -term can not be improved, at least not for the energy difference  $\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u)$ : imagine given data  $f$  and  $\Omega$  and just extend  $u_\varepsilon$  and  $u$  to  $B_R(0)$  with  $R > 0$  large enough. Then you will obtain

$$|B_R(0) \setminus \Omega| \varepsilon_-^p \lesssim \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u).$$

Of course, this argument does not hold for the energy difference  $\mathcal{J}(u_\varepsilon) - \mathcal{J}(u)$ .

However, the second summand can be estimated in different ways. Therefore, we will need additional regularity assumptions on  $f$  and  $\Omega$ . These estimates depend on the decay of the level sets of  $|\nabla u|$ . This is perfectly described by the weak  $L^p$  spaces as a special case of Lorentz spaces. We set

$$\|w\|_{L^{q,\infty}(\Omega)} := \sup_{t>0} \|t\chi_{\{|w|>t\}}\|_{L^q(\Omega)} \quad (3.5)$$

and

$$L^{q,\infty}(\Omega) := \{w \in L^1_{\text{loc}}(\Omega) : \|w\|_{L^{q,\infty}(\Omega)} < \infty\}. \quad (3.6)$$

With

$$\|w\|_{L^{q,1}(\Omega)} := q \int_0^\infty \|t\chi_{\{|w|>t\}}\|_{L^q(\Omega)} \frac{dt}{t} \quad (3.7)$$

and

$$L^{q,1}(\Omega) := \{w \in L^1_{\text{loc}}(\Omega) : \|w\|_{L^{q,1}(\Omega)} < \infty\} \quad (3.8)$$

we have  $(L^{q,1}(\Omega))^* \simeq L^{q',\infty}(\Omega)$  for  $q \in (1, \infty)$  and  $\frac{1}{q} + \frac{1}{q'} = 1$  (see [Gra08, Theorem 1.4.17]). With that we can prove the following lemma.

**Lemma 3.10.** *Let  $|\nabla u| \in L^{q,\infty}(\Omega)$  for some  $q > p$ . Then,*

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) \lesssim \varepsilon_-^p + \varepsilon_+^{-(q-p)} \|\nabla u\|_{L^{q,\infty}(\Omega)}.$$

*Proof.* From [CUMP04, Theorem 1.1] we deduce that  $M : L^{q,\infty}(\Omega) \rightarrow L^{q,\infty}(\Omega)$  is bounded. In particular,

$$\begin{aligned} \|\nabla u\|_{L^{q,\infty}(\Omega)} &\gtrsim \|M(|\nabla u|)\|_{L^{q,\infty}(\Omega)} \\ &= \|\lambda\chi_{\{M(|\nabla u|)>\lambda\}}\|_{L^q(\Omega)} \\ &= \lambda|\mathcal{O}_\lambda(u)|^{\frac{1}{q}} \end{aligned}$$

and hence  $|\mathcal{O}_\lambda(u)| \lesssim \|\nabla u\|_{L^{q,\infty}(\Omega)}^q \lambda^{-q}$ . Furthermore we have

$$\begin{aligned} \|\chi_{\{|\chi_{\mathcal{O}_\lambda(u)}| > t\}}\|_{L^{\frac{q}{q-p}}(\Omega)}^{\frac{q}{q-p}} &= |\{|\chi_{\mathcal{O}_\lambda(u)}| > t\}| \\ &= \begin{cases} |\mathcal{O}_\lambda(u)| & \text{for } t \in [0, 1) \text{ and} \\ 0 & \text{for } t \geq 1 \end{cases} \end{aligned}$$

implying that

$$\begin{aligned} \|\chi_{\mathcal{O}_\lambda(u)}\|_{L^{\frac{q}{q-p},1}(\Omega)} &= \frac{q}{q-p} \int_0^\infty \|t\chi_{\{|\chi_{\mathcal{O}_\lambda(u)}| > t\}}\|_{L^{\frac{q}{q-p}}(\Omega)} \frac{dt}{t} \\ &\approx \int_0^1 |\mathcal{O}_\lambda(u)|^{\frac{q-p}{q}} dt \\ &= |\mathcal{O}_\lambda(u)|^{\frac{q-p}{q}}. \end{aligned}$$

Now, the last equivalence combined with  $|\mathcal{O}_\lambda(u)|^{\frac{q-p}{q}} \lesssim \|\nabla u\|_{L^{q,\infty}(\Omega)}^{q-p} \lambda^{-(q-p)}$  yields

$$\begin{aligned} \int_{\mathcal{O}_\lambda(u)} |\nabla u|^p dx &\lesssim \| |\nabla u|^p \|_{L^{\frac{q}{p},\infty}(\Omega)} \|\chi_{\mathcal{O}_\lambda(u)}\|_{L^{\frac{q}{q-p},1}(\Omega)} \\ &\approx \|\nabla u\|_{L^{q,\infty}(\Omega)}^p |\mathcal{O}_\lambda(u)|^{\frac{q-p}{q}} \\ &\lesssim \|\nabla u\|_{L^{q,\infty}(\Omega)}^q \lambda^{-(q-p)}. \end{aligned}$$

Applying Lemma 3.9 with  $\lambda := \frac{\varepsilon_\pm}{c}$  yields the statement.  $\square$

Now that we know that  $L^{q,\infty}$ -regularity of  $|\nabla u|$  can be directly transferred into estimates for  $\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u)$  we give to examples how to derive rates for  $\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u)$  based on the regularity results in [CM10] and [Ebm02].

The first result can be found in [CM10, Theorem 1.3 and Theorem 1.4].

**Theorem 3.11.** *Let  $\Omega \subset \mathbb{R}^d$  be convex or let its boundary  $\partial\Omega \in W^2L^{d-1,1}$  (for example  $\partial\Omega \in C^2$  suffices) and additionally  $f \in L^{d,1}(\Omega)$ . Then,  $\nabla u \in L^\infty(\Omega)$ .*

An other example is not as directly applicable as this one.

**Theorem 3.12.** *Let  $\Omega$  be a polyhedral domain where the inner angle is strictly less than  $2\pi$  and  $f \in L^{p'}(\Omega)$  and  $\frac{1}{p} + \frac{1}{p'} = 1$ . Then  $\nabla u \in L^{\frac{pd}{d-1},\infty}$ .*



*Proof.* It is shown in [Ebm02, (4.3)] that  $|\nabla u|_{\frac{p}{2}} \in \mathcal{N}^{\frac{1}{2},2}(\Omega)$  which denotes the Nikolskii space. We will not recapitulate the definitions of the occurring function spaces but show the embedding  $\mathcal{N}^{\frac{1}{2},2}(\Omega) \hookrightarrow L^{\frac{2d}{d-1},\infty}(\Omega)$ . First of all, we use

$$\mathcal{N}^{\frac{1}{2},2}(\Omega) = B_{2,\infty}^{\frac{1}{2}}(\Omega)$$

as stated in [KOF77, Remark 8.4.5], where  $B_{p,q}^s(\Omega)$  denotes the standard Besov spaces. With [Tri78, Theorems 1 and 2 in 4.3.1] we find the interpolation couple  $\{B_{2,1}^{\frac{1}{4}}(\Omega), B_{2,1}^{\frac{3}{4}}(\Omega)\}$  to see that

$$B_{2,\infty}^{\frac{1}{2}}(\Omega) = (B_{2,1}^{\frac{1}{4}}(\Omega), B_{2,1}^{\frac{3}{4}}(\Omega))_{\frac{1}{2},\infty}$$

holds. The embeddings (see [EEK06] respectively [Pee66])

$$B_{2,1}^{\frac{1}{4}}(\Omega) \hookrightarrow L^{\frac{4d}{2d-1}}(\Omega) \quad \text{and} \quad B_{2,1}^{\frac{3}{4}}(\Omega) \hookrightarrow L^{\frac{4d}{2d-3}}(\Omega)$$

yield

$$(B_{2,1}^{\frac{1}{4}}(\Omega), B_{2,1}^{\frac{3}{4}}(\Omega))_{\frac{1}{2},\infty} \hookrightarrow (L^{\frac{4d}{2d-1}}(\Omega), L^{\frac{4d}{2d-3}}(\Omega))_{\frac{1}{2},\infty}.$$

Finally, by [Tri78, Theorem 2 in 1.18.6] we get

$$(L^{\frac{4d}{2d-1}}(\Omega), L^{\frac{4d}{2d-3}}(\Omega))_{\frac{1}{2},\infty} = L^{\frac{2d}{d-1},\infty}(\Omega).$$

Hence,  $|\nabla u|_{\frac{p}{2}} \in L^{\frac{2d}{d-1},\infty}(\Omega)$ , so

$$\begin{aligned} \infty &> \| |\nabla u|_{\frac{p}{2}} \|_{L^{\frac{2d}{d-1},\infty}(\Omega)} \\ &= \sup_{t>0} \| t^{\frac{p}{2}} \chi_{\{|\nabla u|_{\frac{p}{2}} > t^{\frac{p}{2}}\}} \|_{L^{\frac{2d}{d-1}}(\Omega)} \\ &= \sup_{t>0} \| t \chi_{\{|\nabla u| > t\}} \|_{L^{\frac{pd}{d-1}}(\Omega)}^{\frac{p}{2}} \\ &= \| \nabla u \|_{L^{\frac{pd}{d-1},\infty}(\Omega)}^{\frac{p}{2}}. \quad \square \end{aligned}$$

So we may finish this chapter with the next two corollaries.

**Corollary 3.13.** *Let  $\Omega \subset \mathbb{R}^d$  be convex or let its boundary  $\partial\Omega \in W^2L^{d-1,1}$  (for example  $\partial\Omega \in C^2$  suffices) and additionally  $f \in L^{d,1}(\Omega)$ . Then for  $\varepsilon_+$  large enough one gets*

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) \lesssim \varepsilon_-^p.$$

*Proof.* As shown in Theorem 3.11 in that case  $|\nabla u| \in L^\infty(\Omega)$  so we have that  $M(|\nabla u|) \in L^\infty(\Omega)$ , too, by using (3.1). Therefore,  $\mathcal{O}_\lambda(u) = \emptyset$  for  $\lambda$  large enough so an application of Lemma 3.9 yields the statement.  $\square$

**Corollary 3.14.** *Let  $\Omega$  be a polyhedral domain where the inner angle is strictly less than  $2\pi$  and  $f \in L^{p'}(\Omega)$  and  $\frac{1}{p} + \frac{1}{p'} = 1$ . Then,*

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) \lesssim \varepsilon_-^p + \varepsilon_+^{-\frac{p}{d-1}}.$$

*Proof.* Combining Theorem 3.12 and Lemma 3.10 yields the statement directly.  $\square$

# 4 The Kačanov Iteration

In this chapter, we discuss the convergence of the Kačanov iteration and see a very important but academic example in the second part. As we already have used in the last section we have

$$\mathcal{J}(u_\varepsilon) - \mathcal{J}(u) \leq \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u)$$

so it obviously suffices to show bounds for the latter term.

In this section we study the convergence of the Kačanov-iteration for fixed relaxation parameter  $\varepsilon = (\varepsilon_-, \varepsilon_+)$ .

## 4.1 Exponential Convergence

We recapitulate: for  $v_0 \in W_0^{1,2}(\Omega)$  arbitrary we calculate recursively  $v_{n+1}$  as solution to

$$\int_{\Omega} \Pi_\varepsilon(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi \, dx = \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} \nabla v_{n+1} \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega)$$

as introduced in Section 2.3.

We will show that  $u_n$  converges to the minimizer  $u_\varepsilon$  of the relaxed energy  $\mathcal{J}_\varepsilon$  with exponential decay of the energy error  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon)$ . The proof is based on the following theorem that says that in each iteration we reduce the energy by a certain part of the remaining energy error.

**Theorem 4.1.** *There is a constant  $c > 1$  such that*

$$\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(v_{n+1}) \geq \delta (\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon))$$

holds for  $\delta := \frac{1}{c} \left( \frac{\varepsilon_-}{\varepsilon_+} \right)^{2-p}$ .

*Proof.* Using Theorem 2.42, the equation for  $u_\varepsilon$  and  $ab \leq \frac{1}{2\gamma}a^2 + \frac{\gamma}{2}b^2$  for arbitrary  $a, b \geq 0$  and  $\gamma > 0$

$$\begin{aligned}
\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) &\leq \int_{\Omega} (A_\varepsilon(v_n) - A_\varepsilon(u_\varepsilon)) \nabla(v_n - u_\varepsilon) dx \\
&= \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} \nabla(v_n - v_{n+1}) \nabla(v_n - u_\varepsilon) dx \\
&\leq \underbrace{\frac{1}{\gamma} \frac{1}{2} \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - v_{n+1})|^2 dx}_{=: I} \\
&\quad + \underbrace{\gamma \frac{1}{2} \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - u_\varepsilon)|^2 dx}_{=: II}.
\end{aligned}$$

We use (2.19), the equation for  $\nabla v_{n+1}$  twice, Lemma 2.33 and the fact that  $\Pi_\varepsilon$  is a projection to deduce

$$\begin{aligned}
I &= \frac{1}{2} \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - v_{n+1})|^2 dx \\
&= \frac{1}{2} \int_{\Omega} \Pi_\varepsilon(|\nabla v_n|)^{p-2} (|\nabla v_n|^2 - |\nabla v_{n+1}|^2 - 2\nabla v_n \nabla v_{n+1} + 2|\nabla v_{n+1}|^2) dx \\
&= \int_{\Omega} \frac{1}{2} \Pi_\varepsilon(|\nabla v_n|)^{p-2} |\nabla v_n|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) \Pi_\varepsilon(|\nabla v_n|)^{p-2} dx - \langle f, v_n \rangle \\
&\quad - \int_{\Omega} \frac{1}{2} \Pi_\varepsilon(|\nabla v_n|)^{p-2} |\nabla v_{n+1}|^2 + \left(\frac{1}{p} - \frac{1}{2}\right) \Pi_\varepsilon(|\nabla v_n|)^{p-2} dx + \langle f, v_{n+1} \rangle \\
&= \mathcal{J}_\varepsilon^s(v_n, |\nabla v_n|) - \mathcal{J}_\varepsilon^s(v_{n+1}, |\nabla v_n|) \\
&\leq \mathcal{J}_\varepsilon^s(v_n, |\nabla v_n|) - \mathcal{J}_\varepsilon^s(v_{n+1}, \Pi_\varepsilon(|\nabla v_{n+1}|)) \\
&= \mathcal{J}_\varepsilon^s(v_n, |\nabla v_n|) - \mathcal{J}_\varepsilon^s(v_{n+1}, |\nabla v_{n+1}|) \\
&= \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(v_{n+1}).
\end{aligned}$$

Since  $\frac{\varphi'_\varepsilon(t)}{t} = \Pi_\varepsilon(t)^{p-2}$  is decreasing we have  $\varepsilon_+^{p-2} \leq \frac{\varphi'_\varepsilon(t)}{t} \leq \varepsilon_-^{p-2}$  for all  $t$ . Therefore,  $\frac{\varphi'_\varepsilon(t)}{t} \leq \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p} \frac{\varphi'_\varepsilon(s)}{s}$  for all  $s, t \geq 0$ . Using this, Lemma 2.41 and Theorem 2.42

we obtain that there is  $c' = \frac{c_1 c_2}{2}$  independent of  $\varepsilon$  such that

$$\begin{aligned}
II &= \frac{1}{2} \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - u_\varepsilon)|^2 dx \\
&= \frac{1}{2} \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p} \int_{\Omega} \frac{\varphi'_\varepsilon(|\nabla v_n| |\nabla u_\varepsilon|)}{|\nabla v_n| |\nabla u_\varepsilon|} |\nabla(v_n - u_\varepsilon)|^2 dx \\
&\leq \frac{c_1}{2} \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p} \int_{\Omega} (A_\varepsilon(\nabla v_n) - A_\varepsilon(\nabla u_\varepsilon)) \nabla(v_n - u_\varepsilon) dx \\
&\leq \frac{c_1 c_2}{2} \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p} (\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon)) \\
&\leq c' \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p} (\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon))
\end{aligned}$$

Putting all estimates together we get

$$\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) \leq \frac{1}{\gamma} (\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(v_{n+1})) + \gamma c' \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p} (\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon)),$$

so

$$\gamma(1 - c\gamma \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p}) (\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon)) \leq \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(v_{n+1}).$$

Now,  $\max_{\gamma>0} \gamma(1 - c\gamma \left(\frac{\varepsilon_+}{\varepsilon_-}\right)^{2-p}) = \frac{1}{4c'} \left(\frac{\varepsilon_-}{\varepsilon_+}\right)^{2-p}$  yields the statement.  $\square$

A direct consequence of the last theorem is the following corollary.

**Corollary 4.2.** *For  $v_n$  generated by the algorithm,  $u_\varepsilon$  being the minimizer of  $\mathcal{J}_\varepsilon$  and  $\delta := \frac{1}{c} \left(\frac{\varepsilon_-}{\varepsilon_+}\right)^{2-p}$  we get*

$$\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) \leq (1 - \delta)^n (\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)).$$

*Proof.* We prove this by induction. The case when  $n = 0$  is clear. So let the statement be true for  $n - 1$ . Then, with Theorem 4.1 we get

$$\begin{aligned}
\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) &= \mathcal{J}_\varepsilon(v_{n-1}) - \mathcal{J}_\varepsilon(u_\varepsilon) + \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(v_{n-1}) \\
&\leq (1 - \delta) (\mathcal{J}_\varepsilon(v_{n-1}) - \mathcal{J}_\varepsilon(u_\varepsilon)) \\
&\leq (1 - \delta) (1 - \delta)^{n-1} (\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)) \\
&= (1 - \delta)^n (\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)). \quad \square
\end{aligned}$$

As for the limit with respect to the relaxation parameter  $\varepsilon$  we want to state the convergence with to the underlying space  $W_0^{1,\varphi_\varepsilon}(\Omega)$ .

**Corollary 4.3.**  $v_n \xrightarrow{n \rightarrow \infty} u_\varepsilon$  in  $W_0^{1,\varphi_\varepsilon}(\Omega)$ .

*Proof.* As in the proof of Corollary 3.8 we deduce for any  $t \geq 0$ ,  $\gamma > 0$  and  $P \in \mathbb{R}^d$  the estimate

$$\varphi_\varepsilon(t) \leq c_\gamma \varphi_{\varepsilon,|P|}(t) + \gamma \varphi_\varepsilon(|P|).$$

Again with Lemma 2.41 and Lemma 3.1 we get

$$\begin{aligned} \int_{\Omega} \varphi_\varepsilon(|\nabla(v_n - u_\varepsilon)|) dx &\leq c_\gamma \int_{\Omega} \varphi_{\varepsilon,|\nabla u_\varepsilon|}(|\nabla(v_n - u_\varepsilon)|) dx + \gamma \int_{\Omega} \varphi_\varepsilon(|\nabla u_\varepsilon|) dx \\ &\lesssim \int_{\Omega} (A_\varepsilon(\nabla v_n) - A_\varepsilon(\nabla u_\varepsilon)) \nabla(v_n - u_\varepsilon) dx + \gamma \int_{\Omega} \varphi_\varepsilon(|\nabla u_\varepsilon|) dx \\ &\lesssim \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) + \gamma \int_{\Omega} \varphi_\varepsilon(|\nabla u_\varepsilon|) dx. \end{aligned}$$

Choosing  $\gamma$  sufficiently small and  $n$  sufficiently large yields convergence with respect to the modulus, that is  $\int_{\Omega} \varphi_\varepsilon(|\nabla(v_n - u_\varepsilon)|) dx \xrightarrow{n \rightarrow \infty} 0$ . With Theorem 2.27, 2. we get the statement.  $\square$

## 4.2 An Example

As in the previous section we fix  $\varepsilon$  such that  $1 \in \varepsilon$ . Furthermore, let  $\Omega := B_1(0)$  and  $f \in (W_0^{1,p}(\Omega))^*$  such that  $u(x) = 1 - |x|$ . That is achieved by the distribution  $f(x) = -\operatorname{div}(\frac{x}{|x|}) \notin L^1(\Omega)$ . Since  $|\nabla u| \equiv 1$ , the factor  $|\nabla u|^{p-2}$  in the non-linear weight does not appear for the minimizer:

$$\int_{\Omega} \nabla u \nabla \xi dx = \int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,\varphi_\varepsilon}(\Omega). \quad (4.1)$$

Furthermore, for all  $w \in W_0^{1,\varphi_\varepsilon}(\Omega)$  we get

$$\begin{aligned} \mathcal{J}_\varepsilon(u) &= \int_{\Omega} \kappa_\varepsilon(|\nabla u|) dx - \langle f, u \rangle \\ &= \int_{\Omega} \frac{1}{p} |\nabla u|^p dx - \langle f, u \rangle \\ &= \mathcal{J}(u) \\ &\leq \mathcal{J}(w) \\ &\leq \mathcal{J}_\varepsilon(w). \end{aligned}$$

So in this case  $u$  also minimizes every  $\mathcal{J}_\varepsilon$  with  $1 \in \varepsilon$ , so  $u_\varepsilon = u$ .

We want to elaborate how the Kačanov iteration performs in this academic example for  $v_0 = 0$ . We show directly, that we have  $v_n = a_n u$  where  $a_n$  is recursively defined via

$$a_0 := 0 \quad \text{and} \quad a_{n+1} := \Pi_\varepsilon(a_n)^{2-p}. \quad (4.2)$$

This follows with (4.1) and since then

$$\begin{aligned} \int_{\Omega} \Pi_\varepsilon(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi dx &= \int_{\Omega} \Pi_\varepsilon(a_n |\nabla u|)^{p-2} a_{n+1} \nabla u \nabla \xi dx \\ &= \int_{\Omega} \Pi_\varepsilon(a_n)^{p-2} \Pi_\varepsilon(a_n)^{2-p} \nabla u \nabla \xi dx \\ &= \int_{\Omega} \nabla u \nabla \xi dx \\ &= \langle f, \xi \rangle. \end{aligned}$$

Note that this representation also holds true for  $p \geq 2$ . We will discuss this case later.

So let  $p < 2$ . Then for  $1 \in \varepsilon$  and any  $t \in \bar{\varepsilon}$  we get  $\varepsilon_- \leq \varepsilon_-^{2-p} \leq t^{2-p} \leq \varepsilon_+^{2-p} \leq \varepsilon_+$ . So if  $a_n \in \bar{\varepsilon}$  we directly get  $a_{n+1} \in \bar{\varepsilon}$ . With  $a_0 = 0$  we get  $a_1 = \Pi_\varepsilon(0)^{2-p} = \varepsilon_-^{2-p} \in \bar{\varepsilon}$  since  $\varepsilon_-^{2-p} \geq \varepsilon_-$  and  $\varepsilon_-^{2-p} \leq 1 \leq \varepsilon_+$ . In particular we get

$$a_0 = 0 \quad \text{and} \quad a_n = \varepsilon_-^{(2-p)^n} \in \bar{\varepsilon} \quad \text{for } n \geq 1$$

directly by (4.2).

With that we can calculate for  $n \geq 1$

$$\begin{aligned}
\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) &= \mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u) \\
&= \int_{\Omega} \kappa_\varepsilon(|\nabla v_n|) - \frac{1}{p} |\nabla u|^p dx - \langle f, v_n - u \rangle \\
&= \int_{\Omega} \kappa_\varepsilon(a_n) - \frac{1}{p} dx - \int_{\Omega} \nabla u \nabla(v_n - u) dx \\
&= \left(\frac{\alpha_n^p}{p} - \frac{1}{p}\right) \int_{\Omega} 1 dx - (\alpha_n - 1) \int_{\Omega} \underbrace{\nabla u \nabla u}_{=1} dx \\
&= \frac{1}{p} |B_1(0)| (\alpha_n^p - 1 - p(\alpha_n - 1)) \\
&= \frac{1}{p} |B_1(0)| (\varepsilon_-^{p(2-p)^n} - 1 - p(\varepsilon_-^{(2-p)^n} - 1)).
\end{aligned}$$

and similar to above

$$\begin{aligned}
\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) &= \mathcal{J}_\varepsilon(v_0) - \mathcal{J}(u) \\
&= \int_{\Omega} \kappa_\varepsilon(0) - \frac{1}{p} |\nabla u|^p dx + \langle f, u \rangle \\
&= |B_1(0)| \left( \left(\frac{1}{p} - \frac{1}{2}\right) \varepsilon_-^p + 1 - \frac{1}{p} \right)
\end{aligned}$$

for completion.

The interesting part is of course the (asymptotic) behaviour for  $n \geq 1$ . It is stated in the following lemma.

**Lemma 4.4.** *The estimate*

$$\frac{1}{p} (\varepsilon_-^{p(2-p)^n} - 1 - p(\varepsilon_-^{(2-p)^n} - 1)) \leq \frac{p-1}{2} \ln(\varepsilon_-)^2 (2-p)^{2n}$$

holds true and is asymptotic in the sense that

$$\frac{\frac{1}{p} (\varepsilon_-^{p(2-p)^n} - 1 - p(\varepsilon_-^{(2-p)^n} - 1))}{\frac{p-1}{2} \ln(\varepsilon_-)^2 (2-p)^{2n}} \xrightarrow{n \rightarrow \infty} 1.$$

*Proof.* Let  $t \in (0, 1]$ . We introduce the functions  $g, h : (0, 1] \rightarrow \mathbb{R}$  via

$$g(t) := \frac{p-1}{2} \ln(t)^2 - \frac{1}{p} (t^p - 1) + (t - 1) \quad \text{and} \quad h(t) := 1 - \ln(t) - t^p.$$



Since  $h'(t) = -\frac{1}{t} - pt^{p-1} \leq 0$  and  $h(1) = 0$  we get that  $h(t) \geq 0$ . The first and second derivative of  $g$  are calculated as

$$\begin{aligned} g'(t) &= (p-1)\ln(t)\frac{1}{t} - t^{p-1} + 1 \text{ and} \\ g''(t) &= (p-1)\left(\frac{1}{t^2} - \ln(t)\frac{1}{t^2} - t^{p-2}\right) \\ &= (p-1)t^{-2}(1 - \ln(t) - t^p) \\ &= (p-1)t^{-2}h(t) \\ &\geq 0. \end{aligned}$$

So  $g$  is convex. This implies  $\frac{g(1)-g(t)}{1-t} \leq g'(1)$ . Now by  $g(1) = g'(1) = 0$  and  $1-t \geq 0$  we get  $g(t) \geq 0$ . In particular,

$$\begin{aligned} 0 &\leq g(\varepsilon_-^{(2-p)^n}) \\ &= \frac{p-1}{2} \ln(\varepsilon_-^{(2-p)^n})^2 - \frac{1}{p}(\varepsilon_-^{p(2-p)^n} - 1) + (\varepsilon_-^{(2-p)^n} - 1) \\ &= \frac{p-1}{2}(2-p)^{2n} \ln(\varepsilon_-)^2 - \frac{1}{p}((\varepsilon_-^{p(2-p)^n} - 1) - p(\varepsilon_-^{(2-p)^n} - 1)) \end{aligned}$$

which shows the estimate.

The asymptotic behaviour follows from applying l'Hôpital's rule twice and choosing  $t_n = \varepsilon_-^{(2-p)^n} \xrightarrow{n \rightarrow \infty} 1$  in

$$\begin{aligned} \lim_{t \rightarrow 1} \frac{\frac{1}{p}(t^p - 1 - p(t-1))}{\frac{p-1}{2} \ln(t)^2} &= \lim_{t \rightarrow 1} \frac{t^{p-1} - 1}{(p-1)\ln(t)\frac{1}{t}} \\ &= \lim_{t \rightarrow 1} \frac{t^p - t}{(p-1)\ln(t)} \\ &= \lim_{t \rightarrow 1} \frac{pt^{p-1} - 1}{(p-1)\frac{1}{t}} \\ &= 1. \end{aligned} \quad \square$$

So indeed, the energy differences  $\mathcal{J}(v_n) - \mathcal{J}(u) = \mathcal{J}_\varepsilon(u_n) - \mathcal{J}(u)$  asymptotically behave like  $\frac{1}{2}|B_1(0)|(p-1)\ln(\varepsilon_-)^2(2-p)^{2n}$  for large  $n$ . In particular, the estimate

$$\begin{aligned} \mathcal{J}(v_n) - \mathcal{J}(u) &= \mathcal{J}_\varepsilon(u_n) - \mathcal{J}(u) \\ &= \mathcal{J}_\varepsilon(u_n) - \mathcal{J}(u) \\ &\leq \frac{1}{2}|B_1(0)|(p-1)\ln(\varepsilon_-)^2(2-p)^{2n} \end{aligned}$$

is sharp.

This asymptotic shows that in this particular case

$$\mathcal{J}_\varepsilon(u_n) - \mathcal{J}_\varepsilon(u_\varepsilon) \leq c_\varepsilon (1 - \delta)^n$$

with  $\delta = 1 - (2 - p)^2 < 1$  independent of  $\varepsilon$ . Therefore, it remains open if such an estimate holds in the general case.

Note that even the energy difference  $\mathcal{J}(v_n) - \mathcal{J}(u)$  truly depends on  $\varepsilon$  – where  $\mathcal{J}(w) - \mathcal{J}(u)$  in general does not. Furthermore, it is possible to deduce the following theorem with this example. Note especially that it even holds for  $\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)$  arbitrarily small (but not equal to zero), so the  $\varepsilon$ -dependence of  $\delta$  is not due to a bad choice of the initial value.

**Theorem 4.5.** *There is no  $\gamma > 0$  such that there is a factor  $\delta \in (0, 1)$  independent of  $\varepsilon$  satisfying*

$$\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(v_{n+1}) \geq \delta(\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon))$$

for all  $v_0 \in X_\gamma := \{w \in W_0^{1,p(\cdot)}(\Omega) : \mathcal{J}_\varepsilon(w) - \mathcal{J}_\varepsilon(u_\varepsilon) < \gamma\}$  and  $n \geq 0$ .

*Proof.* We assume the existence of such a  $\gamma > 0$ , define  $w_0 := 0$  and  $(w_n)$  as the sequence generated by our algorithm. We already know that  $\mathcal{J}_\varepsilon(w_n) \xrightarrow{n \rightarrow \infty} \mathcal{J}_\varepsilon(w)$ , so there is  $N_\gamma \in \mathbb{N}_{\geq 1}$  such that  $w_{N_\gamma} \in X_\gamma$ . Now we choose  $v_0 := w_{N_\gamma}$ . It is clear that this implies  $v_n = w_{n+N_\gamma} = \varepsilon_-^{(2-p)^{n+N_\gamma}}$ . With that we deduce

$$\begin{aligned} \delta &\leq \frac{\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(v_1)}{\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)} \\ &= \frac{\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon) - (\mathcal{J}_\varepsilon(v_1) - \mathcal{J}_\varepsilon(u_\varepsilon))}{\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)} \\ &= 1 - \frac{\mathcal{J}_\varepsilon(v_1) - \mathcal{J}_\varepsilon(u_\varepsilon)}{\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)}. \end{aligned}$$

Since  $\delta$  is independent of  $\varepsilon$  we get

$$\begin{aligned} \delta &\leq \lim_{\varepsilon_- \searrow 0} 1 - \frac{\mathcal{J}_\varepsilon(v_1) - \mathcal{J}_\varepsilon(u_\varepsilon)}{\mathcal{J}_\varepsilon(v_0) - \mathcal{J}_\varepsilon(u_\varepsilon)} \\ &= 1 - \lim_{\varepsilon_- \searrow 0} \frac{\mathcal{J}_\varepsilon(w_{1+N_\gamma}) - \mathcal{J}_\varepsilon(u_\varepsilon)}{\mathcal{J}_\varepsilon(w_{N_\gamma}) - \mathcal{J}_\varepsilon(u_\varepsilon)} \\ &= 1 - \lim_{\varepsilon_- \searrow 0} \frac{\varepsilon_-^{p(2-p)^{1+N_\gamma}} - 1 - p(\varepsilon_-^{p(2-p)^{1+N_\gamma}} - 1)}{\varepsilon_-^{p(2-p)^{N_\gamma}} - 1 - p(\varepsilon_-^{p(2-p)^{N_\gamma}} - 1)} \\ &= 0. \end{aligned}$$

This obviously contradicts  $\delta \in (0, 1)$ . □

Also note that this statement holds true when replacing the the relaxed energies by the  $p$ -Poisson energy

$$\mathcal{J}(v_n) - \mathcal{J}(v_{n+1}) \geq \delta(\mathcal{J}(v_n) - \mathcal{J}(u_\varepsilon)),$$

since the gradients of the constructed sequence satisfy  $|\nabla v_n| \in \varepsilon$ .

Although our considerations are all under the assumption  $1 < p \leq 2$  it is interesting to check how our algorithm performs in the case  $p > 2$  for our example. Note that (4.2) holds true for any  $p > 1$ .

First, we consider the case  $p \geq 3$ . We chose  $\varepsilon_+ := \frac{1}{\varepsilon_-}$  for some arbitrary  $\varepsilon_- < 1$ . We show by induction, that

$$a_0 = 0 \quad \text{and} \quad a_n = \varepsilon_-^{(-1)^n(p-2)} \quad \text{for } n \geq 1.$$

Note that  $p \geq 3$  implies  $\varepsilon_-^{2-p} \geq \varepsilon_-^{-1} = \varepsilon_+$  and  $\varepsilon_-^{p-2} \leq \varepsilon_-$ . Then, the base clause is given by

$$\begin{aligned} a_0 &= 0, \\ a_1 &= \Pi_\varepsilon(a_0)^{2-p} = \Pi_\varepsilon(0)^{2-p} = \varepsilon_-^{2-p} \quad \text{and} \\ a_2 &= \Pi_\varepsilon(a_1)^{2-p} = \Pi_\varepsilon(\varepsilon_-^{2-p})^{2-p} = \varepsilon_+^{2-p} = \varepsilon_-^{p-2}. \end{aligned}$$

Now, let  $n$  be an even number. By induction hypothesis we get  $a_n = \varepsilon_-^{(-1)^n(p-2)} = \varepsilon_-^{p-2}$ . Therefore,

$$a_{n+1} = \Pi_\varepsilon(a_n)^{2-p} = \Pi_\varepsilon(\varepsilon_-^{p-2})^{2-p} = \varepsilon_-^{2-p} = \varepsilon_-^{(-1)^{n+1}(p-2)}.$$

If  $n$  is an odd number we get  $a_n = \varepsilon_-^{(-1)^n(p-2)} = \varepsilon_-^{2-p}$  by induction hypothesis. Then,

$$a_{n+1} = \Pi_\varepsilon(a_n)^{2-p} = \Pi_\varepsilon(\varepsilon_-^{2-p})^{2-p} = \varepsilon_+^{2-p} = \varepsilon_-^{p-2} = \varepsilon_-^{(-1)^{n+1}(p-2)}.$$

To study the convergence of the sequence  $(v_n)$  the energy is not a suitable tool anymore since it deduced for  $p \in (1, 2)$ . So consider  $\|\cdot\|$  to be a norm on  $W_0^{1,p}(\Omega)$ ,  $W_0^{1,2}(\Omega)$  or on  $W_0^{1,\varphi_\varepsilon}(\Omega)$ . Then, for  $n \geq 1$

$$\begin{aligned} \|v_n - u_\varepsilon\| &= \|v_n - u\| \\ &= |1 - \varepsilon_-^{(-1)^n(2-p)}| \|u\| \\ &= \begin{cases} (\varepsilon_-^{2-p} - 1) \|u\| & \text{for } n \text{ even and} \\ (1 - \varepsilon_-^{p-2}) \|u\| & \text{for } n \text{ odd} \end{cases} \\ &= (-1)^n (\varepsilon_-^{(-1)^n(2-p)} - 1) \|u\|. \end{aligned}$$

Hence, we can formulate the following corollary.

**Corollary 4.6.** For  $p \geq 3$  and  $\varepsilon_- \in (0, 1)$ , the by

$$\int_{\Omega} \Pi_{\varepsilon}(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega)$$

recursively defined sequence does not converge in general.

*Proof.* This follows directly by  $\varepsilon_-^{2-p} - 1 \neq 1 - \varepsilon_-^{p-2}$  for  $\varepsilon_- \neq 1$ .  $\square$

For  $p \in (2, 3)$  we show that

$$a_0 = 0 \quad \text{and} \quad a_n = \varepsilon_-^{(2-p)^n} \quad \text{for } n \geq 1,$$

too. It is easy to see that  $a_1 = \Pi_{\varepsilon}(a_0)^{2-p} = \Pi_{\varepsilon}(0)^{2-p} = \varepsilon_-^{(2-p)^1}$ . For the induction step we use  $\varepsilon_- \leq \varepsilon_-^{(2-p)^n} \leq \varepsilon_-^{-1} = \varepsilon_+$  which is due to the fact that  $-1 \leq (2-p)^n \leq 1$  for all  $n \geq 1$  and calculate

$$a_{n+1} = \Pi_{\varepsilon}(a_n)^{2-p} = \Pi_{\varepsilon}(\varepsilon_-^{(2-p)^n})^{2-p} = \varepsilon_-^{(2-p)^{n+1}}.$$

As already mentioned, the relaxed energy is not suitable for  $p > 2$ . But since  $|v_n| \in \varepsilon$  for  $n \geq 1$  the relaxed energies coincide with the unrelaxed ones which are suitable indeed. So

$$\begin{aligned} \mathcal{J}(v_n) - \mathcal{J}(u) &= \mathcal{J}_{\varepsilon}(v_n) - \mathcal{J}(u) \\ &= \mathcal{J}_{\varepsilon}(v_n) - \mathcal{J}_{\varepsilon}(u_{\varepsilon}) \\ &= \int_{\Omega} \kappa_{\varepsilon}(\varepsilon_-^{(2-p)^n} |\nabla u|) - \frac{1}{p} |\nabla u|^p \, dx - \langle f, v_n - u \rangle \\ &= \frac{1}{p} (\varepsilon_-^{(2-p)^n} - 1) \int_{\Omega} 1 \, dx - (\varepsilon_-^{(2-p)^n} - 1) \int_{\Omega} \nabla u \nabla u \, dx \\ &= \frac{1}{p} |B_1(0)| (\varepsilon_-^{p(2-p)^n} - 1 - p((\varepsilon_-^{(2-p)^n} - 1))). \end{aligned}$$

for  $n \geq 1$ . Note that still  $\varepsilon_-^{(2-p)^n} \xrightarrow{n \rightarrow \infty} 1$ , hence as in Lemma 4.4 we get

$$\lim_{n \rightarrow \infty} \frac{\mathcal{J}(v_n) - \mathcal{J}(u)}{\frac{p-1}{2} \ln(\varepsilon_-)^2 (2-p)^{2n}} = 1$$

for  $p \in (2, 3)$ , too. In particular, the algorithm converges for this example also for these choices of  $p$ .

# 5 Overall Convergence Analysis

## 5.1 An Algebraic Rate

As we learned in the last section the Kačanov Iteration converges, but the rate depends badly on the choice of the relaxation interval  $\varepsilon = (\varepsilon_-, \varepsilon_+)$ . Furthermore, we have algebraic decay of the error induced by the relaxation – at least under certain regularity assumptions. We will assume that  $|\nabla u| \in L^{q,\infty}(\Omega)$  for some  $q > p$  and combine the results for both errors to deduce an algebraic rate over all.

To do so we introduce a quantity  $\mathcal{G}_n$  that carries both errors and satisfies

$$(\mathcal{G}_n - \mathcal{G}_\infty) \leq (\mathcal{G}_1 - \mathcal{G}_\infty) \prod_{i=1}^{n-1} (1 - \delta_i).$$

As we are going to use Theorem 4.1 we will have the dependency  $\delta_n \approx (\frac{\varepsilon_{-,n}}{\varepsilon_{+,n}})^{2-p}$ . Hence, we will need to ensure that this ration tends to zero slow enough.

To use the convergence in the relaxation parameter we assume – as already mentioned – that  $|\nabla u| \in L^{q,\infty}(\Omega)$  for some  $q > p$ . Lemma 3.10 ensures the existence of a constant  $c_R > 0$  such that

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}(u) \leq c_R(\varepsilon_-^p + \varepsilon_+^{-(q-p)}). \quad (5.1)$$

Hence, it is reasonable to define

$$\mathcal{G}_n := \mathcal{J}_{\varepsilon_n}(v_n) + M(\varepsilon_{-,n}^p + \varepsilon_{+,n}^{-(q-p)}) \quad \text{and} \quad \mathcal{G}_\infty := \mathcal{J}(u), \quad (5.2)$$

where  $M > 0$  will be determined later. This definition allows to work with the relaxed energy to use the decay of the Kačanov iteration and implies

$$\begin{aligned} \mathcal{J}(v_n) - \mathcal{J}(u) &\leq \mathcal{J}_{\varepsilon_n}(v_n) - \mathcal{J}(u) \\ &\lesssim \mathcal{G}_n - \mathcal{G}_\infty, \end{aligned}$$

so any decay rate for  $\mathcal{G}_n - \mathcal{G}_\infty$  is a decay rate of the global energy error of the iteration. The first result of this section shows how the main statement of Theorem 4.1 can be transferred to  $\mathcal{G}_n$  for a special choice of  $(\varepsilon_{-,n})$  and  $(\varepsilon_{+,n})$ . To be more precise, we choose  $\alpha, \beta > 0$  with  $\alpha + \beta = \frac{1}{2-p}$  and  $\varepsilon_n = (n^{-\alpha}, n^\beta)$ . We will see later, why this coupling is necessary.

With our choice of the relaxation parameters the algorithm reads as follows.

**Algorithm:** The non-adaptive relaxed  $p$ -Kačanov algorithm

**Data:**  $\Omega \subset \mathbb{R}^d$  and  $f \in (W_0^{1,p}(\Omega))^*$  such that  $|\nabla u| \in L^{q,\infty}(\Omega)$  for  $q > p$ ;  
 $\alpha, \beta > 0$  such that  $\alpha + \beta = \frac{1}{2-p}$ ;

**Result:** Approximate solution of the  $p$ -Poisson problem (2.5).

$n := 1$ ;

**while** desired accuracy is not achieved **yet do**

    Calculate  $v_{n+1}$  as solution of

$$\int_{\Omega} \Pi_{(n^{-\alpha}, n^\beta)}(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega);$$

    Update  $n \rightsquigarrow n + 1$ ;

**end**

Note that it is not necessary to define  $v_1$ , since for any  $v_1 \in W_0^{1,2}(\Omega)$  we have  $\Pi_{(1^{-\alpha}, 1^\beta)}(|\nabla v_1|)^{p-2} = 1$ , so  $v_2$  is the solution to the related 2-Poisson problem. For an easier readability of the proof we first prove the following lemma. It shows that for an algebraic sequence, its difference with its index shifted version decays at most with an additional factor of  $n^{-1}$ .

**Lemma 5.1.** *Let  $\gamma > 0$ . Then for all  $n \geq 1$  we have*

$$n^{-\gamma} - (n+1)^{-\gamma} \geq n^{-\gamma-1} \min\left\{\frac{\gamma}{2}, 1 - 2^{-\gamma}\right\}.$$

*Proof.* We define  $h : [0, \frac{1}{2}] \rightarrow \mathbb{R}$  via  $h(t) := 1 - (1-t)^\gamma$ . Note that  $h'(t) = \gamma(1-t)^{\gamma-1}$  and  $h''(t) = \gamma(1-\gamma)(1-t)^{\gamma-2}$ . For  $\gamma \geq 1$  this implies that  $h$  is concave, so

$$\begin{aligned} h(t) &\geq t \left( \frac{h(\frac{1}{2}) - h(0)}{\frac{1}{2}} \right) \\ &= 2th\left(\frac{1}{2}\right) \\ &= 2(1 - 2^{-\gamma})t \end{aligned}$$

On the other hand, if  $\gamma \in (0, 1)$ , the function  $h$  is convex. Therefore,

$$\begin{aligned} h(t) &\geq h(0) + th'(0) \\ &= \gamma t. \end{aligned}$$

Since  $2(1 - 2^{-\gamma})$  is concave in  $\gamma$  and since  $2(1 - 2^{-0}) = 0$  as well as  $2(1 - 2^{-1}) = 1$  we see

$$\min\{\gamma, 2(1 - 2^{-\gamma})\} = \begin{cases} \gamma & \text{for } \gamma \in (0, 1) \text{ and} \\ 2(1 - 2^{-\gamma}) & \text{for } \gamma \geq 1. \end{cases}$$

Overall, this implies  $h(t) \geq \min\{\gamma, 2(1 - 2^{-\gamma})\}t$ . Therefore, we get

$$\begin{aligned} n^{-\gamma} - (n+1)^{-\gamma} &= n^{-\gamma}(1 - (\frac{n+1}{n})^{-\gamma}) \\ &= n^{-\gamma}(1 - (1 - \frac{1}{n+1})^\gamma) \\ &= n^{-\gamma}h(\frac{1}{n+1}) \\ &\geq n^{-\gamma} \min\{\gamma, 2(1 - 2^{-\gamma})\} \frac{1}{n+1} \\ &\geq n^{-\gamma-1} \min\{\frac{\gamma}{2}, 1 - 2^{-\gamma}\}. \end{aligned} \quad \square$$

With that, we can prove an estimate similar to Theorem 4.1. Therefore, we choose  $M$  in the definition of  $\mathcal{G}_n$  – see (5.2) – such that

$$\frac{M + c_R}{Mc_K} \geq \max\{\frac{\alpha p}{2}, \frac{\beta p}{2}, 1 - 2^{-\alpha p}, 1 - 2^{-\beta p}\}.$$

Since  $\frac{M+c_R}{Mc_K}$  is not bounded as  $M > 0$  gets small, this is always possible.

**Theorem 5.2.** *Let  $|\nabla u| \in L^{q,\infty}(\Omega)$  for some  $q > p$ . Then, the sequence  $(v_n)$  generated by the algorithm on page 62 satisfies*

$$\mathcal{G}_n - \mathcal{G}_{n+1} \geq \frac{1}{nc_K}(\mathcal{G}_n - \mathcal{G}_\infty) \quad (5.3)$$

where  $c_K$  is the constant of Theorem 4.1.

*Proof.* With Lemma 5.1 for  $\gamma \in \{\alpha p, \beta(q-p)\}$  and our choice of  $M$  we get

$$Mn^{-\gamma} - M(n+1)^{-\gamma} \geq \frac{M+c_R}{c_K}n^{-\gamma-1}$$

or for  $\delta_n := \frac{1}{nc_K}$

$$-c_R\delta_n n^{-\gamma} \geq M\delta_n n^{-\gamma} - M(n^{-\gamma} - (n+1)^{-\gamma}),$$

respectively. Summing up this inequality for the two choices of  $\gamma$  we get

$$\begin{aligned} -c_R \delta_n (n^{-\alpha p} + n^{-\beta(q-p)}) &\geq M \delta_n (n^{-\alpha p} + n^{-\beta(q-p)}) \\ &\quad - M (n^{-\alpha p} + n^{-\beta(q-p)}) \\ &\quad + M ((n+1)^{-\alpha p} + (n+1)^{-\beta(q-p)}). \end{aligned} \quad (5.4)$$

We write  $\varepsilon_n := (n^{-\alpha}, n^\beta)$ . For the next estimate, we use Theorem 4.1 and our choice of  $\alpha, \beta > 0$  with  $\alpha + \beta = \frac{1}{2-p}$ . That implies  $\delta_n = \frac{1}{nc_K}$ , so together with (5.1) we directly deduce the estimate

$$\begin{aligned} \mathcal{J}_{\varepsilon_n}(v_n) - \mathcal{J}_{\varepsilon_{n+1}}(v_{n+1}) &\geq \mathcal{J}_{\varepsilon_n}(v_n) - \mathcal{J}_{\varepsilon_n}(v_{n+1}) \\ &\geq \delta_n (\mathcal{J}_{\varepsilon_n}(v_n) - \mathcal{J}_{\varepsilon_n}(u_{\varepsilon_n})) \\ &= \delta_n (\mathcal{J}_{\varepsilon_n}(v_n) - \mathcal{J}(u)) - \delta_n (\mathcal{J}_{\varepsilon_n}(u_{\varepsilon_n}) - \mathcal{J}(u)) \\ &\geq \delta_n (\mathcal{J}_{\varepsilon_n}(v_n) - \mathcal{J}(u)) - \delta_n c_R (n^{-\alpha p} + n^{-\beta(q-p)}). \end{aligned}$$

Now the statement follows by applying (5.4) to the last estimate.  $\square$

With that we can perform similar steps as in the proof of Corollary 4.2 to obtain the following result.

**Corollary 5.3.** *Let  $|\nabla u| \in L^{q,\infty}(\Omega)$  for some  $q > p$  (for example  $f \in L^p(\Omega)$  and  $\Omega$  rectangular with inner angle less than  $2\pi$ , see Corollary 3.14). Then, the non-adaptive algorithm on page 62 satisfies*

$$\mathcal{G}_n - \mathcal{G}_\infty \leq (\mathcal{G}_1 - \mathcal{G}_\infty) n^{-c_K^{-1}}$$

where  $c_K$  is the constant of Theorem 4.1. In particular,

$$\mathcal{J}(v_n) - \mathcal{J}(u) \lesssim n^{-c_K^{-1}}.$$

*Proof.* From Theorem 5.2 we directly deduce  $\mathcal{G}_{n+1} - \mathcal{G}_\infty \leq (1 - \delta_n)(\mathcal{G}_n - \mathcal{G}_\infty)$  for  $n \geq 1$ . This inductively implies

$$\mathcal{G}_n - \mathcal{G}_\infty \leq (\mathcal{G}_1 - \mathcal{G}_\infty) \prod_{i=1}^{n-1} (1 - \delta_i).$$

With the integral test for convergence of series we get

$$\begin{aligned} \sum_{i=1}^{n-1} \delta_i &\geq \int_1^n \frac{1}{tc_K} dt \\ &= \ln(n^{c_K^{-1}}). \end{aligned}$$



Using this in

$$\begin{aligned} \prod_{i=0}^{n-1} (1 - \delta_i) &= \exp\left(\sum_{i=0}^{n-1} \ln(1 - \delta_i)\right) \\ &\leq \exp\left(-\sum_{i=0}^{n-1} \delta_i\right) \\ &\leq \exp(-\ln(n^{c_K^{-1}})) \\ &= n^{-c_K^{-1}}. \end{aligned}$$

we get the statement.  $\square$

However, with this technique one can not proof better results. This is due to the fact that the estimate

$$n^{-\gamma} - (n+1)^{-\gamma} \gtrsim n^{-1}n^{-\gamma}$$

is needed to cover

$$n^{-\alpha p} - (n+1)^{-\alpha p} \gtrsim \delta_n n^{-\alpha p}.$$

Hence, with this technique  $\delta_n$  may not tend to zero slower than  $n^{-1}$ . This is also the reason, why we need the relation  $\alpha + \beta = \frac{1}{2-p}$ .

Also note that the estimate

$$\mathcal{G}_n - \mathcal{G}_{n+1} \geq \frac{1}{nc_K}(\mathcal{G}_n - \mathcal{G}_\infty)$$

does not improve when  $q$  becomes larger. Anyway, if  $q$  is known, we can optimize the stabilization of

$$\mathcal{G}_n = \mathcal{J}_{\varepsilon_n}(v_n) + M(n^{-\alpha p} + n^{-\beta(q-p)})$$

in the sense that  $\alpha p = \beta(q-p)$ . Since we need to ensure  $\alpha + \beta = \frac{1}{2-p}$ , too, this reads as

$$\begin{pmatrix} 1 & 1 \\ p & p-q \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \frac{1}{2-p} \\ 0 \end{pmatrix}$$

which is solved by

$$\begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \frac{1}{-q} \begin{pmatrix} p-q & -1 \\ -p & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2-p} \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{q-p}{q(2-p)} \\ \frac{p}{q(2-p)} \end{pmatrix}.$$

Then,  $\mathcal{G}_n = \mathcal{J}_{\varepsilon_n}(v_n) + 2Mn^{-\frac{p(q-p)}{q(2-p)}}$ . Hence, the impact of the stabilization decreases as  $q$  increases.

As we have seen before, we also get convergence in the norm on  $W_0^{1,p}(\Omega)$  and state this here for completion.

**Corollary 5.4.** *Under the assumptions of Corollary 5.3, the sequence  $(v_n)$  of the non-adaptive algorithm on page 62 satisfies  $v_n \xrightarrow{n \rightarrow \infty} u$  in  $W_0^{1,p}(\Omega)$ .*

*Proof.* With the same estimate as in Corollary 3.8 we get for any  $\delta > 0$  the estimate

$$\int_{\Omega} |\nabla(v_n - u)|^p dx \leq c_{\delta}(\mathcal{J}(v_n) - \mathcal{J}(u)) + \delta \int_{\Omega} |\nabla u|^p dx.$$

Choosing  $\delta$  sufficiently small and  $n$  sufficiently large the term on the right hand side becomes arbitrary small.  $\square$

## 5.2 Outlook On Adaptive Strategies for the Relaxation Parameter

In the last section we saw an example for a strategy to couple a certain behaviour of the relaxation parameter to the Kačanov iteration to deduce an algebraic rate. Of course we do not recommend to implement this version of the algorithm. Instead, one should prefer a generalization of the adaptive finite element method. The well known adaptive finite element method always follows the loop

$$\text{SOLVE} \longrightarrow \text{ESTIMATE} \longrightarrow \text{MARK} \longrightarrow \text{REFINE}$$

where one solves the problem on the current discretization, estimates the error, locates the error in mark and refines the discretization in the last step.

We suggest a generalization in the following sense: In the estimate step one does not only calculate error estimators of the spatial discretization but of the errors coming from not running the Kačanov iteration long enough and from choosing the relaxation interval not large enough. After comparing these four estimators

1.  $\eta_{\varepsilon_-}^2 \approx \mathcal{J}_{(\varepsilon_-, \varepsilon_+)}(u_{(\varepsilon_-, \varepsilon_+)}) - \mathcal{J}_{(0, \varepsilon_+)}(u_{(0, \varepsilon_+)})$ ,
2.  $\eta_{\varepsilon_+}^2 \approx \mathcal{J}_{(\varepsilon_-, \varepsilon_+)}(u_{(\varepsilon_-, \varepsilon_+)}) - \mathcal{J}_{(\varepsilon_-, \infty)}(u_{(\varepsilon_-, \infty)})$ ,

## 5.2 Outlook On Adaptive Strategies for the Relaxation Parameter 67

3.  $\eta_h^2 \approx \mathcal{J}_\varepsilon(u_\varepsilon^h) - \mathcal{J}_\varepsilon(u_\varepsilon)$  where  $u_\varepsilon^h$  is the minimizer of  $\mathcal{J}_\varepsilon$  on a finite-dimensional solution space  $X_h \subset W_0^{1,\varphi_\varepsilon}(\Omega)$  and
4.  $\eta_{\text{Kač}}^2 \approx \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon)$ .

one decides either to

1. decrease  $\varepsilon_-$  or to
2. increase  $\varepsilon_+$  or to
3. refine the discrete solution space or to
4. do a further Kačanov step without changing anything.

The decision will be done according to the largest error estimator. Note that all quantities above are not computable. We will derive reliable error estimators in Section 5.3. Unfortunately, we can not prove efficiency for all of the estimators, in particular for the Kačanov estimator.

We start with an estimate that will be used for  $\eta_{\varepsilon_-}^2$  and  $\eta_{\varepsilon_+}^2$ . Let  $0 < \varepsilon_- < 1 < \varepsilon_+ < \infty$  and  $\delta := (\sigma\varepsilon_-, \theta^{-1}\varepsilon_+)$  for some  $\sigma, \theta \in (0, 1]$ . Then,

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) \approx \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon)$$

where the constants depend on  $\theta$  and  $\sigma$ . This equivalence basically states that the gain of switching from  $\mathcal{J}_\varepsilon(u_\varepsilon)$  to  $\mathcal{J}_\delta(u_\delta)$  behaves like the gain of just switching the energy from  $\mathcal{J}_\varepsilon$  to  $\mathcal{J}_\delta$  with the solution  $u_\varepsilon$ .

To prove this, we start with the following lemma. It basically states that the in the sense of Subsection 2.2.2 shifted N-function of  $\varphi_\varepsilon$  is just another N-function with changed relaxation interval.

**Lemma 5.5.** *For  $a > 0$  we have  $\varphi_{\varepsilon,a}(t) = \varphi_{\tilde{\varepsilon}}(t)$ , where  $\tilde{\varepsilon} = (\varepsilon_- \vee a \wedge \varepsilon_+, \varepsilon_+) = (\Pi_\varepsilon(a), \varepsilon_+)$ .*

*Proof.* Comparing

$$\Pi_\varepsilon(a \vee t) = \varepsilon_- \vee (a \vee t) \wedge \varepsilon_+$$

	$t \leq \varepsilon_-$	$t \in [\varepsilon_-, \varepsilon_+]$	$t \geq \varepsilon_+$
$a \leq \varepsilon_-$	$\varepsilon_-$	$t$	$\varepsilon_+$
$a \in [\varepsilon_-, \varepsilon_+]$	$a$	$a \vee t$	$\varepsilon_+$
$a \geq \varepsilon_+$	$\varepsilon_+$	$\varepsilon_+$	$\varepsilon_+$

to

$$\frac{\Pi_\varepsilon(a) \vee t \wedge \varepsilon_+ = (\varepsilon_- \vee a \wedge \varepsilon_+) \vee t \wedge \varepsilon_+}{\begin{array}{c|ccc} & t \leq \varepsilon_- & t \in [\varepsilon_-, \varepsilon_+] & t \geq \varepsilon_+ \\ \hline a \leq \varepsilon_- & \varepsilon_- & t & \varepsilon_+ \\ a \in [\varepsilon_-, \varepsilon_+] & a & a \vee t & \varepsilon_+ \\ a \geq \varepsilon_+ & \varepsilon_+ & \varepsilon_+ & \varepsilon_+ \end{array}}$$

we see that  $\varepsilon_- \vee (a \vee t) \wedge \varepsilon_+ = (\varepsilon_- \vee a \wedge \varepsilon_+) \vee t \wedge \varepsilon_+$ . Hence, for  $\tilde{\varepsilon} = (\varepsilon_- \vee a \wedge \varepsilon_+, \varepsilon_+)$  we get

$$\begin{aligned} \frac{\varphi'_{\varepsilon,a}(t)}{t} &= \frac{\varphi'_{\tilde{\varepsilon}}(a \vee t)}{a \vee t} = (\varepsilon_- \vee (a \vee t) \wedge \varepsilon_+)^{p-2} = (\varepsilon_- \vee (a \vee t) \wedge \varepsilon_+)^{p-2} \\ &= ((\varepsilon_- \vee a \wedge \varepsilon_+) \vee t \wedge \varepsilon_+)^{p-2} = \frac{\varphi'_{\tilde{\varepsilon}}(t)}{t}, \end{aligned}$$

so  $\varphi'_{\varepsilon,a} = \varphi'_{\tilde{\varepsilon}}$ . Since additionally  $\varphi_{\varepsilon,a}(0) = \varphi_{\tilde{\varepsilon}}(0)$  we get the statement.  $\square$

Note that this nice representation only holds with the shift as defined in this work but not with the shifts introduced in [DE08] or [RD07].

We will also need a representation of the complementary N-function  $\varphi_\varepsilon$ , which is presented in the next lemma.

**Lemma 5.6.** *The complementary N-function of  $\varphi_\varepsilon$  is given by*

$$\varphi_\varepsilon^*(t) = \begin{cases} \frac{1}{2}\varepsilon_-^{2-p}t^2 & \text{for } t \leq \varepsilon_-^{p-1} \\ \frac{1}{p^*}t^{p^*} + (\frac{1}{2} - \frac{1}{p^*})\varepsilon_-^p & \text{for } \varepsilon_-^{p-1} \leq t \leq \varepsilon_+^{p-1} \\ \frac{1}{2}\varepsilon_+^{2-p}t^2 - (\frac{1}{2} - \frac{1}{p^*})(\varepsilon_+^p - \varepsilon_-^p) & \text{for } t \geq \varepsilon_+^{p-1}. \end{cases}$$

*Proof.*  $\varphi'_\varepsilon$  is strictly monotone and therefore invertible. We recall

$$\varphi'_\varepsilon(t) = \begin{cases} \varepsilon_-^{p-2}t & \text{for } t \leq \varepsilon_- \\ t^{p-1} & \text{for } \varepsilon_- \leq t \leq \varepsilon_+ \\ \varepsilon_+^{p-2}t & \text{for } t \geq \varepsilon_+ \end{cases}$$

to see

$$(\varphi'_\varepsilon)^{-1}(t) = \begin{cases} \varepsilon_-^{2-p}t & \text{for } t \leq \varepsilon_-^{p-1} \\ t^{\frac{1}{p-1}} & \text{for } \varepsilon_-^{p-1} \leq t \leq \varepsilon_+^{p-1} \\ \varepsilon_+^{2-p}t & \text{for } t \geq \varepsilon_+^{p-1}. \end{cases}$$

With  $\varphi_\varepsilon^*(t) = \int_0^t (\varphi'_\varepsilon)^{-1}(\tau) d\tau$  we get the representation.  $\square$

Combining the last two lemmas we end up with the following technical lemma.

**Lemma 5.7.** *For  $\sigma, \theta \in (0, 1)$ ,  $\varepsilon_- \leq 1 \leq \varepsilon_+$ ,  $\delta_- = \sigma\varepsilon_-$  and  $\varepsilon_+ = \theta\delta_+$  we have*

$$(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|) \lesssim \kappa_\varepsilon(t) - \kappa_\delta(t)$$

where the constant depends on the choices of  $\sigma$  and  $\theta$ .

*Proof.* By the last Lemmas we get

$$(\varphi_{\delta,t})^*(s) = \begin{cases} \frac{1}{2}(\delta_- \vee t \wedge \delta_+)^{2-p}s^2 & \text{for } s \leq (\delta_- \vee t \wedge \delta_+)^{p-1} \\ \frac{1}{p^*}s^{p^*} + (\frac{1}{2} - \frac{1}{p^*})(\delta_- \vee t \wedge \delta_+)^p & \text{for } (\delta_- \vee t \wedge \delta_+)^{p-1} \leq s \leq \delta_+^{p-1} \\ \frac{1}{2}\delta_+^{2-p}s^2 - (\frac{1}{2} - \frac{1}{p^*})(\delta_+^p - (\delta_- \vee t \wedge \delta_+)^p) & \text{for } s \geq \delta_+^{p-1}. \end{cases}$$

Case  $t \in [0, \delta_-]$ : Under this assumption,

$$|\varphi'_\delta(t) - \varphi'_\varepsilon(t)| = \delta_-^{p-2}t - \varepsilon_-^{p-2}t \leq \delta_-^{p-1} = (\delta_- \vee t \wedge \delta_+)^{p-1}.$$

Hence,

$$\begin{aligned} (\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|) &= \frac{1}{2}\delta_-^{2-p}(\delta_-^{p-2} - \varepsilon_-^{p-2})^2t^2 \\ &\leq \frac{1}{2}\delta_-^{2-p}(\delta_-^{p-2} - \delta_-^{p-2}\sigma^{2-p})^2\delta_-^2 \\ &= \frac{1}{2}(1 - \sigma^{2-p})^2\delta_-^p. \end{aligned}$$

On the other hand,

$$\begin{aligned} \kappa_\varepsilon(t) - \kappa_\delta(t) &= \frac{1}{2}(\varepsilon_-^{p-2} - \delta_-^{p-2})t^2 + (\frac{1}{p} - \frac{1}{2})(\varepsilon_-^p - \delta_-^p) \\ &\geq \frac{1}{2}(\varepsilon_-^{p-2} - \delta_-^{p-2})\delta_-^2 + (\frac{1}{p} - \frac{1}{2})(\varepsilon_-^p - \delta_-^p) \\ &\geq ((\frac{1}{p} - \frac{1}{2})(\sigma^{-p} - 1) - \frac{1}{2}(1 - \sigma^{2-p}))\delta_-^p. \end{aligned}$$

This implies  $(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|) \leq c(\kappa_\varepsilon(t) - \kappa_\delta(t))$  with

$$c = \frac{(1 - \sigma^{2-p})^2}{2((\frac{1}{p} - \frac{1}{2})(\sigma^{-p} - 1) - \frac{1}{2}(1 - \sigma^{2-p}))} = \frac{(\sigma^p - \sigma^2)^2}{2\sigma^p((\frac{1}{p} - \frac{1}{2})(1 - \sigma^p) - \frac{1}{2}(\sigma^p - \sigma^2))}.$$

Case  $t \in [\delta_-, \varepsilon_-]$ : Here,

$$|\varphi'_\delta(t) - \varphi'_\varepsilon(t)| = (t^{p-1} - \varepsilon_-^{p-2}t) \leq t^{p-1} = (\delta_- \vee t \wedge \delta_+)^{p-1}.$$

Therefore,

$$\begin{aligned}
\sup_{t \in [\delta_-, \varepsilon_-]} \frac{(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|)}{\kappa_\varepsilon(t) - \kappa_\delta(t)} &= \sup_{t \in [\delta_-, \varepsilon_-]} \frac{\frac{1}{2}t^{2-p}(t^{p-2} - \varepsilon_-^{p-2})2t^2}{\frac{1}{2}\varepsilon_-^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_-^p - \frac{1}{p}t^p} \\
&= \sup_{\tau \in [\sigma, 1]} \frac{\frac{1}{2}\tau^{2-p}\varepsilon_-^{2-p}(\tau^{p-2}\varepsilon_-^{p-2} - \varepsilon_-^{p-2})2\tau^2\varepsilon_-^2}{\frac{1}{2}\varepsilon_-^p\tau^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_-^p - \frac{1}{p}\tau^p\varepsilon_-^p} \\
&= \sup_{\tau \in [\sigma, 1]} \frac{p\tau^{4-p}(\tau^{p-2} - 1)^2}{p\tau^2 + (2-p) - 2\tau^p} \\
&= \sup_{\tau \in [\sigma, 1]} \underbrace{\frac{p\tau^p - 2p\tau^2 + p\tau^{4-p}}{p\tau^2 + (2-p) - 2\tau^p}}_{:=h(\tau)}
\end{aligned}$$

Note that

$$\begin{aligned}
\lim_{\tau \rightarrow 1} h(\tau) &= \lim_{\tau \rightarrow 1} \frac{p\tau^p - 2p\tau^2 + p\tau^{4-p}}{p\tau^2 + (2-p) - 2\tau^p} \\
&= \lim_{\tau \rightarrow 1} \frac{p\tau^{p-1} - 4\tau + (4-p)\tau^{3-p}}{2\tau - 2\tau^{p-1}} \\
&= \lim_{\tau \rightarrow 1} \frac{p(p-1)\tau^{p-2} - 4 + (4-p)(3-p)\tau^{2-p}}{2 - 2(p-1)\tau^{p-2}} \\
&= \frac{p^2 - p - 4 + 12 - 3p - 4p + p^2}{4 - 2p} \\
&= \frac{p^2 - 4p + 4}{2 - p} = 2 - p.
\end{aligned}$$

Hence,  $h$  is continuous on  $[\sigma, 1]$  and therefore admits a maximum.

Case  $t \in [\varepsilon_-, \varepsilon_+]$ : Since  $|\varphi'_\delta(t) - \varphi'_\varepsilon(t)| = 0$  we get  $(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|) = 0$ .

Case  $t \in [\varepsilon_+, \delta_+]$ : In this case, we have

$$|\varphi'_\delta(t) - \varphi'_\varepsilon(t)| = \varepsilon_+^{p-2}t - t^{p-1} = (\varepsilon_+^{p-2} - t^{p-2})t$$

which is strictly increasing and hence  $|\varphi'_\delta(t) - \varphi'_\varepsilon(t)| \in [0, (\theta^{p-2} - 1)\delta_+^{p-1}]$ . We need to consider subcases.

Subcase  $0 \leq (\varepsilon_+^{p-2} - t^{p-2})t \leq t^{p-1}$ : Since  $t \geq \varepsilon_+$  we can write  $t = \tau\varepsilon_+$  for some  $\tau \geq 1$ . Hence,

$$(\varepsilon_+^{p-2} - t^{p-2})t \leq t^{p-1} \iff (1 - \tau^{p-2})\tau \leq \tau^{p-1} \iff 2^{\frac{1}{2-p}} \leq \tau.$$

Therefore,

$$\begin{aligned}
 \sup_{t \in [\varepsilon_+, 2^{\frac{1}{2-p}} \varepsilon_+]} \frac{(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|)}{\kappa_\varepsilon(t) - \kappa_\delta(t)} &= \sup_{t \in [\varepsilon_+, 2^{\frac{1}{2-p}} \varepsilon_+]} \frac{\frac{1}{2}t^{2-p}(\varepsilon_+^{p-2} - t^{p-2})^2 t^2}{\frac{1}{2}\varepsilon_+^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}t^p} \\
 &= \sup_{\tau \in [1, 2^{\frac{1}{2-p}}]} \frac{\frac{1}{2}\tau^{2-p}\varepsilon_+^{2-p}(\varepsilon_+^{p-2} - \tau^{p-2}\varepsilon_+^{p-2})^2 \tau^2 \varepsilon_+^2}{\frac{1}{2}\varepsilon_+^{p-2}\tau^2 \varepsilon_+^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}\tau^p \varepsilon_+^p} \\
 &= \sup_{\tau \in [1, 2^{\frac{1}{2-p}}]} \frac{\frac{1}{2}\tau^{2-p}\varepsilon_+^{2-p}(\varepsilon_+^{p-2} - \tau^{p-2}\varepsilon_+^{p-2})^2 \tau^2 \varepsilon_+^2}{\frac{1}{2}\varepsilon_+^{p-2}\tau^2 \varepsilon_+^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}\tau^p \varepsilon_+^p} \\
 &= \sup_{\tau \in [1, 2^{\frac{1}{2-p}}]} \frac{\frac{1}{2}\tau^{2-p}(1 - \tau^{p-2})^2 \tau^2}{\frac{1}{2}\tau^2 + (\frac{1}{p} - \frac{1}{2}) - \frac{1}{p}\tau^p} \\
 &= \sup_{\tau \in [1, 2^{\frac{1}{2-p}}]} \frac{p\tau^p - 2p\tau^2 + p\tau^{4-p}}{p\tau^2 + (2-p) - 2\tau^p}
 \end{aligned}$$

Note that this is the supremum over the same function as in the case  $t \in [\delta_-, \varepsilon_-]$ . Hence, we already know  $\lim_{\tau \rightarrow 1} h(\tau) = 1$ . Additionally, the denominator satisfies  $p\tau^2 + (2-p) - 2\tau^p > 0$  for  $\tau > 1$  since it is strictly increasing in  $\tau$  and is zero for  $\tau = 1$ . Therefore,  $h$  is continuous and attains a maximum on  $[1, 2^{\frac{1}{2-p}}]$ .

Subcase  $t^{p-1} \leq (\varepsilon_+^{p-2} - t^{p-2})t \leq \delta_+^{p-1}$ : Again, we use  $t = \tau\varepsilon_+$  for  $\theta \geq 1$  to see that

$$(\varepsilon_+^{p-2} - t^{p-2})t \leq \delta_+^{p-1} \leq \delta_+^{p-1} \iff (1 - \tau^{p-2})\tau \leq \theta^{1-p}.$$

Since  $(1 - \tau^{p-2})\tau$  is increasing as the product of two non-negative increasing functions this is equivalent to the existence of a  $T_\theta$  such that  $\tau \leq T_\theta$ . Hence,

$$\begin{aligned}
 \sup_{t \text{ w.r.t. subcase}} \frac{(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|)}{\kappa_\varepsilon(t) - \kappa_\delta(t)} &= \sup_{\tau \in [2^{\frac{1}{2-p}}, T_\theta]} \frac{\frac{1}{p'}(\varepsilon_+^{p-2} - t^{p-2})^{p'} t^{p'} + (\frac{1}{2} - \frac{1}{p'})t^p}{\frac{1}{2}\varepsilon_+^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}t^p} \\
 &= \sup_{\tau \in [2^{\frac{1}{2-p}}, T_\theta]} \frac{\frac{1}{p'}(\varepsilon_+^{p-2} - \varepsilon_+^{p-2}\tau^{p-2})^{p'} \varepsilon_+^{p'} \tau^{p'} + (\frac{1}{2} - \frac{1}{p'})\varepsilon_+^p \tau^p}{\frac{1}{2}\varepsilon_+^{p-2}\varepsilon_+^2 \tau^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}\tau^p \varepsilon_+^p} \\
 &= \sup_{\tau \in [2^{\frac{1}{2-p}}, T_\theta]} \frac{2p((\tau - \tau^{p-1})^{p'} + (\frac{1}{2} - \frac{1}{p'})\tau^p)}{\frac{1}{2}\tau^2 + (\frac{1}{p} - \frac{1}{2}) - \frac{1}{p}\tau^p}.
 \end{aligned}$$

Again, the denominator is not zero. Hence, the fraction is continuous in  $\tau$  and therefore admits a maximum depending on  $p$  and  $\theta$  only.

Subcase  $\delta_+^{p-1} \leq (\varepsilon_+^{p-2} - t^{p-2})t \leq (\theta^{p-2} - 1)\delta_+^{p-1}$ : With the same argumentation as in the last subcase we write  $t = \tau\varepsilon_+$  for  $\tau \geq 1$  and get that the subcase condition is equivalent to  $\tau \in [T_\theta, R_\theta]$ . Hence,

$$\begin{aligned}
& \sup_{t \text{ w.r.t. subcase}} \frac{(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|)}{\kappa_\varepsilon(t) - \kappa_\delta(t)} \\
&= \sup_{t \text{ w.r.t. subcase}} \frac{\frac{1}{2}\delta_+^{2-p}(\varepsilon_+^{p-2} - t^{p-2})^2 t^2 - (\frac{1}{2} - \frac{1}{p})(\delta_+^p - t^p)}{\frac{1}{2}\varepsilon_+^{p-2}t^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}t^p} \\
&= \sup_{\tau \in [T_\theta, R_\theta]} \frac{\frac{1}{2}\varepsilon_+^{2-p}\theta^{p-2}(\varepsilon_+^{p-2} - \varepsilon_+^{p-2}\tau^{p-2})^2 \varepsilon_+^2 \tau^2 - (\frac{1}{2} - \frac{1}{p})(\varepsilon_+^p \theta^{-p} - \varepsilon_+^p \tau^p)}{\frac{1}{2}\varepsilon_+^{p-2}\tau^2 \varepsilon_+^2 + (\frac{1}{p} - \frac{1}{2})\varepsilon_+^p - \frac{1}{p}\tau^p \varepsilon_+^p} \\
&= \sup_{\tau \in [T_\theta, R_\theta]} \frac{\frac{1}{2}\theta^{p-2}(1 - \tau^{p-2})^2 \tau^2 - (\frac{1}{2} - \frac{1}{p})(\theta^{-p} - \tau^p)}{\frac{1}{2}\tau^2 + (\frac{1}{p} - \frac{1}{2}) - \frac{1}{p}\tau^p}
\end{aligned}$$

Again, the fraction is continuous in  $\tau$  on a compact domain and therefore admits a maximum depending on  $p$  and  $\theta$ .

Case  $t \in [\delta_+, \infty)$ : Now,

$$|\varphi'_\delta(t) - \varphi'_\varepsilon(t)| = \varepsilon_+^{p-2}t - \delta_+^{p-2}t = t\delta_+^{p-2}(\theta^{p-2} - 1) \geq \delta_+^{p-1}$$

and so

$$\begin{aligned}
\sup_{t \geq \delta_+} \frac{(\varphi_{\delta,t})^*(|\varphi'_\delta(t) - \varphi'_\varepsilon(t)|)}{\kappa_\varepsilon(t) - \kappa_\delta(t)} &= \sup_{t \geq \delta_+} \frac{\frac{1}{2}\delta_+^{p-2}(\theta^{p-2} - 1)^2 t^2}{\frac{1}{2}(\theta^{p-2} - 1)\delta_+^{p-2}t^2 - (\frac{1}{p} - \frac{1}{2})(1 - \theta^p)\delta_+^p} \\
&= \sup_{t \geq \delta_+} \frac{\frac{1}{2}\delta_+^{p-2}(\theta^{p-2} - 1)^2}{\frac{1}{2}(\theta^{p-2} - 1)\delta_+^{p-2} - (\frac{1}{p} - \frac{1}{2})(1 - \theta^p)\delta_+^p t^{-2}} \\
&\leq \frac{\frac{1}{2}\delta_+^{p-2}(\theta^{p-2} - 1)^2}{\frac{1}{2}(\theta^{p-2} - 1)\delta_+^{p-2} - (\frac{1}{p} - \frac{1}{2})(1 - \theta^p)\delta_+^p \delta_+^{-2}} \\
&= \frac{\frac{1}{2}(\theta^{p-2} - 1)^2}{\frac{1}{2}(\theta^{p-2} - 1) - (\frac{1}{p} - \frac{1}{2})(1 - \theta^p)}. \quad \square
\end{aligned}$$

With that, we can prove the following theorem.

**Theorem 5.8.** *Let  $0 < \varepsilon_- < 1 < \varepsilon_+ < \infty$  and  $\delta := (\sigma\varepsilon_-, \theta^{-1}\varepsilon_+)$  for some  $\sigma, \theta \in (0, 1]$ . Then,*

$$\mathcal{J}_\delta(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) \lesssim \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon)$$

where the constants additionally depend on  $\sigma$  and  $\theta$ .



*Proof.* We use Theorem 2.42, the equations for  $u_\varepsilon$  and  $u_\delta$  and Young's Inequality to deduce for  $\gamma > 0$  the estimate

$$\begin{aligned}
 \mathcal{J}_\delta(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) &\approx \int_{\Omega} (A_\delta(\nabla u_\varepsilon) - A_\delta(\nabla u_\delta)) \nabla(u_\varepsilon - u_\delta) \, dx \\
 &= \int_{\Omega} (A_\delta(\nabla u_\varepsilon) - A_\varepsilon(\nabla u_\varepsilon)) \nabla(u_\varepsilon - u_\delta) \, dx \\
 &\leq \gamma \int_{\Omega} \varphi_{\delta, |\nabla u_\varepsilon|} (|\nabla(u_\varepsilon - u_\delta)|) \, dx \\
 &\quad + c_\gamma \int_{\Omega} (\varphi_{\delta, |\nabla u_\varepsilon|})^* (|\varphi'_\delta(|\nabla u_\varepsilon|) - \varphi'_\varepsilon(|\nabla u_\varepsilon|)|) \, dx
 \end{aligned}$$

Choosing  $\gamma$  small enough, applying Lemma 2.41, Theorem 2.42 and Lemma 5.7 we get

$$\begin{aligned}
 \mathcal{J}_\delta(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) &\lesssim \int_{\Omega} \kappa_\varepsilon(|\nabla u_\varepsilon|) - \kappa_\delta(|\nabla u_\varepsilon|) \, dx \\
 &= \int_{\Omega} \kappa_\varepsilon(|\nabla u_\varepsilon|) - \kappa_\delta(|\nabla u_\varepsilon|) \, dx - \langle f, u_\varepsilon - u_\delta \rangle \\
 &= \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon). \quad \square
 \end{aligned}$$

We use the last theorem to deduce a corollary describing the relation as stated on page 67.

**Corollary 5.9.** *Let  $0 < \varepsilon_- < 1 < \varepsilon_+ < \infty$  and  $\delta := (\sigma\varepsilon_-, \theta^{-1}\varepsilon_+)$  for some  $\sigma, \theta \in (0, 1]$ . Then,*

$$\mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) \approx \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon)$$

where the constants additionally depend on  $\sigma$  and  $\theta$ .

*Proof.* With Theorem 5.8 and the minimization property of  $u_\delta$  we directly get

$$\begin{aligned}
 \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) &= \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon) + \mathcal{J}_\delta(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) \\
 &\lesssim \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon) + \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon) \\
 &\approx \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon) \\
 &\leq \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\delta). \quad \square
 \end{aligned}$$

## 5.3 Numerical Examples

### 5.3.1 Estimators and the Adaptive Algorithm

In this subsection we want to give the heuristics for the estimators we used to run the fully adaptive Kačanov algorithm as stated on page 62. We are interested in three types of estimators. They shall cover the distances between the solution  $u$ , the solution to the relaxed problem  $u_\varepsilon$ , the finite element solution of the relaxed problem  $u_\varepsilon^h$  and the current approximation to the solution computed by the algorithm  $v_n$ :

$$v_n \xleftarrow{\eta_{\text{Ka}\bar{c}}^2} u_\varepsilon^h \xleftarrow{\eta_h^2} u_\varepsilon \xleftarrow{\eta_{\varepsilon-}^2 \text{ and } \eta_{\varepsilon+}^2} u$$

For the error introduced by the relaxation interval we use heuristics based on Corollary 5.9 (for an interpretation see page 67) stating

$$\begin{aligned} \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\delta) &\approx \mathcal{J}_\varepsilon(u_\varepsilon) - \mathcal{J}_\delta(u_\varepsilon) \\ &= \int_{\{|\nabla u_\varepsilon| < \varepsilon_-\}} \kappa_\varepsilon(|\nabla u_\varepsilon|) - \kappa_\delta(|\nabla u_\varepsilon|) dx \\ &\quad + \int_{\{|\nabla u_\varepsilon| > \varepsilon_+\}} \kappa_\varepsilon(|\nabla u_\varepsilon|) - \kappa_\delta(|\nabla u_\varepsilon|) dx, \end{aligned}$$

because  $\langle f, u_\varepsilon - u_\varepsilon \rangle = 0$  and  $\kappa_\varepsilon$  and  $\kappa_\delta$  coincide on  $\varepsilon$ . Unfortunately, this is not computable during the iteration since  $u_\varepsilon$  is not known. We replace it by the best approximation that is available – namely the current iterated which approximates  $u_\varepsilon$  actually better than  $u$  – and define

$$\eta_{\varepsilon-}^2(v) := \int_{\{|\nabla v| < \varepsilon_-\}} \kappa_\varepsilon(|\nabla v|) - \kappa_\delta(|\nabla v|) dx$$

and

$$\eta_{\varepsilon+}^2(v) := \int_{\{|\nabla v| > \varepsilon_+\}} \kappa_\varepsilon(|\nabla v|) - \kappa_\delta(|\nabla v|) dx.$$

Note that with these definitions we get

$$\mathcal{J}_\varepsilon(v) - \mathcal{J}_\delta(v) = \eta_{\varepsilon-}^2(v) + \eta_{\varepsilon+}^2(v).$$

For the calculation of the estimators  $\eta_{\varepsilon_-}^2(v)$  and  $\eta_{\varepsilon_+}^2(v)$  we used  $\delta_- := 10^{-5}\varepsilon_-$  and  $\delta_+ := 10^5\varepsilon_+$ .

For the discretization error we use the error estimators presented in [DK08]. Note that the results in this paper not only apply to the  $p$ -Poisson equation, but more general to the  $\varphi$ -Poisson equation. In particular, the error estimators are valid for the error between  $u_\varepsilon$  and  $u_\varepsilon^h$  where  $u_\varepsilon^h$  is the minimizer of  $\mathcal{J}_\varepsilon$  on a discrete space  $X_h \subset W_0^{1,p}(\Omega)$ . We will use the space of piecewise linear functions  $\mathcal{P}_1$ . For each triangle  $T$  of the discretization the definition of the error estimator for a function  $v_h \in X_h$  reads as

$$\eta_h^2(v_h, T) := \int_T (\varphi_{\varepsilon, |\nabla v_h|})^* (h_T |f|) dx + \sum_{\gamma \subset \partial T} \int_\gamma h_\gamma |V_\varepsilon(\nabla v_h)|^2 dx$$

where the sum over  $\gamma \subset \partial T$  describes the sum over all faces of  $T$ ,  $h_T$  being the diameter of  $T$  and  $h_\gamma$  being the diameter of the face  $\gamma$ . Finally, we use

$$\eta_h^2(v_h) := \sum_{T \in \mathcal{T}_h} \eta_h^2(v_h, T)$$

as an estimator for the global discretization error. For a more detailed description we refer to [DK08]. Note that this estimator requires that  $f$  is a function but not just a functional.

Finally, we need an estimator for the gain that can be achieved by just doing a Kačanov step without changing  $\varepsilon_-$  or  $\varepsilon_+$ . So let us assume that  $v_{n+1}$  admits

$$\int_\Omega \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} \nabla v_{n+1} \nabla \xi dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega).$$

Using this, the equation of  $u_\varepsilon$ , Theorem 2.42 and Young's Inequality we obtain for any  $\gamma > 0$  the estimate

$$\begin{aligned} \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) &\leq \int_\Omega (A_\varepsilon(\nabla v_n) - A_\varepsilon(\nabla u_\varepsilon)) \nabla(v_n - u_\varepsilon) dx \\ &= \int_\Omega \left( \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} \nabla v_n - \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} \nabla v_{n+1} \right) \nabla(v_n - u_\varepsilon) dx \\ &\leq \gamma \int_\Omega \varphi_{\varepsilon, |\nabla v_n|} (|\nabla(v_n - u_\varepsilon)|) dx \\ &\quad + c_\gamma \int_\Omega (\varphi_{\varepsilon, |\nabla v_n|})^* \left( \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - v_{n+1})| \right) dx. \end{aligned}$$

Choosing  $\gamma$  small enough and applying Lemma 2.41 we deduce

$$\mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon) \lesssim \int_{\Omega} (\varphi_{\varepsilon, |\nabla v_n|})^* \left( \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - v_{n+1})| \right) dx. \quad (5.5)$$

Note that  $\mathcal{J}_\varepsilon(v_{n+1}) - \mathcal{J}_\varepsilon(u_\varepsilon) \leq \mathcal{J}_\varepsilon(v_n) - \mathcal{J}_\varepsilon(u_\varepsilon)$  and that the right hand side of (5.5) is computable. Hence, up to a constant we have a computable upper bound for the gain that can be achieved by just performing Kačanov steps arbitrarily often. This gives reason to define

$$\eta_{\text{Kač}}^2(v_{n+1}) := \int_{\Omega} (\varphi_{\varepsilon, |\nabla v_n|})^* \left( \frac{\varphi'_\varepsilon(|\nabla v_n|)}{|\nabla v_n|} |\nabla(v_n - v_{n+1})| \right) dx.$$

Of course, the calculation of  $\eta_{\text{Kač}}^2(v_{n+1})$  requires a function  $v_n$ . Due to the lack of  $v_{-1}$  we define  $\eta_{\text{Kač}}^2(v_0) := 0$  manually.

Now having estimators for the relaxation parameters, the remaining Kačanov error and the discretization we suggest the following algorithm.

**Algorithm:** The adaptive relaxed  $p$ -Kačanov algorithm

**Data:**  $\Omega \subset \mathbb{R}^d$ ;  $f \in L^1_{\text{loc}}(\Omega)$ ; finite-dimensional  $X_h \subset \mathcal{P}_1 \subset W_0^{1,2}(\Omega)$ ;

**Result:** Approximate solution of the  $p$ -Poisson problem (2.5).

$n := 0$ ; costs = 0;  $\varepsilon_- := 1$ ;  $\varepsilon_+ := 1$ ;

Calculate  $v_0$  as solution of

$$\int_{\Omega} \nabla v_0 \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega);$$

**while** desired accuracy is not achieved yet **do**

    Calculate  $\eta_{\text{Kač}}^2(v_n)$ ,  $\eta_{\varepsilon_-}^2(v_n)$ ,  $\eta_{\varepsilon_+}^2(v_n)$ ,  $\eta_h^2(v_n)$  as defined above;

**if**  $\eta_{\varepsilon_-}^2(v_n) = \max\{\eta_{\text{Kač}}^2(v_n), \eta_{\varepsilon_-}^2(v_n), \eta_{\varepsilon_+}^2(v_n), \eta_h^2(v_n)\}$  **then**

        Update  $\varepsilon_- \rightsquigarrow 0.8 \cdot \varepsilon_-$ ;

**if**  $\eta_{\varepsilon_+}^2(v_n) = \max\{\eta_{\text{Kač}}^2(v_n), \eta_{\varepsilon_-}^2(v_n), \eta_{\varepsilon_+}^2(v_n), \eta_h^2(v_n)\}$  **then**

$\varepsilon_+ \rightsquigarrow 1.25 \cdot \varepsilon_+$ ;

**if**  $\eta_h^2(v_n) = \max\{\eta_{\text{Kač}}^2(v_n), \eta_{\varepsilon_-}^2(v_n), \eta_{\varepsilon_+}^2(v_n), \eta_h^2(v_n)\}$  **then**

        Perform Dörfler marking with refine fraction 0.2;

        Perform red/green refinement to obtain finer  $X_h$ ;

    Calculate  $v_{n+1}$  as solution of

$$\int_{\Omega} \Pi_{\varepsilon}(|\nabla v_n|)^{p-2} \nabla v_{n+1} \nabla \xi \, dx = \langle f, \xi \rangle \quad \forall \xi \in W_0^{1,2}(\Omega);$$

    Update costs  $\rightsquigarrow$  costs + degrees of freedom of  $X_h$ ;

    Update  $n \rightsquigarrow n + 1$ ;

**end**

To ensure a reasonable performance, we multiplied the estimators  $\eta_{\varepsilon_+}^2$  by  $10^2$  and  $\eta_h^2$  by  $10^{-2}$  in all the following experiments. The algorithm was implemented using DUNE PDELab combined with DUNE's grid manager UG grid. The Dörfler marking strategy goes back to the fundamental paper [Dör96]. One can find a description of the red/green refinement in [Cer04].

Furthermore we want to note that in each experiment we chose a very coarse initial triangulation of  $\Omega$  corresponding to a very small space  $X_h$  on purpose. The reason to point out this explicitly is that at least the theory for adaptive schemes sometimes requires a sufficiently fine initial triangulation.

Usually, the  $x$ -axis is marked with the degrees of freedom. However, this is not suitable in this setting, since each iteration that is not a refinement produces computation costs but does not increase the degrees of freedom. Hence, we think it is more adequate to mark the  $x$ -axis with the computational costs for solving the linear systems. As defined in the algorithm, the costs are the sum over all degrees of freedom calculated up to the recent iteration step.

The next table gives the mathematical expressions for the terms in the legends of the graphs in the examples.

legend	term
Jeps	$\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u)$
est_decLo	$\eta_{\varepsilon_-}^2(v_n)$
est_incUp	$\eta_{\varepsilon_+}^2(v_n)$
est_refine	$\eta_h^2(v_n)$
est_kacanov	$\eta_{\text{Kac}}^2(v_n)$
lowerShift	$\varepsilon_-$
upperShift	$\varepsilon_+$
inv_costs	$\text{costs}^{-1}$

The quantity  $\text{costs}^{-1}$  is additionally plotted in the figures to make it possible to compare the data to linear decay in the double-logarithmic scaling. It is multiplied with a suitable constant to translate the image to a suitable height. The factor is chosen identically over the range of  $p$  for each example to keep the graphs comparable easier.

For every example we use  $p \in \{\frac{3}{2}, \frac{4}{3}, \frac{8}{7}, \frac{16}{15}\}$  since this corresponds to the choices  $p' \in \{3, 4, 8, 16\}$ .

### 5.3.2 Example: Bump

The first example is given on  $\Omega := (-1, 1)^2$  by the solution “bump”

$$\begin{aligned} u : \Omega &\rightarrow \mathbb{R} \\ x &\mapsto (x_1^2 - 1)(x_2^2 - 1). \end{aligned}$$

This is a relatively nice example meaning that  $u \in C^\infty(\overline{\Omega})$  and  $u = 0$  on  $\partial\Omega$ . Nevertheless, we have  $|\nabla u(0)| = 0$ , so the lower constraint will strike anyway. One can see very well that even for  $p = \frac{16}{15} = 1.0\overline{6}$  the algorithm almost produces a linear decay of  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u)$ , where  $\mathcal{J}(u)$  was calculated numerically as the integral over the projection of the exact solution to a finer finite-dimensional space.

The linear parts of  $\eta_{\varepsilon_+}^2(v_n)$  (for example in Figure 5.1 between  $10^2$  and  $10^4$ ) always indicate that the quantity is equal to zero which can not be represented due to the logarithmic scale of the  $y$ -axis.

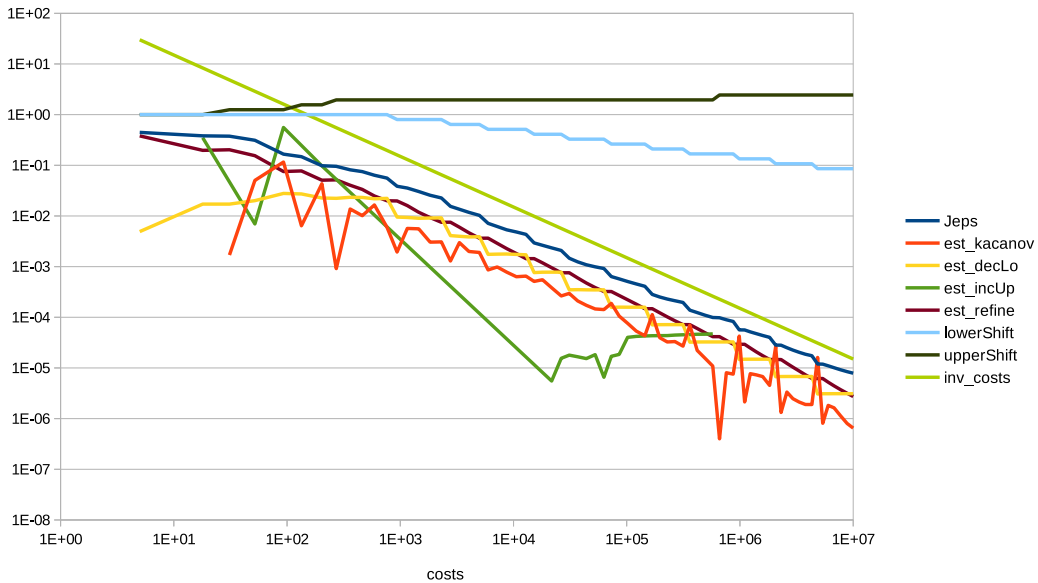


Figure 5.1: Performance of the algorithm for the bump and  $p = \frac{3}{2}$ .

Note that  $\eta_{\text{Kac}}^2(v_n)$  always jumps, when  $\varepsilon_-$  is decreased. This is due to the fact that after changing  $\varepsilon_-$ , the next iterated approximates a different  $u_\varepsilon$  than the previous iterated. Hence, the difference between those to solutions and therefore  $\eta_{\text{Kac}}^2(v_n)$  is relatively large. This effect will occur in all examples.

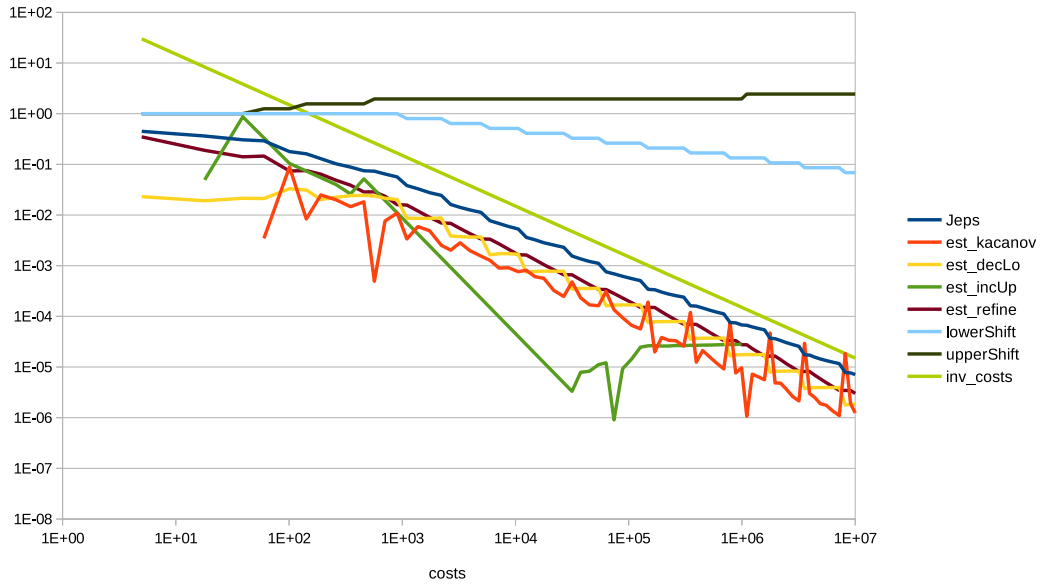


Figure 5.2: Performance of the algorithm for the bump and  $p = \frac{4}{3}$ .

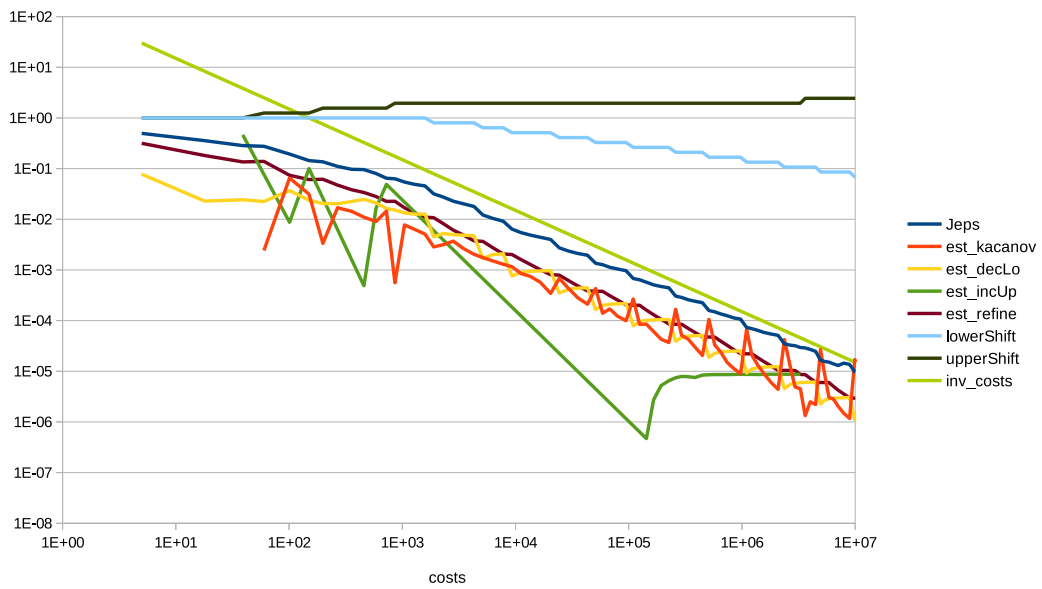


Figure 5.3: Performance of the algorithm for the bump and  $p = \frac{8}{7}$ .



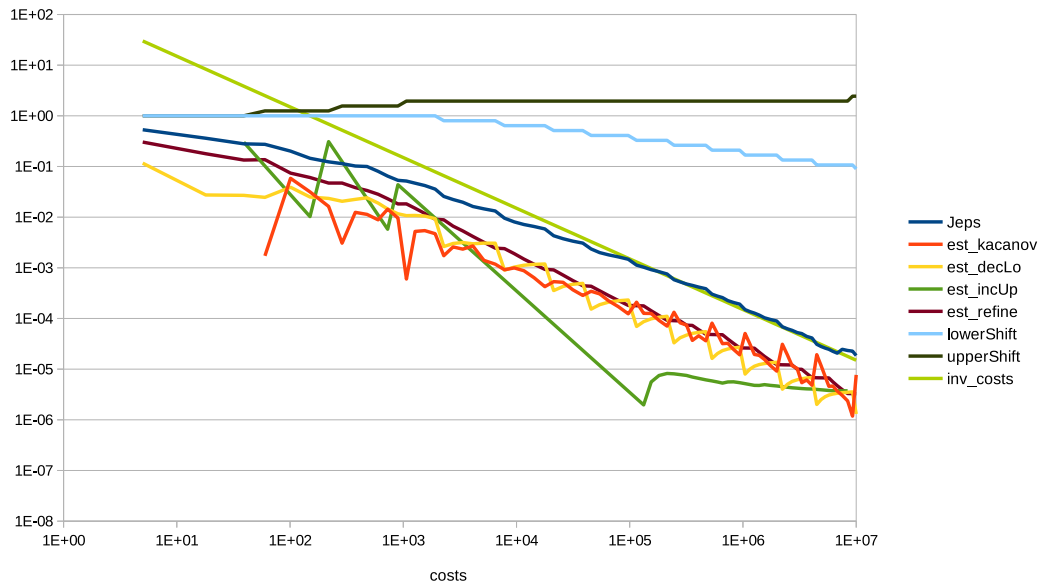


Figure 5.4: Performance of the algorithm for the bump and  $p = \frac{16}{15}$ .

### 5.3.3 Example: Needle

For  $\Omega := (-1, 1)^2$  we define the function “needle”

$$\begin{aligned} u : \Omega &\rightarrow \mathbb{R} \\ x &\mapsto |x|^{1-\frac{1}{p}} - 1. \end{aligned}$$

Note that in this case  $u$  does not admit zero boundary values. Hence, this example can be interpreted as a test for the conjecture, whether the algorithm does need zero boundary values or not.

The exponent is chosen depending on  $p$ . This is done to reach the borderline case for  $|V(\nabla u)| \approx |x|^{-\frac{1}{2}} \in W^{\frac{1}{2}}L^{2,\infty}(\Omega)$  meaning that a half derivative of  $|V(\nabla u)|$  is still in the Lorentz space  $L^{2,\infty}(\Omega)$ . This refers to an example in [BDK12] where it is shown in Figure 3 that adaptive mesh refinement admits the optimal convergence rate in terms of degrees of freedoms for this regularity of the solution. Note that this is the same regularity as the  $p$ -harmonic function on the slit domain admits. Again,  $\mathcal{J}(u)$  was calculated by integrating the projection of  $u$  onto a finer finite element space.

Furthermore, it is interesting to test our algorithm with this example since  $|\nabla u| \approx |x|^{-\frac{1}{p}} \notin L^\infty(\Omega)$ . Therefore,  $\varepsilon_+$  plays an important role.

One can see over all chosen values of  $p$  that  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u)$  roughly behaves like  $\text{costs}^{-1}$ . Since we have  $\text{costs}^{-1} \leq \text{dofs}^{-1}$  this implies that our algorithm converges with the same speed as the optimal adaptive finite element method. This is a great achievement since – contrarily to the test performed in [BDK12] – our algorithm does not need to solve a non-linear system in each step.

Additionally one can see in the graphs that  $\varepsilon_+$  needs to be adjusted regularly. Note that  $\varepsilon_+$  starts to be raised later if  $p$  gets smaller (compare Figure 5.9 to Figure 5.12). Indeed, the singularity of  $|\nabla u| \approx |x|^{-\frac{1}{p}}$  becomes worse, but the domain where  $|\nabla u|$  is large becomes smaller as  $p$  gets smaller. Hence, the impact of the singularity of  $|\nabla u|$  to the energy decreases as  $p$  gets smaller. Furthermore, the non-linearity gets stronger as  $p$  decreases such that it might be necessary to perform more Kačanov steps.

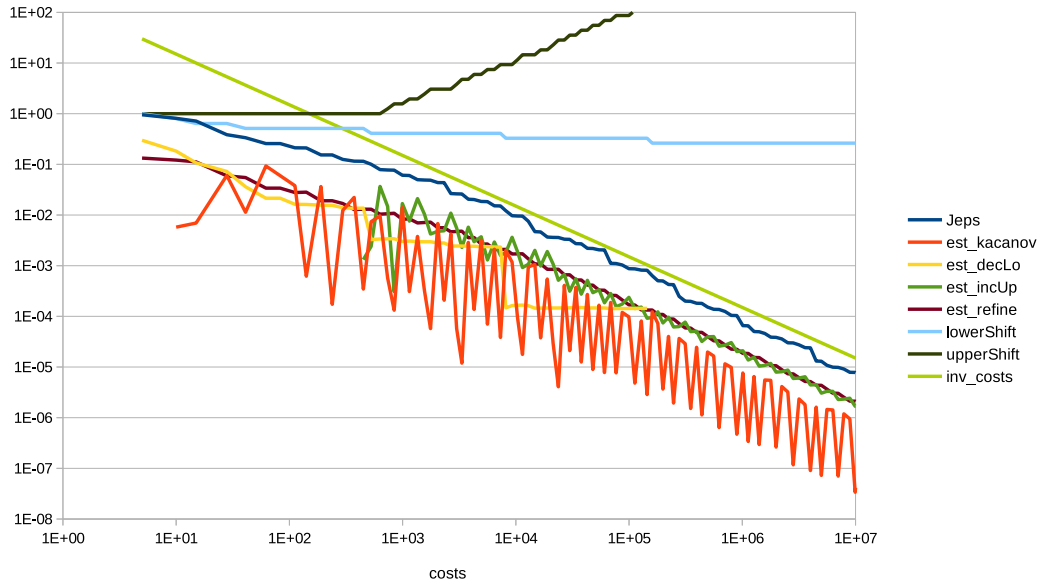


Figure 5.5: Performance of the algorithm for the needle and  $p = \frac{3}{2}$ .

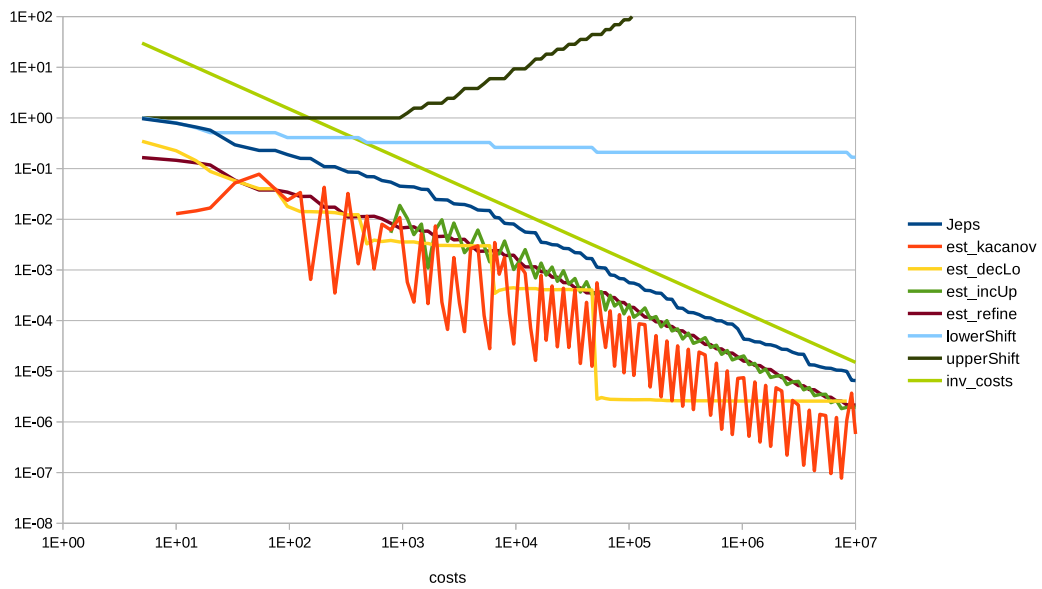


Figure 5.6: Performance of the algorithm for the needle and  $p = \frac{4}{3}$ .

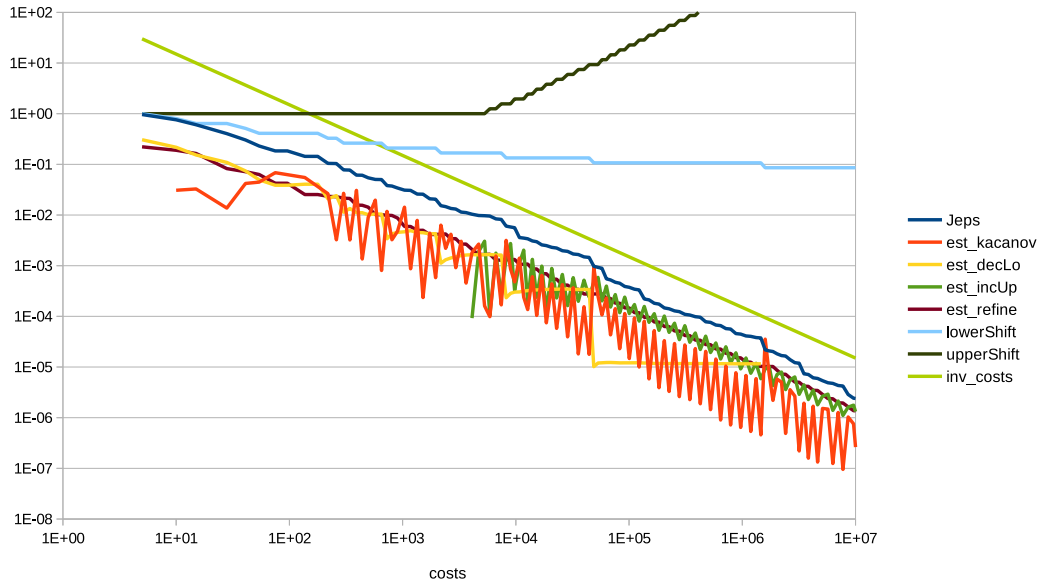


Figure 5.7: Performance of the algorithm for the needle and  $p = \frac{8}{7}$ .

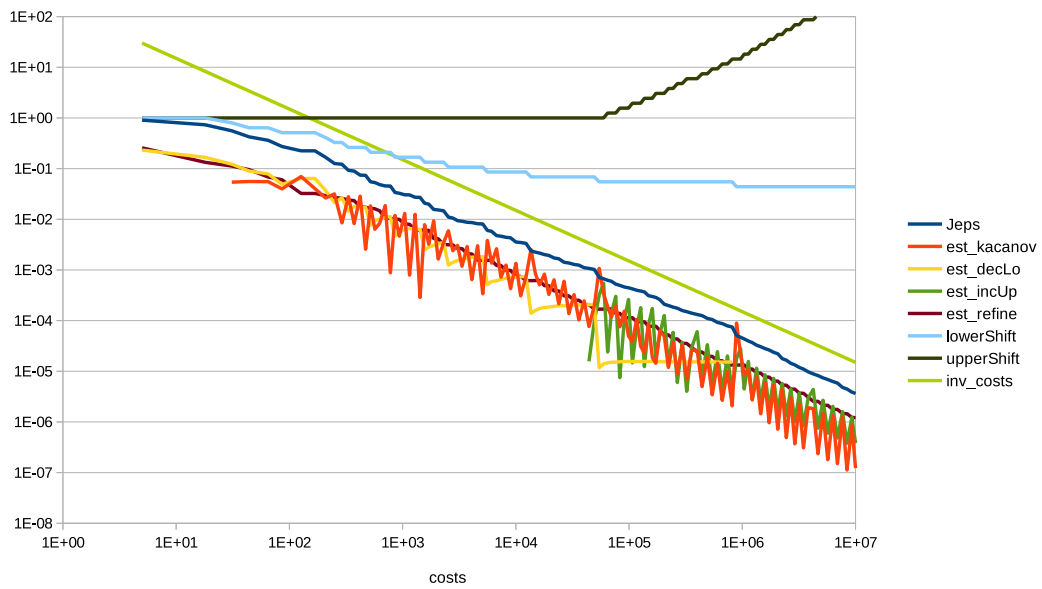


Figure 5.8: Performance of the algorithm for the needle and  $p = \frac{16}{15}$ .

### 5.3.4 Example: Constant Force

For  $\Omega := (-1, 1)^2 \setminus [0, 1]^2$  we define  $u \in W_0^{1,p}(\Omega)$  as the solution of

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \nabla \xi \, dx = \langle 2, \xi \rangle \quad \forall \xi \in W_0^{1,p}(\Omega).$$

This is outstanding because of two reasons. Firstly, we chose  $\Omega$  as the L-shaped domain and secondly,  $u$  is not known for this choice of  $f$ . The energy  $\mathcal{J}(u)$  was estimated by running the algorithm up to costs =  $2 \cdot 10^7$ .

This model problem is widely used in the numerics of elliptic partial differential equations. It is known that the gradient of the solution admits a singularity in the corner  $0 \in \mathbb{R}^2$ .

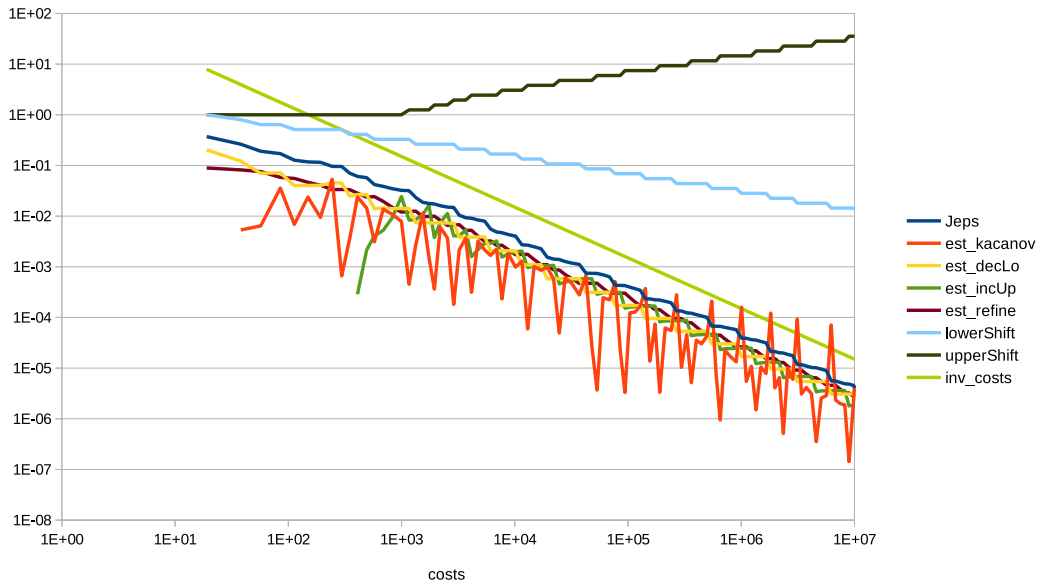


Figure 5.9: Performance of the algorithm for constant force and  $p = \frac{3}{2}$ .

The sequence of peaks of the Kačanov estimator behaves interestingly, too. It seems to decrease linearly in the double logarithmic scale but again much slower as  $p$  gets smaller. Furthermore one could think that  $\eta_{\text{Kač}}^2(v)$  overestimates every time  $\varepsilon_-$  was updated.

Compared to the example with the needle, once more the impact of  $\varepsilon_+$  seems to decrease as  $p$  gets small. An extreme example is shown for constant force and  $p = \frac{16}{15}$  in Figure 5.12.

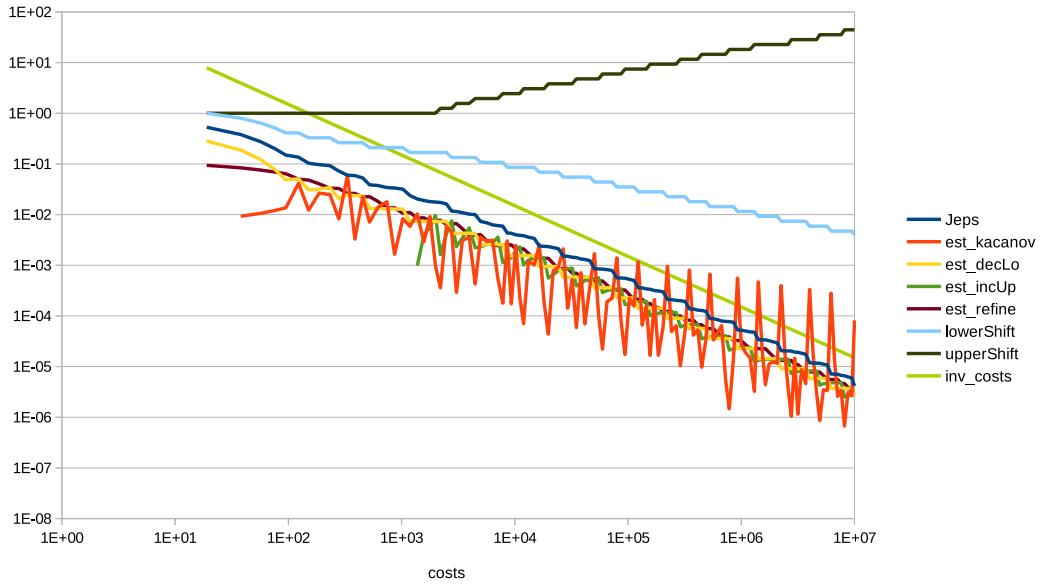


Figure 5.10: Performance of the algorithm for constant force and  $p = \frac{4}{3}$ .

In this example,  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u)$  behaves linear on the double logarithmic scale very nice but  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u) \lesssim \text{costs}^{-1}$  does not seem to hold (whereas the estimate  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u) \lesssim \text{costs}^{-\alpha}$  could still be possible). However, marking the  $x$ -axis with dofs instead of costs it seems that  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u) \lesssim \text{dofs}^{-1}$  does hold true (compare Figure 5.12 to Figure 5.13), so neglecting the steps where  $X_h$  is not refined compensates this effect.

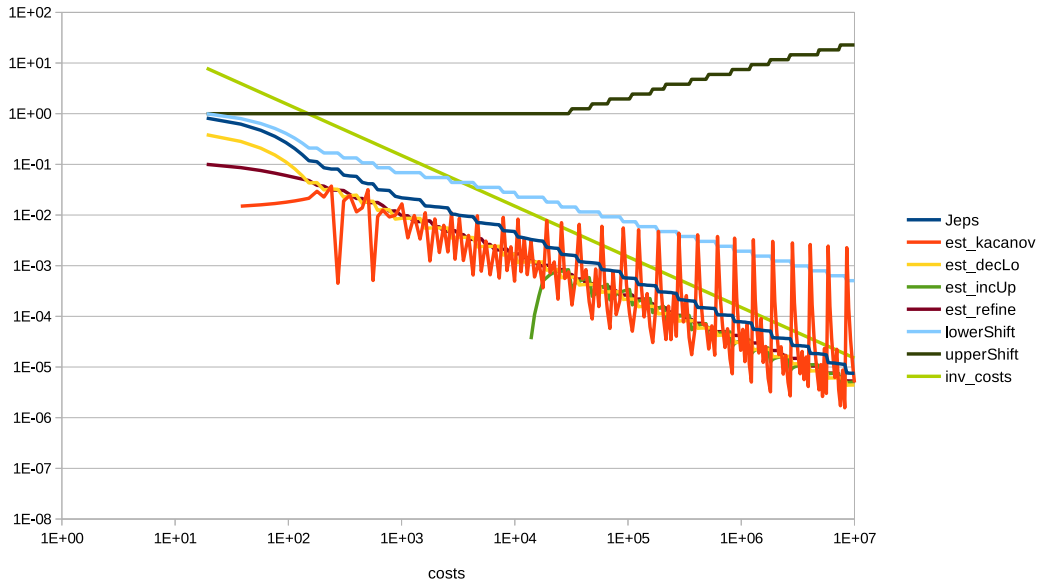


Figure 5.11: Performance of the algorithm for constant force and  $p = \frac{8}{7}$ .

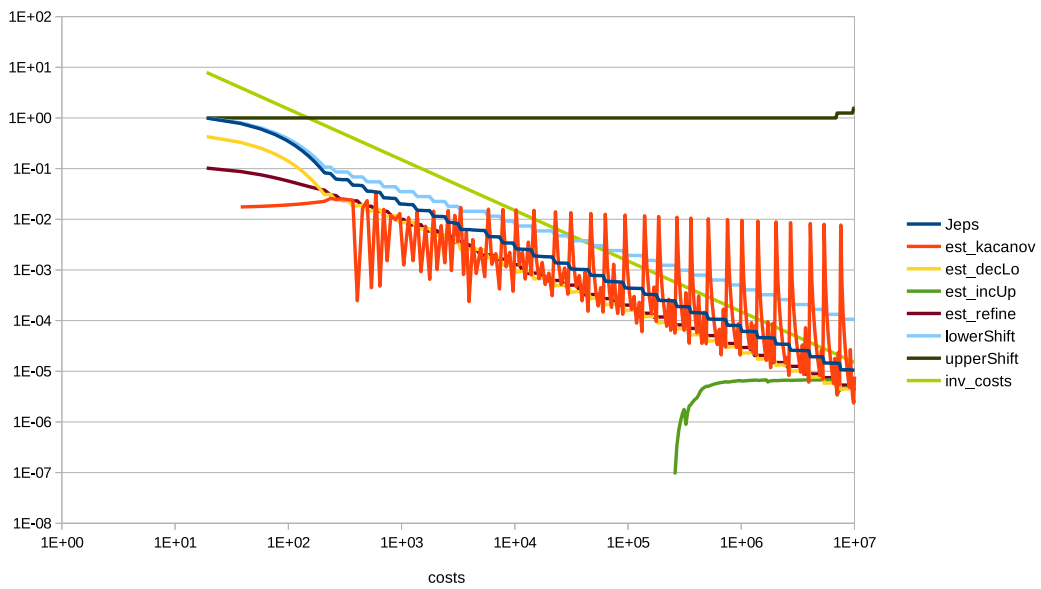


Figure 5.12: Performance of the algorithm for constant force and  $p = \frac{16}{15}$ .

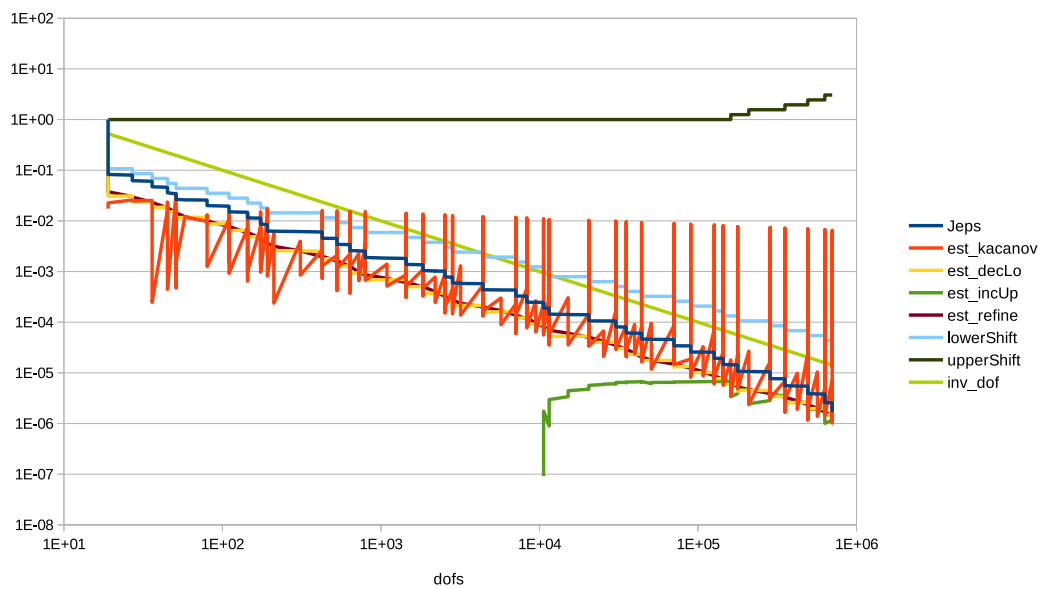


Figure 5.13: Performance of the algorithm for constant force and  $p = \frac{16}{15}$ . In this plot, the  $x$ -axis is marked with dofs instead of costs.



### 5.3.5 Short Summary of the Experiments

The example “bump” shows that for nice data – that is  $\Omega$  convex,  $u \in C^\infty(\overline{\Omega})$  and  $u = 0$  on  $\partial\Omega$  – the fully adaptive Kačanov iteration yields linear convergence of  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u)$  on the double logarithmic scale.

The effect of  $|\nabla u| \notin L^\infty(\Omega)$  can be studied in the example “needle”. Also in this experiment we observe optimal rates.

In the last example “constant force” on the L-shaped domain two problems arise. Firstly, (although the energy difference still admits linear decay on the double logarithmic scale) the estimate  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u) \lesssim \text{costs}^{-1}$  does not seem to hold true anymore. However, the estimate  $\mathcal{J}_\varepsilon(v_n) - \mathcal{J}(u) \lesssim \text{dofs}^{-1}$  seems to hold. This indicates that the number of non-refinement steps between two refinements grows. Secondly, the Kačanov estimator always jumps to a large value when  $\varepsilon_-$  is decreased. It seems that in both cases the Kačanov estimator overestimates. Therefore, a future goal is to find a more suitable Kačanov error estimator.



# Bibliography

- [AF03] Robert A. Adams and John J. F. Fournier, *Sobolev spaces*, vol. 140, Academic press, Amsterdam Boston, 2003.
- [Arn07] Anton Arnold, *Variationsrechnung*, <http://www.asc.tuwien.ac.at/~arnold/lehre/variationsrechnung/var-rechn.pdf>, 2007.
- [BDK12] Liudmila Belenki, Lars Diening, and Christian Kreuzer, *Optimality of an adaptive finite element method for the  $p$ -Laplacian equation*, IMA Journal of Numerical Analysis **32** (2012), no. 2, 484–510.
- [Cer04] Johann Cervenka, *Three-dimensional mesh generation for device and process simulation*, <http://www.iue.tuwien.ac.at/phd/cervenka/node14.html>, 2004.
- [CM10] Andrea Cianchi and Vladimir G. Maz'ya, *Global Lipschitz Regularity for a Class of Quasilinear Elliptic Equations*, Communications in Partial Differential Equations **36** (2010), no. 1, 100–133.
- [CUMP04] David Cruz-Urbe, José Martell, and Carlos Pérez, *Extrapolation from  $A_\infty$  weights and applications*, J. Funct. Anal. **213** (2004), no. 2, 412–439.
- [DE08] Lars Diening and Frank Ettwein, *Fractional estimates for non-differentiable elliptic systems with general growth*, Forum Mathematicum **20** (2008), no. 3, 523–556.
- [DK08] Lars Diening and Christian Kreuzer, *Linear convergence of an adaptive finite element method for the  $p$ -Laplacian equation*, SIAM Journal on Numerical Analysis **46** (2008), no. 2, 614–638.
- [DKS13] Lars Diening, Christian Kreuzer, and Endre Süli, *Finite Element Approximation of Steady Flows of Incompressible Fluids with Implicit Power-Law-Like Rheology*, SIAM Journal on Numerical Analysis **51** (2013), no. 2, 984–1015.

- [dM93] Gianni dal Maso, *An Introduction to  $\Gamma$ -convergence*, Birkhäuser, Boston, MA, 1993.
- [Dör96] Willy Dörfler, *A convergent adaptive algorithm for poisson's equation*, SIAM Journal on Numerical Analysis **33** (1996), no. 3, 1106–1124.
- [DRS10] Lars Diening, Michael Růžička, and Katrin Schumacher, *A decomposition technique for John domains*, Ann. Acad. Sci. Fenn. Math **35** (2010), no. 1, 87–114.
- [Dzi10] Gerhard Dziuk, *Theorie und Numerik partieller Differentialgleichungen*, De Gruyter, Berlin New York, 2010.
- [Ebm02] Carsten Ebmeyer, *Mixed Boundary Value Problems for Nonlinear Elliptic Systems with  $p$ -Structure in Polyhedral Domains*, Mathematische Nachrichten **236** (2002), no. 1, 91–108.
- [EEK06] David E. Edmunds, W. Desmond Evans, and Georgi E. Karadzhov, *Sharp estimates of the embedding constants for Besov spaces*, Rev. Mat. Complut. **19** (2006), no. 1, 161–182. MR 2219827 (2007a:46027)
- [FK97] Alberto Fiorenza and Miroslav Krbeč, *Indices of Orlicz spaces and some applications*, Commentationes Mathematicae Universitatis Carolinae **38** (1997), no. 3, 433–452.
- [Gra08] Loukas Grafakos, *Classical fourier analysis*, vol. 2, Springer, 2008.
- [HJS97] Weimin Han, Søren Jensen, and Igor Shimansky, *The kačanov method for some nonlinear problems*, Applied Numerical Mathematics **24** (1997), no. 1, 57–79.
- [KOF77] Alois Kufner, John Oldrich, and Svatopluk Fučík, *Function spaces*, Noordhoff International Pub, Leyden, 1977.
- [KR61] Mark A. Krasnosel'skiĭ and Yakov B. Rutickii, *Convex functions and Orlicz spaces*, Cambridge Univ Press, 1961.
- [Lin06] Peter Lindqvist, *Notes on the  $p$ -laplace equation*, <http://www.math.ntnu.no/~lqvist/p-laplace.pdf>, 2006.
- [MS64] Norman Meyers and James Serrin,  *$H = W$* , Proc. Nat. Acad. Sci USA **51** (1964), 1055–1056.
- [Pee66] Jaak Peetre, *Espaces d'interpolation et théorème de Soboleff*, Ann. Inst. Fourier (Grenoble) **16** (1966), no. fasc. 1, 279–317. MR 0221282 (36 #4334)

- [RD07] Michael Růžička and Lars Diening, *Non-Newtonian fluids and function spaces*, NAFSA 8—Nonlinear analysis, function spaces and applications. Vol. 8, Czech. Acad. Sci., Prague, 2007, pp. 94–143.
- [Růž06] Michael Růžička, *Nichtlineare Funktionalanalysis: Eine Einführung*, Springer-Verlag, 2006.
- [Sim64] Igor Borisovich Simonenko, *Interpolation and extrapolation of linear operators in Orlicz spaces*, *Matematicheskii Sbornik* **105** (1964), no. 4, 536–553.
- [SM93] Elias Stein and Timothy Murphy, *Harmonic analysis: real-variable methods, orthogonality, and oscillatory integrals*, Princeton mathematical series, Princeton, N.J. Princeton University Press, 1993.
- [Tri78] Hans Triebel, *Interpolation theory, function spaces, differential operators*, North-Holland Pub. Co, Amsterdam New York, 1978.
- [Wan13] Maximilian Wank, *A Caccioppoli Inequality for Energy Minimising Maps*, Master thesis, Ludwig-Maximilians-Universität München, September 2013.
- [Wei95] Juraj Weisz, *A posteriori error estimate of approximate solutions to a special nonlinear elliptic boundary value problem*, *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik* **75** (1995), no. 1, 79–81.



# Acknowledgements

Firstly, I would like to thank Prof. Dr. Lars Diening. His outstanding supervision arose in many facets, such as his ability to explain complex mathematics in a surprisingly simple way, his patience when doing so and that he always held a protecting hand over me – no matter whether on a professional or a bureaucratic level. Furthermore, his supervision was very sustained as our path was the same for a half of a decade.

Moreover, I would also like to extend thanks to my entire work group. First of all to Roland Tomasi, who impressed me a lot with his intuition for mathematics, experience in programming and knowledge in physics. I am very glad, that besides his expertise he became a good friend of mine. Of course it was not only him, but the hole group of applied analysis containing Prof. Dr. Stefan Kunis, Ulrich von der Ohe and Dominik Nagel as well as over a long time Dr. Thomas Peter, Dr. Ines Melzer and Toni Scharle. Additionally, I want to thank the members of the former group in Munich, consisting of Prof. Dr. Dominic Breit, Dr. Sebastian Schwarzacher, Dr. Parth Soneji and Franz-Xaver Gmeineder, for constructing the basis of my interest in numerical analysis.

Luckily, I was not only supported by my academic environment. I want to express my deepest gratitude to my dearly beloved wife Michaela Wank. As the person knowing me the best she always found the right words to motivate me and to put gentle pressure into the right directions. Beyond that, I can not thank her enough for her loyalty during the long time I spent in Osnabrück.

Finally, I would like to thank my family, my fellow students and my friends, who are too many to list. However, I want to point out my flatmates Stefan Niemeier, Jens Telljohann and Lukas Pälme for their uncomplicated way of giving me a home and a place for distraction.