



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

# Optimization of Complex Systems

Lecture Notes in WiSe2023

5. Oktober 2022

Winnifried Wollner

---

[winnifried.wollner@uni-hamburg.de](mailto:winnifried.wollner@uni-hamburg.de)

Universität Hamburg,  
MIN Fakultät,  
Fachbereich Mathematik



# Preface

These notes accompany the lecture “Optimization of Complex Systems” in Winter 2022/2023. They have advanced from lecture notes from my lectures on “Optimization in Function Spaces” in Summer 2018 and Winter 2019/2020 at TU Darmstadt. and are based on lecture notes “Optimization in Function Spaces” by Prof. Dr. Stefan Ulbrich, TU Darmstadt and “Optimization in Banach Spaces” by Prof. Dr. Michael Ulbrich, TU München.

They are made available to participants for their private educational use. Please do not distribute them any further.

Be advised that these notes can not replace your own notes taken during the lecture, because neither will these notes contain any sketches nor will all comments from the lecture be included. In addition, these notes may differ from the lecture at some times.

If you wish to broaden some of the topics of this lecture here are some helpful references (although there are many more)

**Calculus of Variations** [Dacorogna \[2008\]](#)

**Sobolev-spaces** [Adams and Fournier \[2003\]](#)

**Functional analysis (and function spaces)** [Alt \[2002\]](#), [Ciarlet \[2013\]](#), [Rudin \[1991\]](#), [Werner \[2005\]](#)

**Optimization with PDE constraints** [De los Reyes \[2015\]](#), [Hinze et al. \[2009\]](#), [Ito and Kunisch \[2008\]](#), [Tröltzsch \[2005\]](#).

It is expected by the participants to have some background in functional analysis, the main results needed in the lecture are summarized in Appendix [A](#) if you need a reminder.

You may send comments and corrections to

[winnifried.wollner@uni-hamburg.de](mailto:winnifried.wollner@uni-hamburg.de)

I will gratefully incorporate these into the next version of the script

Winnifried Wollner

Hamburg, October 2022.



# Inhaltsverzeichnis

<b>1. Introduction</b>	<b>1</b>
1.1. Introduction	1
1.2. Why Should we Consider the Function Space Case?	1
1.2.1. Lavrentiev Phenomenon	2
<b>2. Existence of Solutions</b>	<b>5</b>
2.1. Existence by Compactness	6
2.2. Existence by Weak-Convergence	7
<b>3. Optimality Conditions</b>	<b>13</b>
3.1. Differentiability	13
3.2. First Order Necessary Optimality Conditions	16
3.3. The Karush-Kuhn-Tucker conditions	23
3.3.1. Cone Constraints	28
3.3.2. Slater Condition	29
3.4. Sufficient optimality conditions	34
3.4.1. The convex case	34
3.4.2. Second order sufficient conditions	34
<b>4. Adjoint Approach</b>	<b>39</b>
4.1. Adjoint Approach for the First Derivative	39
4.2. Adjoint Formulas for the Second Derivatives	42
<b>5. Methods</b>	<b>45</b>
5.1. Unconstrained Case	45
5.1.1. Gradient descent	45
5.1.2. Newton Methods	50
5.2. Constrained Case	53
5.2.1. Reformulation of the KKT-conditions as a nonsmooth equation	53
5.2.2. Generalized Newton Methods	55
5.2.3. Semi-smoothness	59
5.2.4. Semi-Smoothness in Finite Dimensions	60
5.2.5. Semi-Smoothness of Nemyzkii-Operators	61
5.2.6. Application	64
5.2.7. Control Constrained Optimal Control	65

<b>6. Discretization</b>	<b>69</b>
6.1. Linear Finite Elements for Elliptic Problems	70
6.1.1. A Priori Error Estimates	71
6.2. Discretization of the Model Problem	72
6.2.1. Variational Discretization	73
<b>A. Functional Analytic Background</b>	<b>77</b>
A.1. Normed Linear Spaces	77
A.2. Convexity	78
A.3. Linear Operators	80
A.4. Adjoints	83
A.5. Weak Convergence	84
A.6. Lebesgue and Sobolev spaces	86
A.6.1. Domain regularity	86
A.6.2. Lebesgue spaces	86
A.6.3. Sobolev spaces	87
A.6.4. Embedding theorems	89

# 1. Introduction

## 1.1. Introduction

In this lecture, we are interested in problems of the type

$$f(x) \rightarrow \min \quad \text{s.t. } x \in X^{\text{ad}} \quad (1.1)$$

where  $f : X \rightarrow \mathbb{R}$  is a functional acting on a Banach space  $X$  and  $X^{\text{ad}} \subset X$  denotes the admissible set.

In many cases the unknown variable denotes a function and the corresponding Banach space encodes the properties of the considered functions. There are many examples where such problems occur, to give only a few of them:

- Modelling of elastic materials, damage, fracture
- Optimal design (shape and topology) of buildings, ships, aircrafts, cars
- Optimal control of robot movements, space craft trajectories, heating and cooling processes (wine fermentation, steel hardening, ...)
- Inverse problems, e.g., seismic inversion, data assimilation in weather and climate models, tomography

## 1.2. Why Should we Consider the Function Space Case?

We will highlight this with the very basic question of existence of solutions. The following is well known from 'Calculus 2'

**Theorem 1.2.1.** *Let  $f : X^{\text{ad}} \subset \mathbb{R}^n \rightarrow \mathbb{R}$  be continuous on the bounded closed set  $X^{\text{ad}}$ . Then (1.1) has at least one solution.*

*Beweis.* We have that  $Q^{\text{ad}}$  is compact by the theorem of Bolzano-Weierstrass. Hence the continuous function  $j$  attains its minimum on  $Q^{\text{ad}}$  due to the extreme value theorem.  $\square$

Unfortunately, this is no longer true in infinite dimensions, as is shown by the following example

## 1. Introduction

**Example 1.2.2.** (Counterexample) Consider the problem

$$\begin{aligned} \min_{x \in C[-1,1]} f(x) &:= \int_{-1}^1 (x(s) - x^d(s))^2 ds = \|x - x^d\|_{L^2(-1,1)}^2 \\ \text{s.t. } -1 &\leq x(s) \leq 1 \quad \forall s \in [-1, 1]. \end{aligned}$$

with

$$x^d(s) = \begin{cases} -1 & s < 0, \\ 1 & s \geq 0. \end{cases}$$

Clearly  $x \mapsto f(x) \geq 0$  is continuous and the set  $\{x \in C[-1, 1] \mid -1 \leq x \leq 1\}$  is bounded and closed. But the sequence

$$x_k(s) = \begin{cases} -1 & s < -1/k, \\ ks & -1/k \leq s \leq 1/k, \\ 1 & s \geq 1/k, \end{cases}$$

satisfies

$$f(x_k) = \int_{-1/k}^{1/k} (1 - k|s|)^2 ds \leq \int_{-1/k}^{1/k} ds \rightarrow 0 \quad (k \rightarrow \infty).$$

Hence it holds

$$\inf_{x \in C[-1,1]} f(x) = 0$$

but  $f(\bar{x}) = 0$  if and only if  $\bar{x} = x^d \notin C[-1, 1]$ .

In particular, we will need some additional tools for the analysis of such problems. We can already see that if instead of functions  $x \in C[-1, 1]$  we would have taken functions in  $L^2[-1, 1]$  (functions that are square integrable) the counterexample would have failed.

Hence, we may expect that the right choice of the space where minimizers are searched for is crucial but also that we need an improved theory for the analysis.

### 1.2.1. Lavrentiev Phenomenon

It should be noted that, in general, not only the existence of a minimizer is depending on the choice of the spaces. Even the infimum values of the functional may depend on the chosen spaces. This is obvious, of course, since on any one-dimensional subspace not containing a minimizing sequence for the functional this is true. But infact, such defects can be a lot more subtle as it is demonstrated by the so called Lavrentiev phenomenon, [Lavrentieff \[1927\]](#) which stated that certain minimal functional values can not be approximated by Lipschitz continuous functions.



## 1.2. Why Should we Consider the Function Space Case?

We illustrate this with the following example. Consider the functional

$$J(u) = \int_0^1 (u(t)^3 - t)^2 u'(t)^6 dt$$

then it holds (with the boundary conditions  $u(0) = 0, u(1) = 1$ )

$$\min_{u \in C^1(0,1)} J(u) = 0 < \inf_{u \in C^{0,1}(0,1)} J(u).$$

To see this, we note that  $\sqrt[3]{t} \in C^1(0,1)$  but not in  $C^{0,1}(0,1)$  and  $J(\sqrt[3]{\cdot}) = 0$ . Now, let  $u \in C^{0,1}(0,1)$  be arbitrary and let  $v = \sqrt[3]{t}/2$ . Then there exists some  $t_0 \in (0,1)$  such that

$$u(t) \leq v(t) \quad \forall t \in [0, t_0]$$

and  $u(t_0) = v(t_0)$  due to the regularity (and boundary values) of  $u$  and  $v$ . **draw sketch** Hence for  $t \in [0, t_0]$  and  $\xi \in \mathbb{R}$  it holds  $(|u^3 - t| \geq |v^3 - t|, \text{ since } u^3 \leq v^3 \leq t/8 \leq t)$

$$(u^3(t) - t)^2 \xi^6 \geq (v^3(t) - t)^2 \xi^6 = \left( \frac{\sqrt[3]{t^3}}{2^3} - t \right)^2 \xi^6 = \frac{7^2}{8^2} t^2 \xi^6.$$

Further, by fundamental theorem of calculus and Hölder inequality, we get

$$\begin{aligned} \frac{\sqrt[3]{t_0}}{2} &= v(t_0) = u(t_0) = \int_0^{t_0} u'(t) dt = \int_0^{t_0} t^{-1/3} (t^{1/3} u'(t)) dt \\ &\leq \left( \int_0^{t_0} t^{-2/5} dt \right)^{5/6} \left( \int_0^{t_0} t^2 u'(t)^6 dt \right)^{1/6} \\ &\leq \left( \frac{5}{3} t^{3/5} \Big|_0^{t_0} \right)^{5/6} \left( \int_0^{t_0} t^2 u'(t)^6 dt \right)^{1/6} \\ &= \left( \frac{5}{3} \right)^{5/6} t_0^{1/2} \left( \int_0^{t_0} t^2 u'(t)^6 dt \right)^{1/6}. \end{aligned}$$

Hence we obtain

$$\int_0^{t_0} t^2 u'(t)^6 dt \geq \left( \frac{t_0^{1/3}}{2} t_0^{-1/2} \left( \frac{3}{5} \right)^{5/6} \right)^6 = \frac{1}{2^6 t_0} \left( \frac{3}{5} \right)^5.$$

Combining this we conclude that for any  $u \in C^{0,1}(0,1)$  it holds

$$\begin{aligned} J(u) &\geq \int_0^{t_0} (u(t)^3 - t)^2 u'(t)^6 dt \\ &\geq \int_0^{t_0} (v(t)^3 - t)^2 u'(t)^6 dt \\ &\geq \frac{7^2}{8^2} \int_0^{t_0} t^2 u'(t)^6 dt \\ &\geq \frac{7^2}{8^2} \frac{3^5}{5^5} \frac{1}{2^6 t_0} \geq \frac{7^2}{8^2} \frac{3^5}{5^5} \frac{1}{2^6} > 0. \end{aligned}$$

Which shows the claim.

## *1. Introduction*

## 2. Existence of Solutions

We will now discuss how to obtain the existence of solutions to minimization problems in Banach spaces.

There are essentially two possibilities to cure the defect seen in Example 1.2.2. One is by adding assumptions on the feasible set while the other adds assumptions for the functions.

In both cases we assume that our problem is given in the form

$$\min f(x) \quad G(x) \in K \quad (\text{P})$$

where

### Assumption 2.0.1.

1.  $G: X \rightarrow Z$  is a continuous operator between two (real) Banach spaces,
2.  $f: X \rightarrow \mathbb{R}$  is lower semicontinuous.
3.  $K \subset Z$  is closed and convex,

as a reminder:

**Definition 2.0.2.** A function  $f: X \rightarrow \mathbb{R}$  is lower semicontinuous, if for any convergent sequence  $x^k \in X$  it holds

$$\liminf_{k \rightarrow \infty} f(x^k) \geq f(\lim_{k \rightarrow \infty} x^k),$$

or, as an equivalent geometric condition, if for any  $\alpha \in \mathbb{R}$  the sets

$$\{x \in X \mid f(x) \leq \alpha\}$$

are (sequentially) closed.

It is clear by this, that the admissible set

$$X^{\text{ad}} = \{x \in X \mid G(x) \in K\}$$

is closed since  $G$  is continuous and  $K$  is closed.

## 2. Existence of Solutions

The notation in (P) contains all the usual conditions, e.g.,

$$h(x) = 0, \quad g(x) \in C, \quad x \in M$$

for continuous operators  $h: X \rightarrow Z_1$ ,  $g: X \rightarrow Z_2$  for Banach spaces  $Z_1, Z_2$ , a closed convex cone  $C \in Z_2$ , and a closed convex set  $M \subset X$  can be modeled in this setting by choosing  $Z = Z_1 \times Z_2 \times X$  with  $K = \{0\} \times C \times M$  and  $G(x) = (h(x), g(x), x)^T$ .

The constraint  $F(x) \in C$  for a closed convex cone is an abstract inequality constraint, as the usual component-wise inequality  $g(x) \leq 0$  in  $\mathbb{R}^m$  from nonlinear optimization can also be written as the inclusion

$$g(x) \in C = (-\infty, 0]^m.$$

### 2.1. Existence by Compactness

Now, we show the first alternative of generalizing the existence Theorem 1.2.1.

**Theorem 2.1.1.** *Let Assumption 2.0.1 be satisfied and assume that there is a point  $x^0 \in X^{\text{ad}}$  such that the level-set*

$$\mathcal{N}(x^0) = \{x \in X^{\text{ad}} \mid f(x) \leq f(x^0)\}$$

*is compact. Then Problem (P) has at least one global solution.*

*Beweis.* Let  $x^k \in \mathcal{N}(x^0)$  be a minimizing sequence, i.e.,

$$\lim_{k \rightarrow \infty} f(x^k) = \inf_{x \in X^{\text{ad}}} f(x) =: f^*.$$

By compactness of  $\mathcal{N}(x^0) \subset X^{\text{ad}}$ , there exists a convergent subsequence  $x^k \rightarrow \bar{x}$ . Lower semicontinuity asserts

$$f^* = \lim_{k \rightarrow \infty} f(x^k) \geq f(\lim_{k \rightarrow \infty} x^k) = f(\bar{x}) \geq f^*$$

and thus  $\bar{x}$  is a global minimizer. □

**Remark 2.1.2.** Notice, that level-set  $\mathcal{N}(x^0)$  is always closed for a lower semicontinuous function, so the essential assumption of Theorem 2.1.1 is that in addition the set  $\mathcal{N}(x^0)$  is precompact.

## 2.2. Existence by Weak-Convergence

While this appears to be very natural (we only use the Bolzano-Weierstrass Theorem in finite-dimensions to see the equivalence to bounded and closed), the characterization of compact sets in infinite dimensional spaces is a bit more cumbersome. For instance on a bounded domain (open and connected)  $\Omega \subset \mathbb{R}^d$  a subset  $M \subset C(\overline{\Omega})$  is relatively compact ( $\text{cl } M$  is compact) if and only if the functions in  $M$  are uniformly bounded, i.e.,

$$\sup_{f \in M} \|f\|_{\infty} < \infty$$

and equicontinuous, i.e., for any  $\varepsilon > 0$  there is some  $\delta$  such that

$$\sup_{f \in M} |f(x) - f(y)| \leq \varepsilon, \quad \forall \|x - y\| \leq \delta$$

as it is asserted by the Arzelà-Ascoli theorem.

## 2.2. Existence by Weak-Convergence

To avoid checking for compactness (which is often not given) an alternative lies in enforcing existence of minimizers by increasing the conditions on the operators  $f$  and  $G$ .

To this end, we recall from functional analysis

**Definition 2.2.1.** We say that a sequence  $x^k \in X$  converges weakly to  $x \in X$  (written as  $x^k \rightharpoonup x$ ) if

$$\langle x^*, x^k \rangle \rightarrow \langle x^*, x \rangle \quad \forall x^* \in X^*$$

where  $X^* = \mathcal{L}(X; \mathbb{R})$  denotes the dual space to  $X$ .

Before we continue, we collect a few useful facts:

**Theorem 2.2.2.** Let  $X$  be a real Banach space, then it holds:

1. any closed (norm-topology) and convex set is weak sequentially closed (Banach-Saks-Mazur theorem),
2. if  $X$  is reflexive, then any bounded sequence  $x^k \in X$  contains a weakly convergent subsequence (Banach-Eberlein-Šmulian theorem),
3. a weakly convergent sequence is bounded,
4. if  $f : X \rightarrow \mathbb{R}$  is convex and lower semicontinuous, then  $f$  is weakly lower semiconti-

## 2. Existence of Solutions

nuous (w.l.s.c), i.e., for any  $x^k \rightharpoonup x$  it holds

$$\liminf_{k \rightarrow \infty} f(x^k) \geq f(x).$$

In particular, it is 'easy' to check whether a sequence contains a weakly convergent subsequence!

With this we can now formulate our second existence result

**Theorem 2.2.3.** Assume that

1.  $G: X \rightarrow Z$  is weakly sequentially continuous on the **reflexive** Banach space  $X$  to the Banach space  $Z$ , i.e.,

$$x_k \rightharpoonup x \quad \text{then} \quad G(x_k) \rightharpoonup G(x),$$

2.  $f: X \rightarrow \mathbb{R}$  is weakly lower semicontinuous,
3.  $K \subset Z$  is closed and convex,
4. there is a point  $x^0 \in X^{\text{ad}}$  such that the level-set

$$\mathcal{N}(x^0) = \{x \in X^{\text{ad}} \mid f(x) \leq f(x^0)\}$$

is bounded.

Then Problem (P) has at least one global solution.

*Beweis.* As in the proof of Theorem 2.1.1, we select a minimizing sequence  $x^k \in \mathcal{N}(x^0)$  with

$$\lim_{k \rightarrow \infty} f(x^k) = \inf_{x \in X^{\text{ad}}} f(x) = f^*.$$

By boundedness of  $\mathcal{N}(x^0)$  and Theorem 2.2.2-2. this sequence contains a subsequence with weak limit  $\bar{x}$ . By the assumed weak continuity  $G(x^k) \rightharpoonup G(\bar{x})$  and since  $K$  is closed and convex Theorem 2.2.2-1. asserts  $\bar{x} \in X^{\text{ad}}$ .

Now,  $f$  is w.l.s.c and thus

$$f^* = \lim_{k \rightarrow \infty} f(x^k) \geq f(\bar{x}) \geq f^*$$

and thus  $\bar{x}$  is a global minimizer. □

This leaves us with the remaining problem to see which functions are weak continuous or lower semicontinuous. Since every convergent sequence is also weakly convergent this is a stronger requirement than continuity. One possibility is utilizing Theorem 2.2.2-4.

Further, in many situations it is helpful to utilize compactness properties of mappings.

**Definition 2.2.4.** A continuous linear operator  $A: X \rightarrow Y$  between two Banach spaces is called compact, if for any bounded sequence  $x^k \in X$  the sequence  $Ax^k \in Y$  contains a convergent subsequence.

In particular, a Banach space  $X$  is compactly embedded in a Banach space  $X_0$  if the embedding  $\text{id}: X \rightarrow X \subset X_0$  is compact.

**Lemma 2.2.5.** If a linear operator  $A: X \rightarrow Y$  is compact, and  $x^k \rightharpoonup x$  in  $X$  then  $Ax^k \rightarrow Ax$  in  $Y$ .

*Beweis.* Exercise! □

To show how these techniques can be applied, we consider the following example.

**Example 2.2.6.** Consider the problem

$$\begin{aligned} \min \quad & J(q, u) = \frac{1}{2} \|u - u^d\|_{L^2(\Omega)}^2 + \frac{\alpha}{2} \|q\|_{L^2(\partial\Omega)}^2 \\ \text{s.t.} \quad & -\Delta u + u^3 = 0 \quad \text{on } \Omega, \\ & \partial_n u + u = q \quad \text{on } \partial\Omega, \\ & a \leq q \leq b \quad \text{on } \partial\Omega, \end{aligned} \tag{2.1}$$

where  $\Omega \subset \mathbb{R}^d$ ,  $d = 2, 3$  a Lipschitz-domain and  $a, b \in L^2(\partial\Omega)$  with w.l.o.g.  $a \leq 0 \leq b$ .

Then we can define  $Q = L^2(\partial\Omega)$ ,  $U = H^1(\Omega)$ , and

$$Q^{\text{ad}} = \{q \in Q \mid a \leq q \leq b\}.$$

For the equality constraint, we consider the weak form to define

$$\begin{aligned} E: (q, u) \in Q \times U &\rightarrow U^* \\ \langle E(q, u), v \rangle &:= (\nabla u, \nabla v)_{L^2(\Omega)^n} + (u^3, v)_{L^2(\Omega)} + (u - q, v)_{L^2(\partial\Omega)}. \end{aligned}$$

Although it can be shown, that for any given  $q \in Q$  the equation  $E(q, u) = 0$  admits a unique solution, we will not use this fact. Instead, we just notice, that the embedding  $H^1(\Omega) \hookrightarrow L^6(\Omega)$  is continuous (Sobolev-embedding theorem) and that  $H^1(\Omega) \hookrightarrow L^5(\Omega)$  is compact (Rellich-Kondrachov) if  $d = 2, 3$ .

We can now start to check the prerequisites of the existence Theorem 2.2.3. To this end, we fix  $X = Q \times U$ ,  $Z = Q \times U^*$ ,  $K = Q^{\text{ad}} \times \{0\}$ ,  $G(x) = G(q, u) = (q, E(q, u))^T$ ,  $f(x) = f(q, u) = J(q, u)$ .

## 2. Existence of Solutions

Now, we check the prerequisites.

1. First, we check the weak continuity of  $G$ . To this end, let  $x^k = (q^k, u^k) \rightharpoonup x = (q, u)$  be a weak convergent sequence. for the first component of  $G$  nothing is to be shown. For the second component, we notice, that all continuous linear terms immediately converge by definition as, e.g.,

$$(\nabla u^k, \nabla v) \rightarrow (\nabla u, \nabla v)$$

by definition. For the nonlinear term, we recall that  $H^1(\Omega) \subset L^5(\Omega)$  is compact, and hence  $u^k \rightarrow u$  (strongly) in  $L^5(\Omega)$ . As a consequence  $(u^k)^3 \rightarrow u^3$  in  $L^{5/3}(\Omega)$ , because the inequality  $(a, b \in \mathbb{R})$

$$|a^3 - b^3| = |a^2 + ab + b^2||a - b| \leq \frac{3}{2}|a^2 + b^2||a - b|$$

shows via Hölder's inequality with  $\frac{2}{5} + \frac{1}{5} = \frac{3}{5}$

$$\|(u^k)^3 - u^3\|_{L^{5/3}} \leq \frac{3}{2} \left( \|(u^k)^2\|_{L^{5/2}} + \|u^2\|_{L^{5/2}} \right) \|u^k - u\|_{L^5} \rightarrow 0.$$

Now,  $U \subset L^{5/2}(\Omega) = L^{5/3}(\Omega)^*$  and thus, for any  $v \in U = H^1(\Omega)$ ,

$$((u^k)^3 - u^3, v)_{L^2(\Omega)} \leq \|(u^k)^3 - u^3\|_{L^{5/3}} \|v\|_{H^1(\Omega)} \rightarrow 0$$

which shows

$$G(x^k) \rightharpoonup G(x).$$

Further,  $X$  is a Hilbert-space and thus reflexive.

2. To see weak lower semi-continuity of  $f$ , we refer to Theorem 2.2.2-4.
3. Clearly  $K$  is closed and convex.
4. To see, the last point, let  $x^0 = (0, 0)$  be given. Clearly  $x^0 \in X^{\text{ad}}$ . Now, to see boundedness, we notice that for  $\|q^k\|_{L^2(\partial\Omega)} \rightarrow \infty$  it is

$$f(x^k) \geq \frac{\alpha}{2} \|q^k\|_{L^2(\partial\Omega)}^2 \rightarrow \infty$$

and hence the  $q$ -component remains bounded. (This is alternatively clear from the bounds  $a \leq q \leq b$ )



## 2.2. Existence by Weak-Convergence

For the  $u$ -component assume, that  $(q, u) \in Q \times U$  solves  $E(q, u) = 0$ . Then testing the weak equation with  $v = u$  one obtains

$$\begin{aligned} \|\nabla u\|_{L^2(\Omega)^n}^2 + \|u\|_{L^4(\Omega)}^4 + \|u\|_{L^2(\partial\Omega)}^2 &= (q, u)_{L^2(\partial\Omega)} \\ &\leq \|q\|_{L^2(\partial\Omega)} \|u\|_{L^2(\partial\Omega)} \\ &\leq \frac{1}{2} \|q\|_{L^2(\partial\Omega)}^2 + \frac{1}{2} \|u\|_{L^2(\partial\Omega)}^2 \end{aligned}$$

and consequently

$$\|\nabla u\|_{L^2(\Omega)^n}^2 + \frac{1}{2} \|u\|_{L^2(\partial\Omega)}^2 \leq \frac{1}{2} \|q\|_{L^2(\partial\Omega)}^2.$$

Using a generalized Poincaré-Friedrichs inequality, i.e.,

$$\|u\|_{L^2(\Omega)}^2 \leq c_\Omega \left( \|\nabla u\|_{L^2(\Omega)^n}^2 + \frac{1}{2} \|u\|_{L^2(\partial\Omega)}^2 \right),$$

shows

$$\|u\|_{H^1(\Omega)}^2 \leq c_\Omega \|q\|_{L^2(\partial\Omega)}^2$$

and thus  $\mathcal{N}(0)$  is bounded.

Consequently, Theorem 2.2.3 asserts the existence of at least one solution.

## *2. Existence of Solutions*

## 3. Optimality Conditions

Analogously to the finite dimensional case, we will try to characterize local minimizers by conditions on the derivatives of the involved functionals. Hence we start with a recollection on differentiability in Banach-spaces.

### 3.1. Differentiability

Let  $X, Y$  be two Banach spaces, and  $F: X \rightarrow Y$  be some mapping.

**Definition 3.1.1.** For given  $x \in X$  and  $\delta x \in X$  assume that the limit

$$F'(x; \delta x) = \lim_{t \downarrow 0} \frac{F(x + t\delta x) - F(x)}{t}$$

exists, then  $F'(x; \delta x)$  is called the *directional derivative* of  $F$  at  $x$  in direction  $\delta x$ .

Note that the directional derivative is not necessarily linear, e.g., consider  $F: \mathbb{R} \rightarrow \mathbb{R}$ ,  $F(x) = |x|$  which is directionally differentiable in  $x = 0$  but  $F'(0; 1) = 1 \neq F'(0; -1)$ .

**Definition 3.1.2.** Assume that  $F'(x; \delta x)$  exists for all  $\delta x \in X$  and that there exists an operator  $A \in \mathcal{L}(X, Y)$  such that

$$A\delta x = F'(x; \delta x) \quad \forall \delta x \in X$$

then  $F$  is called *Gâteaux-differentiable* at  $x$ . The mapping  $A = F'(x)$  is called the *Gâteaux-derivative*.

**Definition 3.1.3.**  $F$  is called *Fréchet-differentiable* at  $x$  if there exists a function  $\omega: X \rightarrow Y$  with  $\omega(\delta x) \in o(\|\delta x\|_X)$  such that

$$F(x + \delta x) = F(x) + F'(x)\delta x + \omega(\delta x).$$

In this case  $F'(x)$  is called the *Fréchet-derivative*.

### 3. Optimality Conditions

The notion  $\omega(\delta x) \in o(\|\delta x\|_X)$  means

$$\lim_{\delta x \rightarrow 0} \frac{\|\omega(\delta x)\|_Y}{\|\delta x\|_X} = 0$$

#### Example 3.1.4.

1. Consider the mapping  $A(x) = A_0x + b$  with  $A_0 \in \mathcal{L}(X, Y)$  and  $b \in Y$ . Then  $A$  is Fréchet-differentiable with  $F'(x) = A_0$  since

$$A(x + \delta x) = A_0x + A_0\delta x + b = A(x) + A_0\delta x.$$

2. Let  $H$  be a Hilbert space with scalar product  $(\cdot, \cdot)$ . Then the squared natural norm  $F(x) = \|x\|^2$  is Fréchet-differentiable with

$$F'(x) = 2x \quad \in \mathcal{L}(H, \mathbb{R}) = H^*.$$

(Exercise)

3. Consider the space  $X = C[0, 1]$  together with the mapping  $\sin: X \rightarrow X$  given by  $x \mapsto \sin(x)$ . (This is the Nemyzkii-operator generated by  $\sin: \mathbb{R} \rightarrow \mathbb{R}$ .) Then the directional derivative is given by

$$\begin{aligned} F'(x; \delta x)(s) &= \lim_{t \downarrow 0} \frac{1}{t} (\sin(x(s) + t\delta x(s)) - \sin(x(s))) \\ &= \cos(x(s))\delta x(s) \end{aligned}$$

noting that the convergence is uniform w.r.t.  $s \in [0, 1]$  due to the uniformly bounded second derivative of  $\sin$ . Now we know from Example A.3.5 that the operator given by pointwise multiplication with  $\cos(x(s))$  is in  $\mathcal{L}(X, X)$  and hence  $\sin$  is Gâteaux-differentiable with  $\sin' = \cos$ . To obtain Fréchet-differentiability we note that

$$|\sin(x(s) + \delta x(s)) - \sin(x(s)) - \cos(x(s))\delta x(s)| = \frac{1}{2} |\sin(\xi_s)| \delta x(s)^2 = |\omega(\delta x)(s)|$$

and hence

$$\frac{\sup_{s \in [0, 1]} |\sin(x(s) + \delta x(s)) - \sin(x(s)) - \cos(x(s))\delta x(s)|}{\|\delta x\|_\infty} \leq \|\delta x\|_\infty,$$

i.e., the mapping is Fréchet-differentiable.

4. We now consider  $\sin: L^p(0, 1) \rightarrow L^p(0, 1)$  for some  $1 \leq p < \infty$ . We note that the remainder in the Taylor expansion of the sin-function in  $x \equiv 0$  can be written pointwise as

$$\omega(\delta x) = \sin(0 + \delta x(s)) - \sin(0) - \cos(0)\delta x(s) = \sin(\delta x(s)) - \delta x(s).$$

Now we can define a sequence of functions

$$\delta x_\varepsilon(s) = \begin{cases} 1 & s \in [0, \varepsilon], \\ 0 & \text{otherwise,} \end{cases} \in L^p(0, 1)$$

and obtain

$$\|\omega(\delta x_\varepsilon)\|_p^p = \int_0^1 |\omega(\delta x_\varepsilon(s))|^p ds = \int_0^\varepsilon (1 - \sin(1))^p ds = (1 - \sin(1))^p \varepsilon.$$

This gives

$$\frac{\|\omega(\delta x_\varepsilon)\|_p}{\|\delta x_\varepsilon\|_p} = \frac{(1 - \sin(1))\varepsilon^{1/p}}{\varepsilon^{1/p}} = (1 - \sin(1)) > 0$$

and hence  $\sin$  is not Fréchet-differentiable between  $L^p(\Omega)$  with  $1 \leq p < \infty$ .

**Remark 3.1.5.** The last example is not a pathological case, in fact, if  $\psi: \mathbb{R} \rightarrow \mathbb{R}$  generates a Nemyzki-operator  $\Psi: W^{1,p}(\Omega) \rightarrow W^{1,p}(\Omega)$  with  $1 \leq p < \infty$ , then this operator is Fréchet-differentiable iff  $\Psi$  is continuous and **affine-linear**, see, e.g., [Appell and Zabrejko \[1990\]](#).

**Theorem 3.1.6.** Let  $X, Y, Z$  be Banach spaces and let  $\mathcal{O}_X \subset X$  and  $\mathcal{O}_Y \subset Y$  be open. Then the following is true:

1. (Sum-rule)

If  $f, g: \mathcal{O}_X \rightarrow Y$  are Gâteaux or Fréchet-differentiable then so is  $\lambda f + g$  with  $\lambda \in \mathbb{R}$  and it holds

$$(\lambda f + g)'(x) = \lambda f'(x) + g'(x).$$

2. (Mean-value theorem)

If  $f: \mathcal{O}_X \rightarrow Y$  is Gâteaux-differentiable and the interval  $I = [x_0, x_1] = \{\lambda x_0 + (1 - \lambda)x_1 \mid \lambda \in [0, 1]\}$  is contained in  $\mathcal{O}_X$ , then

$$\|f(x_0) - f(x_1)\|_Y \leq \sup_{\xi \in I} \|f'(\xi)\|_{\mathcal{L}(X,Y)} \|x_0 - x_1\|_X$$

### 3. Optimality Conditions

3. If  $f : \mathcal{O}_X \rightarrow Y$  is Gâteaux-differentiable and  $f' : \mathcal{O}_X \rightarrow \mathcal{L}(X, Y)$  is continuous, then  $f$  is Fréchet-differentiable and we say that  $f$  is continuously differentiable and write  $f \in C^1(\mathcal{O}_X, Y)$ .

4. (Chain rule)

Let  $f : \mathcal{O}_X \rightarrow Y$  and  $g : \mathcal{O}_Y \rightarrow Z$  with  $f(\mathcal{O}_X) \subset \mathcal{O}_Y$  be Fréchet-differentiable in  $x_0$  and  $f(x_0)$ , then  $g \circ f$  is Fréchet-differentiable in  $x_0$  and it holds

$$(g \circ f)'(x_0) = g'(f(x_0)) \circ f'(x_0).$$

5. (Implicit function theorem)

Let  $F : X \times Y \supset \mathcal{O}_X \times \mathcal{O}_Y \rightarrow Z$  be continuously differentiable with  $F(x_0, y_0) = 0$ . Further, assume that the derivative of  $y \mapsto F(x_0, y)$  is an isomorphism from  $Y$  to  $Z$  at  $y = y_0$ . Then there exists neighborhoods  $N(x_0)$  of  $x_0$  and  $N(y_0)$  of  $y_0$  such that for any  $x \in N(x_0)$  the equation  $F(x, y) = 0$  admits a unique solution  $y = f(x) \in N(y_0)$ . The mapping  $f : N(x_0) \rightarrow Y$  defined by this is continuously differentiable.

## 3.2. First Order Necessary Optimality Conditions

We will now derive some optimality conditions for solutions of (P) and set

$$X^{\text{ad}} = \{x \in X \mid G(x) \in K\}.$$

**Definition 3.2.1.** We say that  $\bar{x} \in X^{\text{ad}}$  is a *local solution* to the problem (P), with respect to the norm on  $X$ , if there exists  $r > 0$  such that

$$f(\bar{x}) \leq j(x) \quad \forall x \in X^{\text{ad}}, \|x - \bar{x}\|_X \leq r.$$

A solution is *global* if this holds for all  $r > 0$ .

Analogous to the finite dimensional case, we will define a cone of tangent directions as

**Definition 3.2.2.** For a set  $X^{\text{ad}} \subset X$  in a Banach space  $X$ , we define the *Bouligand-cone* (or *sequential tangent cone*, or *contingent-cone*) at  $\bar{x}$  as follows

$$\begin{aligned} T(X^{\text{ad}}, \bar{x}) &:= \{d \in X \mid d = \lim_{k \rightarrow \infty} \frac{1}{t_k}(x^k - \bar{x}), t_k \downarrow 0, x^k \in X^{\text{ad}}\} \\ &= \{d \in X \mid d = \lim_{k \rightarrow \infty} \eta_k(x^k - \bar{x}), \eta_k > 0, \bar{x} \leftarrow x^k \in X^{\text{ad}}\}. \end{aligned}$$

### 3.2. First Order Necessary Optimality Conditions

This cone is always closed (Exercise)!

With this, we have the following

**Theorem 3.2.3.** *Let  $f : O_X \subset X \rightarrow \mathbb{R}$  be defined on an open neighborhood of  $X^{\text{ad}} \subset X$ , and let  $\bar{x} \in X^{\text{ad}}$  be a local solution of (P). If  $f$  is Fréchet-differentiable in  $\bar{x}$  then*

$$f'(\bar{x})d \geq 0 \quad \forall d \in T(X^{\text{ad}}, \bar{x}).$$

*Beweis.* Let  $d \in T(X^{\text{ad}}, \bar{x})$  be given and  $t_k$  and  $x^k \in X^{\text{ad}}$  be the approximating sequences. Then it holds

$$f'(\bar{x})d = \lim_{k \rightarrow \infty} \frac{1}{t_k} f'(\bar{x})(x^k - \bar{x}) = \lim_{k \rightarrow \infty} \frac{1}{t_k} (f(x^k) - f(\bar{x}) - \omega(x^k - \bar{x})).$$

Using  $x^k \rightarrow \bar{x}$  and local optimality of  $\bar{x}$ , we get

$$f'(\bar{x})d \geq \lim_{k \rightarrow \infty} \frac{1}{t_k} \omega(x^k - \bar{x}) = \lim_{k \rightarrow \infty} \frac{\|x^k - \bar{x}\|_X}{t_k} \frac{\omega(x^k - \bar{x})}{\|x^k - \bar{x}\|_X} = \|d\|_X \lim_{k \rightarrow \infty} \frac{\omega(x^k - \bar{x})}{\|x^k - \bar{x}\|_X} = 0.$$

□

As in finite-dimensions, this characterization is not considered to be sufficient, since the set  $T(X^{\text{ad}}, \bar{x})$  is hard to describe in general. Thus we will try to replace it by simpler settings.

**Theorem 3.2.4 (The Convex Case).** *Let  $X^{\text{ad}} \subset X$  be convex, and  $f : X^{\text{ad}} \rightarrow \mathbb{R}$  be directionally differentiable. Then for any local solution  $\bar{x} \in X^{\text{ad}}$  it holds variational inequality (VI)*

$$f'(\bar{x}; x - \bar{x}) \geq 0 \quad \forall x \in X^{\text{ad}}. \quad (3.1)$$

*Further, if  $f$  is convex on  $X^{\text{ad}}$  then any  $\bar{x} \in X^{\text{ad}}$  which satisfies (3.1) is a global solution to (P).*

*Beweis.* The proof for the necessary condition is analogous to the proof of Theorem 3.2.3. For any  $x \in X^{\text{ad}}$ , we set  $d = x - \bar{x}$  and observe that for any  $t \in [0, 1]$   $\bar{x} + td = (1-t)\bar{x} + tx \in X^{\text{ad}}$  and thus

$$f'(\bar{x})(x - \bar{x}) = \lim_{t \downarrow 0} \frac{f(\bar{x} + td) - f(\bar{x})}{t} \geq 0.$$

To see the sufficiency, we note that convexity of  $f$  implies that for any  $x \in X^{\text{ad}}$  and any  $\lambda \in (0, 1)$  it holds

$$\frac{f(\bar{x} + \lambda(x - \bar{x})) - f(\bar{x})}{\lambda} \leq f(x) - f(\bar{x})$$

### 3. Optimality Conditions

and thus

$$f(x) - f(\bar{x}) \geq \delta f(\bar{x}; x - \bar{x})$$

which shows the assertion.  $\square$

**Remark 3.2.5.** The convex case is indeed a generalization of Theorem 3.2.3, since for  $\bar{x} \in X^{\text{ad}}$  it holds

$$T(X^{\text{ad}}; \bar{x}) = \text{cl}\{d \in X \mid d = \lambda(x - \bar{x}), \lambda > 0, x \in X^{\text{ad}}\}$$

for any convex set  $X^{\text{ad}}$ .

As in finite-dimensions, for general functionals and sets, we would like to replace  $T(X^{\text{ad}}; \bar{x})$  by a simpler object.

**Definition 3.2.6.** Let  $G: X \rightarrow Z$  be Fréchet-differentiable at  $\bar{x} \in X^{\text{ad}}$ . Then we define the *linearizing cone* at  $\bar{x}$  as

$$T_l(G, K, \bar{x}) = \{d \in X \mid G'(\bar{x})d \in T(K, G(\bar{x}))\}.$$

note, that we keep the Tangent-cone to the convex set  $K$ , as it is 'easy' to calculate following Remark 3.2.5. As in finite-dimensions, we have the inclusion

**Lemma 3.2.7.** Let  $G$  be Fréchet-differentiable at  $\bar{x} \in X^{\text{ad}}$ . Then it holds

$$T(X^{\text{ad}}; \bar{x}) \subset T_l(G, K, \bar{x}).$$

*Beweis.* Let  $d \in T(X^{\text{ad}}; \bar{x})$  be given with corresponding sequences  $t_k \downarrow 0$  and  $X^{\text{ad}} \ni x^k \rightarrow \bar{x}$ . Then

$$\frac{G(x^k) - G(\bar{x})}{t_k} = G'(\bar{x}) \frac{x^k - \bar{x}}{t_k} + \frac{\omega(x^k - \bar{x})}{\|x^k - \bar{x}\|_X} \frac{\|x^k - \bar{x}\|_X}{t_k} \rightarrow G'(\bar{x})d.$$

Further,  $K \ni G(x^k) \rightarrow G(\bar{x}) \in K$ , and thus

$$G'(\bar{x})d = \lim_{k \rightarrow \infty} \frac{G(x^k) - G(\bar{x})}{t_k} \in T(K, G(\bar{x})).$$

$\square$

Hence, in order to replace the Bouligand-cone by the linearizing cone, we need to assume



### 3.2. First Order Necessary Optimality Conditions

**Definition 3.2.8.** A point  $\bar{x} \in X^{\text{ad}}$  satisfies the *Abadie Constraint Qualification* (ACQ), iff it holds

$$T_l(G, K, \bar{x}) \subset T(X^{\text{ad}}, \bar{x}). \quad (\text{ACQ})$$

As in finite dimensions, we would like to have more easily accessible conditions to verify if (ACQ) holds.

**Definition 3.2.9.** A point  $\bar{x} \in X^{\text{ad}}$  is said to satisfy *Robinson's Constraint Qualification* (RCQ) [Robinson \[1976\]](#), iff it holds

$$0 \in \text{int}(G(\bar{x}) + G'(\bar{x})X - K). \quad (\text{RCQ})$$

Before we can show, that (RCQ) is in fact a constraint qualification, we need a preparatory result on so called *metric regularity*, see, e.g., [[Bonnans and Shapiro, 1998](#), Proposition 3.3] for the statement or [[Bonnans and Shapiro, 2000](#), Theorem 2.87] for the proof.

**Lemma 3.2.10.** Let  $G: X \rightarrow Z$  be continuously Fréchet-differentiable near  $\bar{x} \in X^{\text{ad}} = G^{-1}(K)$ , where as before  $K \subset Z$  is closed and convex. Then assuming (RCQ) to hold at  $\bar{x}$ , there exist constants  $c, \delta > 0$  such that

$$\text{dist}(x, G^{-1}(K - z)) \leq c \text{dist}(G(x) + z, K)$$

holds for all  $x \in B_{\delta}^X(\bar{x}) \subset X$  and  $z \in B_{\delta}^Z(0) \subset Z$ . As usual in the above, we have

$$G^{-1}(K - z) = \{x \in X \mid G(x) + z \in K\}, \quad \text{dist}(b, A) = \inf_{a \in A} \|a - b\|.$$

**Remark 3.2.11.** While we do not provide the technical proof of this result –compare [[Bonnans and Shapiro, 2000](#), Theorem 2.87]–, we would like to motivate it in case the stronger (Exercise) constraint qualification that  $G'(\bar{x})$  is an isomorphism holds. As we will only need the case  $z = 0$ , we will fix  $z = 0$  for this. Since  $G'(\bar{x})$  is an isomorphism, we obtain from the implicit function Theorem 3.1.6-5. that there exists neighborhoods  $B_{2\delta_x}^X(\bar{x}) \subset X$  and  $B_{2\delta_z}^Z(G(\bar{x})) \subset Z$  such that  $G: B_{2\delta_x}^X(\bar{x}) \rightarrow B_{2\delta_z}^Z(G(\bar{x}))$  is invertible and that

$$\sup_{x \in B_{2\delta_x}^X(\bar{x})} \|G'(x)\|_{L(X;Z)} + \sup_{z \in B_{2\delta_z}^Z(G(\bar{x}))} \|(G^{-1})'(z)\|_{L(Z;X)} \leq C.$$

**draw sketch** Now, let  $x \in B_{\delta_x}^X(\bar{x})$  and assume (potentially shrinking  $\delta_x$ ) that  $G(B_{\delta_x}^X(\bar{x})) \subset$

### 3. Optimality Conditions

$B_{\delta_z}^Z(G(\bar{x}))$ . Then, taking a minimizing sequence  $z^k \in K$  of the distance to  $G(x)$ , i.e.,

$$\|G(x) - z^k\|_Z \rightarrow \inf_{k \in K} \|G(x) - k\|_Z$$

we notice, that since  $G(\bar{x}) \in K$  we can assume that  $z^k \in B_{2\delta_z}^Z(G(\bar{x}))$  and thus define  $x^k = G^{-1}(z^k)$ . Then we obtain, with the mean-value Theorem 3.1.6-2

$$\begin{aligned} \text{dist}(x, G^{-1}(K)) &\leq \|x - x^k\|_X \\ &= \|G^{-1}(G(x)) - G^{-1}(z^k)\|_X \\ &\leq \sup_{z \in B_{2\delta_z}^Z(G(\bar{x}))} \|(G^{-1})'(z)\|_{L(Z, X)} \|G(x) - z^k\|_Z \\ &\leq C \|G(x) - z^k\|_Z \rightarrow C \text{dist}(G(x), K). \end{aligned}$$

We are now in position to show that (RCQ) is a constraint qualification.

**Theorem 3.2.12.** *If  $G$  is continuously Fréchet-differentiable near  $\bar{x} \in X^{\text{ad}}$  and (RCQ) holds at  $\bar{x}$ . Then (ACQ) is satisfied at  $\bar{x}$ , i.e.,  $T(X^{\text{ad}}, \bar{x}) = T_l(G, K, \bar{x})$ .*

*Beweis.* In view of Lemma 3.2.7, it is sufficient to show  $T_l(G, K, \bar{x}) \subset T(X^{\text{ad}}, \bar{x})$ . Hence, let  $d \in T_l(G, K, \bar{x})$  be given, and let  $y = G'(\bar{x})d \in T(K, G(\bar{x}))$ . By definition there exists  $z^k \in K$  with  $z^k \rightarrow G(\bar{x})$  and  $t_k \downarrow 0$  such that

$$y^k = \frac{z^k - G(\bar{x})}{t_k} \rightarrow y.$$

**draw sketch**

Taylor-expansion gives

$$\begin{aligned} G(\bar{x} + t_k d) &= G(\bar{x}) + t_k G'(\bar{x})d + \omega(t_k d) \\ &= z^k + (G(\bar{x}) - z^k) + t_k y + \omega(t_k d) \\ &= z^k + t_k(y - y^k) + \omega(t_k d), \end{aligned}$$

with  $\omega(t_k d) \in o(t_k \|d\|_X)$ . Taking  $x = \bar{x} + t_k d$  and  $z = 0$  in Lemma 3.2.10 (and  $k$  sufficiently large so that  $x \in B_{\delta}^X(\bar{x})$ ), we get, using  $z^k \in K$ , that

$$\begin{aligned} \text{dist}(\bar{x} + t_k d, X^{\text{ad}}) &\leq c \text{dist}(G(\bar{x} + t_k d), K) \\ &= c \text{dist}(z^k + t_k(y - y^k) + \omega(t_k d), K) \\ &\leq c \|(z^k + t_k(y - y^k) + \omega(t_k d)) - z^k\|_Z \\ &= c \|t_k(y - y^k) + \omega(t_k d)\|_Z. \end{aligned}$$

### 3.2. First Order Necessary Optimality Conditions

Consequently, for  $k$  large enough, there is  $x^k \in X^{\text{ad}}$  such that

$$\|(\bar{x} + t_k d) - x^k\|_X \leq c \|t_k(y - y^k) + \omega(t_k d)\|_Z + \frac{t_k}{k}.$$

Dividing by  $t_k$  asserts

$$\left\| \frac{x^k - \bar{x}}{t_k} - d \right\|_X \leq c \|y - y^k\|_Z + \frac{\|\omega(t_k d)\|_Z}{t_k} + \frac{1}{k} \rightarrow 0.$$

Hence, we have found  $x^k \in X^{\text{ad}}$  with  $\frac{x^k - \bar{x}}{t_k} \rightarrow d$  and thus  $d \in T(X^{\text{ad}}, \bar{x})$ .  $\square$

Before we continue, we would like to state a few additional constraint qualifications that are potentially easier to verify.

**Definition 3.2.13.** We say, that a point  $\bar{x} \in X^{\text{ad}}$  satisfies the *linearized Slater* constraint qualification if there exists  $d \in X$  such that

$$G(\bar{x}) + G'(\bar{x})d \in \text{int}(K). \quad (3.2)$$

Clearly, this condition requires that  $\text{int}(K) \neq \emptyset$ .

**Lemma 3.2.14.** Assume that  $\text{int}(K) \neq \emptyset$ , then (3.2) is equivalent to (RCQ).

*Beweis.* Let (3.2) be satisfied. Then, for  $\delta > 0$  sufficiently small, it is  $G(\bar{x}) + G'(\bar{x})d + B_\delta^Z(0) \subset K \subset Z$ . Consequently, it holds

$$B_\delta^Z(0) \subset G(\bar{x}) + G'(\bar{x})d - K \subset G(\bar{x}) + G'(\bar{x})X - K$$

and thus (RCQ) holds.

Conversely, if (3.2) is violated, then the convex sets  $G(\bar{x}) + G'(\bar{x})X$  and  $\text{int}(K)$  have an empty intersection. By the Hahn-Banach theorem, there exists a separating hyperplane, i.e.,  $z^* \in Z^* \setminus \{0\}$  such that for any  $d \in X$  and  $k \in K$  it holds

$$\langle z^*, G(\bar{x}) + G'(\bar{x})d \rangle_{Z^*, Z} \geq \langle z^*, k \rangle_{Z^*, Z}.$$

Now, let  $z \in Z$  be such that  $\langle z^*, z \rangle < 0$ . Then for any  $t > 0$ ,  $d \in X$  and  $k \in K$  it is

$$\langle z^*, G(\bar{x}) + G'(\bar{x})d - k \rangle_{Z^*, Z} \geq 0 > \langle z^*, tz \rangle.$$

Consequently,  $tz \notin G(\bar{x}) + G'(\bar{x})X - K$  and since  $tz \rightarrow 0$  (RCQ) can not hold.  $\square$

### 3. Optimality Conditions

**Remark 3.2.15.** It is easy to see, that if  $G'(\bar{x}) \in \mathcal{L}(X; Z)$  is surjective, then (RCQ) holds. This corresponds to the well known linear independence constraint qualification.

Finally, we want to see, that (RCQ) corresponds to the well known Mangasarian Fromovitz constraint qualification (MFCQ) in finite dimensions.

To see this, consider the NLP

$$\min f(x) \quad \text{s.t.} \quad g(x) \leq 0, \quad h(x) = 0$$

with differentiable functions  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , and  $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$ . In our notation, this means  $X = \mathbb{R}^n$ ,  $Z = \mathbb{R}^m \times \mathbb{R}^p$ ,  $G(x) = \begin{pmatrix} g(x) \\ h(x) \end{pmatrix}$ ,  $K = (-\infty, 0]^m \times \{0\}^p$ . Now, we want to see, that (RCQ) is equivalent to the MFCQ

$$\text{rank } \nabla h(\bar{x}) = p, \quad \text{there exists } d \in \mathbb{R}^n: \nabla h(\bar{x})^T d = 0, \quad \nabla g_i(\bar{x})d < 0 \text{ whenever } g_i(\bar{x}) = 0.$$

“ $\Rightarrow$ ” Now, let  $\bar{x} \in X^{\text{ad}}$  be given, and (RCQ) be satisfied, i.e.,

$$0 \in \text{int} \left\{ \begin{pmatrix} g(\bar{x}) + \nabla g(\bar{x})^T d - v \\ \nabla h(\bar{x})^T d \end{pmatrix} \mid d \in \mathbb{R}^n, v \in \mathbb{R}^m, v \leq 0 \right\}.$$

The lower line implies that  $\nabla h(\bar{x})$  is surjective, i.e.,  $\text{rank } \nabla h(\bar{x}) = p$ . For the first line, let  $\delta > 0$  be sufficiently small, and set  $w \in \mathbb{R}^m$  with  $w_i = -\delta$  if  $g_i(\bar{x}) = 0$  and  $w_i = 0$  otherwise. Then, by (RCQ) there is  $d \in \mathbb{R}^n$ , and  $v \leq 0$  solving

$$\begin{pmatrix} w \\ 0 \end{pmatrix} = \begin{pmatrix} g(\bar{x}) + \nabla g(\bar{x})^T d - v \\ \nabla h(\bar{x})^T d \end{pmatrix}$$

hence  $\nabla h(\bar{x})^T d = 0$  and for all  $i$  with  $g_i(\bar{x}) = 0$  it is

$$\nabla g_i(\bar{x})^T d = w_i + v_i = -\delta + v_i \leq -\delta < 0$$

and hence MFCQ holds.

“ $\Leftarrow$ ” Conversely, let MFCQ hold. Rescaling  $d$  from MFCQ, we can assert

$$\begin{aligned} \nabla h(\bar{x})^T d &= 0 \\ g_i(\bar{x}) + \nabla g_i(\bar{x})^T d &< -2\delta \end{aligned}$$

for all  $i$  and a sufficiently small  $\delta > 0$ . Now, we can pick  $\rho > 0$  such that for any  $d_0 \in B_\rho^{\mathbb{R}^n}(0)$  it is

$$g_i(\bar{x}) + \nabla g_i(\bar{x})^T (d + d_0) < -\delta.$$

Since  $\nabla h(\bar{x})$  has rank  $p$  there is some  $\varepsilon \in (0, \delta]$  such that

$$B_\varepsilon^{\mathbb{R}^p}(0) \subset \nabla h(\bar{x})^T B_\rho^{\mathbb{R}^n}(0).$$

### 3.3. The Karush-Kuhn-Tucker conditions

Now, let  $z_1 \in \mathbb{R}^m$  and  $z_2 \in \mathbb{R}^p$  with  $\|z_i\| < \varepsilon$  be arbitrary. Then there is  $d_0 \in B_\rho^{\mathbb{R}^n}(0)$  such that

$$z_2 = \nabla h(\bar{x})^T d_0 = \nabla h(\bar{x})^T (d + d_0)$$

and

$$-(z_1)_i + [g_i(\bar{x}) + \nabla g_i(\bar{x})^T (d + d_0)] < \varepsilon - \delta \leq 0.$$

Hence we can set  $v = -z_1 + g(\bar{x}) + \nabla g(\bar{x})^T (d + d_0) \in (-\infty, 0]^m$  and have

$$\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} g(\bar{x}) + \nabla g(\bar{x})^T (d + d_0) - v \\ \nabla h(\bar{x})^T (d + d_0) \end{pmatrix}$$

and thus (RCQ) holds.

### 3.3. The Karush-Kuhn-Tucker conditions

Similar to the finite-dimensional NLP case, we would like to utilize the simpler structure of the linearizing cone  $T_l(G, K, \bar{x})$  to derive optimality conditions of Karush-Kuhn-Tucker (KKT) type.

Before we can state the KKT-conditions, we need to generalize the sign conditions we are used to from NLPs.

**Definition 3.3.1.** Let  $C \subset Z$  be a cone. Then we define the *polar cone*  $C^\circ$  of  $C$  by

$$C^\circ = \{z^* \in Z^* \mid \langle z^*, z \rangle_{Z^*, Z} \leq 0 \ \forall z \in C\}.$$

The polar cone is always closed and convex (exercise!).

**Theorem 3.3.2** (KKT-Conditions). *Let  $X, Z$  be Banach-spaces and  $K \subset Z$  be closed and convex. Further, let  $\bar{x}$  be a local solution of (P) at which (RCQ) is satisfied (and of course where  $f : X \rightarrow \mathbb{R}$  and  $G : X \rightarrow Z$  are Fréchet-differentiable). Then there exists a Lagrange multiplier  $\bar{\lambda} \in Z^*$ , i.e.,*

$$\begin{aligned} f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} &= 0, \\ G(\bar{x}) &\in K, \\ \bar{\lambda} &\in T(K, G(\bar{x}))^\circ. \end{aligned}$$

### 3. Optimality Conditions

Moreover, the set

$$\Lambda(\bar{x}) = \{\lambda \in T(K, G(\bar{x}))^\circ \mid f'(\bar{x}) + G'(\bar{x})^* \lambda = 0\}$$

is closed, convex, and bounded.

Before we come to the proof of the KKT-conditions, we need to collect a few preparatory results.

Clearly, (RCQ) is equivalent to the existence of  $\delta > 0$  such that the ball  $B_\delta^Z(0) \subset Z$  satisfies

$$B_\delta^Z(0) \subset G'(\bar{x})X - (K - G(\bar{x})).$$

Multiplication with  $\lambda > 0$  implies that

$$Z = G'(\bar{x})X - \text{cone}(K - G(\bar{x})) \quad (3.3)$$

where

$$\text{cone}(K - G(\bar{x})) = \{d \in Z \mid d = \lambda(k - G(\bar{x})) \text{ for some } \lambda \geq 0, k \in K\}$$

is the *radial cone* of  $K$  at  $G(\bar{x})$ . Hence (RCQ) implies (3.3). In fact it can be shown, e.g., [Zowe and Kurcyusz, 1979, Theorem 2.1] or [Ito and Kunisch, 2008, Theorem 1.4], that both conditions are equivalent.

Due to Remark 3.2.5 this shows that (RCQ) also implies

$$Z = G'(\bar{x})X - T(K, G(\bar{x})). \quad (3.4)$$

We will later need to use this representation for 'small'  $z \in Z$  and need to assert, that these small values can be reached by small values in  $x \in X$  and  $k \in T(K, G(\bar{x}))$ . To this end, we require the following:

**Theorem 3.3.3** (Generalized Open Mapping Theorem). *Let  $X, Z$  be Banach-spaces and  $\Psi: X \rightrightarrows Z$  be a closed convex multifunction, i.e., its graph*

$$\text{graph } \Psi = \{(x, z) \in X \times Z \mid z \in \Psi(x)\}$$

*is closed and convex. Then for any  $z \in \text{int}(\Psi(X))$ , all  $x \in \Psi^{-1}(z)$  and any  $\delta > 0$  it is*

$$z \in \text{int } \Psi(B_\delta^X(x)).$$

For a proof, see [Bonnans and Shapiro, 2000, Theorem 2.70].

We can now apply this to the set valued mapping corresponding to (3.4), i.e.,  $\Psi: X \times Z \rightrightarrows Z$  with

$$\Psi(x, k) = \begin{cases} \{G'(\bar{x})x - k\} & k \in T(K, G(\bar{x})), \\ \emptyset & \text{otherwise.} \end{cases}$$

### 3.3. The Karush-Kuhn-Tucker conditions

Indeed, this mapping is closed, since for  $(x^n, k^n, z^n) \in \text{graph } \Psi$  with  $(x^n, k^n, z^n) \rightarrow (x, k, z)$  it holds  $k^n \in T(K, G(\bar{x}))$  and  $z^n = G'(\bar{x})x^n - k^n$ . Thus, since  $T(K, G(\bar{x}))$  is closed,  $k \in T(K, G(\bar{x}))$  and further

$$k = \lim_{n \rightarrow \infty} k^n = \lim_{n \rightarrow \infty} (G'(\bar{x})x^n - z^n) = G'(\bar{x})x - z$$

and thus  $(x, k, z) \in \text{graph } \Psi$ . Convexity of  $\text{graph } \Psi$ , i.e., for any  $(x_i, k_i) \in X \times Z$  ( $i = 1, 2$ ) and any  $t \in [0, 1]$  it is

$$t\Psi(x_1, k_1) + (1-t)\Psi(x_2, k_2) \subset \Psi(t(x_1, k_1) + (1-t)(x_2, k_2)),$$

follows immediately, since  $x, k$  appear linearly in  $\Psi$  whenever  $\Psi(x, k) \neq \emptyset$ , convexity of  $T(K, G(\bar{x}))$ , and  $\emptyset + M = \emptyset$  for any set  $M$ .

Now, if (3.4) holds we have  $\Psi(X \times Z) = Z$ , and  $\Psi(0, 0) = 0$ . The generalized open mapping Theorem 3.3.3 yields  $0 \in \text{int } \Psi(B_1^X(0) \times B_1^Z(0))$ . More explicitly written, this is

$$0 \in \text{int} \left( G'(\bar{x})B_1^X(0) - (T(K, G(\bar{x})) \cap B_1^Z(0)) \right). \quad (3.5)$$

It is also true, that (3.5) implies (3.4) which can be seen by multiplication with  $t > 0$ .

Hence, we have

$$(\text{RCQ}) \iff (3.3) \Rightarrow (3.4) \iff (3.5)$$

*Beweis.* [Of Theorem 3.3.2] As usual, existence of Lagrange multipliers is based on a Hahn-Banach separation type argument, where the separating hyperplane provides the desired multiplier.

To this end, we define  $M \subset \mathbb{R} \times Z$  as

$$M = \{(t, z) \in \mathbb{R} \times Z \mid \exists d \in X : t \geq f'(\bar{x})d, G'(\bar{x})d - z \in T(K, G(\bar{x}))\}.$$

**draw sketch, and discuss expectation on M in 1d** Since  $t, z, h$  appear only linearly, and  $T(K, G(\bar{x}))$  is closed and convex this set  $M$  is convex and closed.

Further, we notice, that  $(0, 0) \in M$  is a boundary point of  $M$ . To see this, let  $\tau > 0$  be given; and we show  $(-\tau, 0) \notin M$ . Assume for contradiction, that  $(-\tau, 0) \in M$ . Then there would exist some  $d \in X$  such that

$$f'(\bar{x})d \leq -\tau < 0, \quad \text{and} \quad G'(\bar{x})d \in T(K, G(\bar{x})).$$

The latter means  $d \in T_l(G, K, \bar{x})$  and the necessary optimality condition in Theorem 3.2.3 and Theorem 3.2.12 give

$$f'(\bar{x})d \geq 0$$

contradicting our assumption.

Next, we need to see that  $\text{int } M \neq \emptyset$  so that we can utilize Hahn-Banach's theorem to separate  $\text{int } M$  from  $(0, 0)$ . This can be done by (RCQ) and the implied condition (3.5).

### 3. Optimality Conditions

Indeed (3.5) asserts that there is  $\delta > 0$  such that for any  $z \in B_\delta^Z(0)$  there exists  $d \in B_1^X(0)$  and  $k \in T(K, G(\bar{x})) \cap B_1^Z(0)$  such that

$$G'(\bar{x})d - z = k \in T(K, G(\bar{x})).$$

Moreover, for any  $t \geq \|f'(\bar{x})\|_{X^*}$  it is by definition

$$f'(\bar{x})d \leq \|f'(\bar{x})\|_{X^*}\|d\|_X \leq \|f'(\bar{x})\|_{X^*} \leq t$$

and hence

$$[\|f'(\bar{x})\|_{X^*}, \infty) \times B_\delta^Z(0) \subset M.$$

Now, Hahn-Banach asserts the existence of  $(\alpha, \lambda) \in (\mathbb{R} \times Z^*) \setminus \{(0, 0)\}$  such that

$$\alpha t + \langle \lambda, z \rangle_{Z^*, Z} \geq 0 \quad \forall (t, z) \in M.$$

In more detail this means, for any  $\tau > 0$ ,  $d \in X$ ,  $k \in T(K, G(\bar{x}))$  it holds

$$\alpha(f'(\bar{x})d + \tau) + \langle \lambda, G'(\bar{x})d - k \rangle_{Z^*, Z} \geq 0.$$

We now need to show that  $\alpha > 0$ .

We start by fixing  $d = 0$ ,  $k = 0$ ,  $\tau \rightarrow \infty$  to get  $\alpha \geq 0$ . Now, to exclude the case  $\alpha = 0$ , assume for contradiction that  $\alpha = 0$  we have

$$\langle \lambda, G'(\bar{x})d - k \rangle_{Z^*, Z} \geq 0 \quad \forall d \in X, k \in T(K, G(\bar{x}))$$

by (3.4)  $G'(\bar{x})d - k$  can take any value in  $Z$  and thus  $\lambda = 0$  contradicting  $(\alpha, \lambda) \neq (0, 0)$ .

Now, we can rescale the multiplier, to get  $(1, \bar{\lambda}) = (1, \frac{1}{\alpha}\lambda)$ , to have

$$(f'(\bar{x})d + \tau) + \langle \bar{\lambda}, G'(\bar{x})d - k \rangle_{Z^*, Z} \geq 0 \quad \forall \tau \geq 0, d \in X, k \in T(K, G(\bar{x})).$$

Taking  $d = 0$ ,  $\tau = 0$ , and  $k \in T(K, G(\bar{x}))$  arbitrary gives

$$\bar{\lambda} \in T(K, G(\bar{x}))^\circ.$$

Now, taking  $k = 0$ ,  $\tau = 0$  and  $d \in X$  arbitrary shows

$$f'(\bar{x})d + \langle \bar{\lambda}, G'(\bar{x})d \rangle_{Z^*, Z} \geq 0 \quad \forall d \in X$$

or equivalently, by definition of  $G'(\bar{x})^*$

$$\langle f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda}, d \rangle_{Z^*, Z} \geq 0 \quad \forall d \in X.$$

Since  $d \in X$  is arbitrary this shows

$$f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} = 0$$



### 3.3. The Karush-Kuhn-Tucker conditions

and thus  $\bar{\lambda}$  is the desired multiplier.

To see, the claim on  $\Lambda(\bar{x})$ , we note that closedness and convexity follow immediately from these properties of  $T(K, G(\bar{x}))^\circ$  and the fact that  $\lambda$  appears only linearly in the confining equation.

To see boundedness, let  $\lambda \in \Lambda(\bar{x})$  be given, picking  $(1, \lambda) \in \mathbb{R} \times Z^*$  it holds for arbitrary  $d \in X$ ,  $k \in T(K, G(\bar{x}))$ , and  $\tau \geq 0$  that

$$\begin{aligned} (f'(\bar{x})d + \tau) + \langle \lambda, G'(\bar{x})d - k \rangle_{Z^*, Z} \\ = \langle f'(\bar{x}) + G'(\bar{x})^* \lambda, d \rangle_{Z^*, Z} + \tau - \langle \lambda, k \rangle_{Z^*, Z} \\ = \tau - \langle \lambda, k \rangle_{Z^*, Z} \geq 0 \end{aligned}$$

by definition of  $T(K, G(\bar{x}))^\circ$ .

Hence the normal  $(1, \lambda)$  supports  $M$  at  $(0, 0)$ . Further, we know that for some  $\delta > 0$  we have  $[\|f'(\bar{x})\|_{X^*}, \infty) \times B_\delta^Z(0) \subset M$ . Now, choosing  $z_\lambda \in B_\delta^Z(0)$  with  $\langle \lambda, z_\lambda \rangle_{Z^*, Z} \leq -\frac{\delta}{2} \|\lambda\|_{Z^*}$  and  $\tau = \|f'(\bar{x})\|_{X^*}$  we obtain

$$0 \leq 1 \cdot \|f'(\bar{x})\|_{X^*} + \langle \lambda, z_\lambda \rangle_{Z^*, Z} \leq \|f'(\bar{x})\|_{X^*} - \frac{\delta}{2} \|\lambda\|_{Z^*}$$

and thus

$$\|\lambda\|_{Z^*} \leq \frac{2}{\delta} \|f'(\bar{x})\|_{X^*}.$$

□

**Remark 3.3.4.** Utilizing the *Lagrange* function

$$\mathcal{L} : X \times Z^* \rightarrow \mathbb{R}; \quad \mathcal{L}(x, \lambda) = f(x) + \langle \lambda, G(x) \rangle_{Z^*, Z}$$

the KKT-conditions can be written in very short notation since

$$\mathcal{L}_x(x, \lambda) = f'(x) + G'(x)^* \lambda.$$

In fact, not only can we characterize solutions to

$$f'(\bar{x})d \geq 0 \quad \forall d \in T_l(G, K, \bar{x}).$$

by the KKT-conditions provided (RCQ) holds in such a solution, but also the converse holds true.

**Lemma 3.3.5.** If  $\bar{x}$  is a KKT-point of (P), then

$$f'(\bar{x})d \geq 0 \quad \forall d \in T_l(G, K, \bar{x}).$$

### 3. Optimality Conditions

*Beweis.* If  $d \in T_l(G, K, \bar{x})$  then  $G'(\bar{x})d \in T(K, G(\bar{x}))$  and consequently, for any  $\lambda \in \Lambda(\bar{x}) \subset T(K, G(\bar{x}))^\circ$  it holds

$$f'(\bar{x})d = -\langle \lambda, G'(\bar{x})d \rangle_{Z^*, Z} \geq 0.$$

□

#### 3.3.1. Cone Constraints

We will discuss the case, if we add the assumption, that the convex, closed set  $K \subset Z$  is in addition a *cone*, i.e., for any  $z \in K$  and  $\lambda \geq 0$  it is

$$\lambda z \in K.$$

If this cone is, in addition, *salient* (im deutschen 'regulär' oder 'spitz'), i.e.,

$$z \in K \text{ and } -z \in K \iff z = 0,$$

then we can define a partial ordering by

$$z \leq_K 0 \iff z \in K,$$

or for any  $z_1, z_2 \in Z$

$$z_1 \leq_K z_2 \iff z_1 - z_2 \in K$$

Hence, in this case the constraint  $G(x) \in K$  can be seen as a generalized inequality constraint. In order to rewrite the KKT-conditions 3.3.2 in this case, we note the following.

**Lemma 3.3.6.** *If  $K \subset Z$  is a closed, convex cone, then for any  $\bar{z} \in K$  it is*

$$T(K, \bar{z})^\circ = K^\circ \cap \{\bar{z}\}^\perp$$

*where the annihilator is given by*

$$\{\bar{z}\}^\perp = \{\lambda \in Z^* \mid \langle \lambda, \bar{z} \rangle_{Z^*, Z} = 0\}.$$

*Beweis.* “ $\supset$ ” Let  $\lambda \in K^\circ \cap \{\bar{z}\}^\perp$  be given. Then for any  $d \in T(K, \bar{z})$ , with corresponding sequence  $t_k \downarrow 0$ ,  $z^k \in K$ , it is

$$\langle \lambda, d \rangle_{Z^*, Z} \leftarrow \langle \lambda, \frac{(z^k - \bar{z})}{t_k} \rangle_{Z^*, Z} = \frac{1}{t_k} \langle \lambda, z^k \rangle_{Z^*, Z} - \frac{1}{t_k} \langle \lambda, \bar{z} \rangle_{Z^*, Z} = \frac{1}{t_k} \langle \lambda, z^k \rangle_{Z^*, Z} \leq 0.$$

Hence,  $\lambda \in T(K, \bar{z})^\circ$ .

” $\subset$ ” Conversely, let  $\lambda \in T(K, \bar{z})^\circ$  be given.

### 3.3. The Karush-Kuhn-Tucker conditions

We start by showing  $\lambda \in K^\circ$ . For any  $z \in K$  it is

$$z + \bar{z} = 2\left(\frac{1}{2}z + \frac{1}{2}\bar{z}\right) \in K$$

since  $K$  is a convex cone. By Remark 3.2.5 it is thus

$$z = (z + \bar{z}) - \bar{z} \in T(K, \bar{z})$$

hence  $\lambda \in K^\circ$ .

Now, to see that  $\lambda \in \{\bar{z}\}^\perp$ , we note that analogously,  $\bar{z} = 2\bar{z} - \bar{z} \in T(K, \bar{z})$  and  $-\bar{z} = 0 - \bar{z} \in T(K, \bar{z})$ . As a consequence

$$\langle \lambda, \pm \bar{z} \rangle_{Z^*, Z} \leq 0$$

and thus  $\lambda \in \{\bar{z}\}^\perp$ . □

We have thus learned, that for a closed, convex cone, the second part of the KKT-conditions

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in T(K, G(\bar{x}))^\circ$$

can equivalently be written as a *complementarity condition*

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in K^\circ, \quad \langle \lambda, G(\bar{x}) \rangle_{Z^*, Z} = 0.$$

#### 3.3.2. Slater Condition

**Definition 3.3.7.** Let  $K \subset Z$  be a closed, convex cone, and  $G: X \rightarrow Z$  be convex with respect to, the partial order given by  $K$ , i.e., for any  $t \in [0, 1]$  and all  $x_1, x_2 \in X$  it holds

$$G((1-t)x_1 + tx_2) - (1-t)G(x_1) - tG(x_2) \in K.$$

We say that the *Slater-condition* holds if there exists some  $\hat{x} \in X^{\text{ad}}$  such that

$$G(\hat{x}) \in \text{int}(K).$$

**Theorem 3.3.8.** Let  $K$  be a closed, convex cone, and  $G: X \rightarrow Z$  Fréchet-differentiable and convex w.r.t  $K$ . Further, assume that (P) satisfies the Slater-condition. Then any  $\bar{x} \in X^{\text{ad}}$  satisfies (RCQ).

*Beweis.* Indeed, by convexity, of  $G$ , we have for any  $t \in (0, 1)$

$$K \ni \frac{G(\bar{x} + t(\hat{x} - \bar{x})) - (1-t)G(\bar{x}) - tG(\hat{x})}{t} = G(\bar{x}) + \frac{G(\bar{x} + t(\hat{x} - \bar{x})) - G(\bar{x})}{t} - G(\hat{x}).$$

### 3. Optimality Conditions

Taking the limit  $t \downarrow 0$  shows

$$G(\bar{x}) + G'(\bar{x})(\hat{x} - \bar{x}) \in K + G(\hat{x}).$$

Now, the Slater-condition asserts, that for some  $\varepsilon > 0$  it is  $G(\hat{x}) + B_\varepsilon^Z(0) \subset K$  and hence

$$G(\bar{x}) + G'(\bar{x})(\hat{x} - \bar{x}) + B_\varepsilon^Z(0) \subset K + G(\hat{x}) + B_\varepsilon^Z(0) \subset K + K = K,$$

where the last equality follows as in the previous Section 3.3.1 from the fact that  $K$  is a convex cone. We have thus shown, that the Slater-condition implies the linearized Slater condition (3.2). Thus we have shown the assertion by Lemma 3.2.14.  $\square$

Before we continue, we will discuss the KKT-conditions for two examples.

**Example 3.3.9.** First, we consider a control constrained optimal control problem:

$$\begin{aligned} \min J(q, u) \\ \text{s.t. } \begin{cases} E(q, u) = 0, \\ q \in Q^{\text{ad}} := \{p \in Q = L^2(\Omega) \mid a \leq p \leq b\}. \end{cases} \end{aligned}$$

Here  $E: Q \times U \rightarrow W$  and  $J: Q \times U \rightarrow \mathbb{R}$  are continuously differentiable,  $U, W$  are Banach-spaces,  $\Omega$  is a bounded domain, and  $a \leq b \in L^2(\Omega)$ .

As we have seen in the exercises, the set  $Q^{\text{ad}}$  has no interior in  $L^2(\Omega)$ . Hence we have to work with the (RCQ) directly. To this end, let us assume that  $(\bar{q}, \bar{u})$  is a local solution in which

$$E'_u(\bar{q}, \bar{u}) \text{ is surjective, i.e., } E'_u(\bar{q}, \bar{u})U = W.$$

Then, (RCQ) is satisfied. To see this, let  $X = Q \times U$ ,  $Z = W \times Q$ ,  $K = \{0\} \times Q^{\text{ad}}$ ,  $\bar{x} = (\bar{q}, \bar{u})$ , and  $G(q, u) = (E(q, u), q)^T \in Z$ . Then we have

$$G(\bar{x}) + G'(\bar{x})X - K = \left\{ \begin{pmatrix} E'_u(\bar{q}, \bar{u})v + E'_q(\bar{q}, \bar{u})p \\ \bar{q} + p - r \end{pmatrix} \mid v \in U, p \in Q, r \in Q^{\text{ad}} \right\}.$$

Now, for arbitrary  $z = (w, q) \in Z$ , let  $p = q$ ,  $r = \bar{q}$  and select  $v \in U$  such that  $E'_u(\bar{q}, \bar{u})v = w - E'_q(\bar{q}, \bar{u})p$ , which is possible since  $E'_u$  is surjective. This shows, that

$$z \in G(\bar{x}) + G'(\bar{x})X - K.$$

Since  $z \in Z$  was arbitrary, (RCQ) is shown.

The KKT-conditions 3.3.2 ensure the existence of  $\bar{\lambda} = (\bar{v}, \bar{\mu}) \in W^* \times Q^* = W^* \times Q$  such that

$$E'_q(\bar{q}, \bar{u})^* \bar{\mu} = -J'_q(\bar{q}, \bar{u}), \quad E'_u(\bar{q}, \bar{u})^* \bar{v} = -J'_u(\bar{q}, \bar{u}), \quad E(\bar{q}, \bar{u}) = 0, \quad \bar{u} \in Q^{\text{ad}}, \quad \bar{\mu} \in T(Q^{\text{ad}}, \bar{q})^\circ.$$

Now, we would like to extract sign conditions from this optimality conditions. To this end, we need the pointwise representation

$$T(Q^{\text{ad}}, \bar{q}) = \{d \in L^2(\Omega) \mid d|_{\bar{q}=a} \geq 0, d|_{\bar{q}=b} \leq 0\}. \quad (3.6)$$

To see that this representation holds, we start by taking any  $d$  of the form on the right of (3.6). We can define

$$d^k = \frac{q^k - \bar{q}}{t_k}$$

with  $t_k \downarrow 0$  and

$$q^k = \min(\max(\bar{q} + t_k d), a), b) \in Q^{\text{ad}}.$$

Now, we note that almost all  $x \in \Omega$  are simultaneous Lebesgue points of  $d$ ,  $a$ ,  $b$ , and  $\bar{q}$ , cf., Theorem A.6.5. Hence for any such  $x$  it is immediately clear that  $d^k(x) \rightarrow d(x)$ . Further, by Lipschitz continuity of min and max we get

$$|d^k(x)| \leq |d(x)|$$

and Lebesgue's dominated convergence theorem shows  $d = \lim_{k \rightarrow \infty} d^k \in T(Q^{\text{ad}}, \bar{q})$ .

The reverse inclusion is simple to see, i.e., if

$$T(Q^{\text{ad}}, \bar{q}) \ni d = \lim_{k \rightarrow \infty} d^k = \lim_{k \rightarrow \infty} \frac{q^k - \bar{q}}{t_k},$$

with  $q^k \in Q^{\text{ad}}$ , then clearly for almost any  $x \in \{x \in \Omega \mid \bar{q}(x) = a(x)\}$  it is  $d^k(x) \geq 0$  and hence the same holds true for  $d$ . Analogously, the sign condition on  $\{x \in \Omega \mid \bar{q}(x) = b(x)\}$  is obtained.

Now, (3.6) implies

$$T(Q^{\text{ad}}, \bar{q})^\circ = \{\mu \in L^2(\Omega) \mid \mu|_{\bar{q}=a} \leq 0, \mu|_{\bar{q}=b} \geq 0, \mu|_{a < \bar{q} < b} = 0\}.$$

This shows the expected sign conditions on the multiplier, i.e.,

$$\mu(x) \begin{cases} \leq 0 & \text{a.e. on } \{\bar{q} = a\}, \\ \geq 0 & \text{a.e. on } \{\bar{q} = b\}, \\ = 0 & \text{otherwise.} \end{cases}$$

### 3. Optimality Conditions

**Example 3.3.10.** In the second example, we want to consider so called state constraints.

$$\begin{aligned} \min J(q, u) \\ \text{s.t. } \begin{cases} Au = Bq + b, \\ u - \psi \in U^{\text{ad}} := \{v \in C(\bar{\Omega}) \mid v \leq 0\}. \end{cases} \end{aligned}$$

For this, we make the following assumptions

1.  $\Omega \subset \mathbb{R}^n$ ,  $n = 1, 2, 3$  is a bounded Lipschitz domain,
2.  $b \in L^2(\Omega)$ ,  $\psi \in C(\bar{\Omega})$ ,  $B \in \mathcal{L}(Q, L^2(\Omega))$ ,  $Q$  a Banach space,
3.  $U = H_0^1(\Omega) \cap H^2(\Omega) \subset C(\bar{\Omega})$  (embedding ok for  $n = 1, 2, 3$ .)
4.  $A \in \mathcal{L}(H_0^1(\Omega), H^{-1}(\Omega))$  is a second order elliptic operator. In particular  $A^{-1} \in \mathcal{L}(H^{-1}(\Omega), H_0^1(\Omega))$  exists.
5.  $A$  is  $H^2$  regular, i.e.,  $A^{-1} \in \mathcal{L}(L^2(\Omega), U)$  (and  $A \in \mathcal{L}(U, L^2(\Omega))$ ). This can be asserted, if the domain has a sufficiently smooth ( $C^{1,1}$  or convex) boundary, and the coefficients in  $A$  are  $W^{1,\infty}(\Omega) = C^{0,1}(\Omega)$ .
6.  $J : Q \times U \rightarrow \mathbb{R}$  is continuously differentiable.

In this situation, we can not use the same trick as in the control constrained case, i.e., pick something suitable for the inequality constrained component in the (RCQ) and then correct in the equality component since  $B$  is not surjective, in general. Hence, we have to work in a stronger topology which allows for interior points – thus we picked  $U \subset C(\bar{\Omega})$ .

Now, let  $(\bar{q}, \bar{u})$  be a local solution in which the Slater-type condition

$$\text{there exists } \hat{q} \in Q, \hat{u} \in U \text{ satisfying } A\hat{u} = B\hat{q} + b, \hat{u} < \psi \text{ on } \bar{\Omega}$$

holds. To see that now (RCQ) holds, we set  $X = Q \times U$ ,  $Z = L^2(\Omega) \times C(\bar{\Omega})$ ,  $K = \{0\} \times U^{\text{ad}}$ , and

$$G(x) = \begin{pmatrix} Au - Bq - b \\ Iu - \psi \end{pmatrix}.$$

Notice, that in contrast to the previous example, the inequality constraint is considered in a different topology than the natural one. Thus we add the embedding  $I : Y \rightarrow C(\bar{\Omega})$ .

To see (RCQ), i.e.,

$$0 \in \text{int} \left\{ \begin{pmatrix} Av - Bp \\ (I\bar{u} - \psi) + Iv - w \end{pmatrix} \mid v \in U, p \in Q, w \in U^{\text{ad}} \right\}$$

we notice, that, by Weierstrass's extreme value theorem, there exists  $\varepsilon > 0$  such that

$$I\hat{u} - \psi \leq -\varepsilon < 0.$$

By continuity of the solution map  $A^{-1}: L^2(\Omega) \rightarrow U$ , we have that for  $\delta > 0$  small enough and any  $p \in B_{\delta}^{L^2(\Omega)}(0)$

$$IA^{-1}(B(\hat{q}) + b + p) = I\hat{u} + IA^{-1}p \in B_{\varepsilon/2}^{C(\bar{\Omega})}(0) + I\hat{u}$$

since  $\hat{u} = A^{-1}B(\hat{q} + b)$ .

Now, let  $z = (q, u) \in Z = L^2(\Omega) \times C(\bar{\Omega})$  with  $\|q\|_Q \leq \delta$  and  $\|u\|_{\infty} \leq \varepsilon/2$  be arbitrary. We can set  $v = A^{-1}(B(\hat{q} - \bar{q}) + q)$  and  $p = \hat{q} - \bar{q}$  asserting  $Av - Bp = q$ . Further, by definition of  $\varepsilon, \delta$ , it holds

$$\begin{aligned} (I\bar{u} - \psi) + Iv &= I\bar{u} - \psi + IA^{-1}(B(\hat{q} - \bar{q}) + b - b + q) \\ &= I\bar{u} - \psi + IA^{-1}(B(\hat{q}) + b + q) - I\bar{u} \\ &= IA^{-1}q + I\hat{u} - \psi \\ &\leq IA^{-1}q - \varepsilon \\ &\leq \frac{\varepsilon}{2} - \varepsilon = -\frac{\varepsilon}{2} < 0 \end{aligned}$$

and hence  $(I\bar{u} - \psi) + Iv \in \text{int } U^{\text{ad}}$  showing that (RCQ) holds.

Hence in the selected minimizer the KKT-conditions 3.3.2 hold, i.e., there exists  $\bar{\lambda} = (\bar{v}, \bar{\mu}) \in L^2(\Omega) \times C(\bar{\Omega})^*$  such that

$$\begin{aligned} A^*\bar{v} + I^*\bar{\mu} &= -J'_u(\bar{q}, \bar{u}), \\ J'_q(\bar{q}, \bar{u}) - B^*\bar{v} &= 0, \\ A\bar{u} &= B\bar{q} + b, \\ I\bar{u} - \psi &\leq 0, \\ \bar{\mu} &\in T(U^{\text{ad}}, I\bar{u} - \psi)^{\circ} = (U^{\text{ad}})^{\circ} \cap (I\bar{u} - \psi)^{\perp}, \end{aligned}$$

where the last identity follows from Section 3.3.1.

To understand what  $I^*$  in the first, so called *adjoint*, equation means, we write for any  $u \in U$

$$\langle I^*\bar{\mu}, u \rangle_{U^*, U} = \langle \bar{\mu}, Iu \rangle_{C^*, C} = \langle \bar{\mu}, u \rangle_{C^*, C}$$

Hence one often simply writes

$$A^*\bar{v} + \bar{\mu} = -J'_u(\bar{q}, \bar{u})$$

dropping the explicit notion of  $I$  and  $I^*$ .

To understand the complementarity condition

$$\bar{u} - \psi \leq 0, \quad \bar{\mu} \in (U^{\text{ad}})^{\circ}, \quad \langle \bar{\mu}, \bar{u} - \psi \rangle_{C^*, C} = 0,$$

### 3. Optimality Conditions

we remark, that it can be shown, that  $C(\bar{\Omega}) \simeq M(\bar{\Omega})$  is the space of real, regular Borel measures on  $\bar{\Omega}$  with

$$\|\bar{\mu}\|_{M(\bar{\Omega})} = |\mu|(\bar{\Omega}) = \sup_{\|v\|_{\infty} \leq 1, v \in C(\bar{\Omega})} \int_{\bar{\Omega}} v \, d\bar{\mu}.$$

With this the complementarity conditions can be written as

$$\bar{u} - \psi \leq 0, \quad \bar{\mu} \geq 0, \quad \bar{\mu}(\{\bar{u} - \psi < 0\}) = 0.$$

As a final remark, in many cases,  $J : Q \times U \rightarrow \mathbb{R}$  is also differentiable as a mapping on a larger space  $Q \times \tilde{U} \supset Q \times U$ . Then improved regularity of  $\bar{v}$  can be shown, however this is limited by the appearance of  $\bar{\mu}$  in the adjoint equation.

## 3.4. Sufficient optimality conditions

### 3.4.1. The convex case

**Theorem 3.4.1.** *Let  $X, Z$  be Banach-spaces,  $f : X \rightarrow \mathbb{R}$  be convex and Fréchet-differentiable, and  $G : X \rightarrow Z$  be Fréchet-differentiable and convex w.r.t the closed, convex cone  $K \subset Z$ . Then, any point  $\bar{x} \in X^{\text{ad}}$  satisfying the KKT-conditions is a global solution of (P).*

*Beweis.* Exercise. □

### 3.4.2. Second order sufficient conditions

As in finite dimensions, it is desirable to obtain sufficient conditions that are close to the necessary ones. It is thus necessary to consider growth conditions only on a very small set of critical directions. This is done by the following definition.

**Definition 3.4.2.** For  $\bar{x} \in X^{\text{ad}}$  we define the *critical cone* by

$$\begin{aligned} C(\bar{x}) &= \{d \in X \mid G'(\bar{x})d \in T(K, G(\bar{x})), f'(\bar{x})d \leq 0\} \\ &= \{d \in T_l(K, G, \bar{x}) \mid f'(\bar{x})d \leq 0\}. \end{aligned}$$

Clearly, if  $\bar{x}$  is a stationary point satisfying (ACQ), then no directions  $d$  with  $f'(\bar{x})d < 0$  exist, and the definition reduces to

$$C(\bar{x}) = \{d \in X \mid G'(\bar{x})d \in T(K, G(\bar{x})), f'(\bar{x})d = 0\}.$$



### 3.4. Sufficient optimality conditions

In finite dimensions, one would now formulate a sufficient condition by requiring

$$f'(\bar{x})d \geq 0 \quad \forall d \in T_l(K, G, \bar{x})$$

and

$$f''(\bar{x})(d, d) > 0 \quad \forall d \in C(\bar{x}) \setminus \{0\}$$

however, this is **not enough** in function spaces! To understand this, consider the following example.

**Example 3.4.3.** Let  $X = Z = l_2$  the Hilbert-space of square summable sequences, i.e.,

$$\|x\|_{l_2}^2 = \sum_{i=0}^{\infty} x_i^2.$$

Consider  $K = \{x \in l_2 \mid x_i \geq 0\}$ ,  $G(x) = x$ , and  $f(x) = (c, x)_{l_2} - (x, x)_{l_2}$  where  $c \in l_2$  with  $c_i > 0$  for all  $i$  is fixed.

Then at  $\bar{x} = 0$  it holds

$$T(K, G(\bar{x})) = T(K, 0) = K.$$

Further, for  $d \in K \setminus \{0\}$  it is

$$f'(\bar{x})d = (c, d)_{l_2} = \sum_{i=0}^{\infty} c_i d_i > 0$$

since  $c_i > 0$ ,  $d_i \geq 0$ , and  $d_j > 0$  at least for one  $j$ .

Consequently,  $C(\bar{x}) = \{0\}$ , i.e., all directions  $d \in T(K, 0) \setminus \{0\}$  are ascent directions. But  $\bar{x}$  is no local minimum of  $f$  on  $X^{\text{ad}} = K$ . To see this, let  $x^k = (2\delta_{ik}c_i)_{i \in \mathbb{N}}$ . Then  $\|x^k - \bar{x}\|_{l_2} = \|x^k\|_{l_2} = 2c_k \rightarrow 0$  as  $k \rightarrow \infty$ , and

$$f(x^k) = 2c_k^2 - 4c_k^2 = -2c_k^2 < 0 = f(\bar{x}).$$

Hence, the above sketched conditions can not be sufficient.

This shows, that the cone of critical directions is too small in infinite dimensions. Thus we consider the enlarged sets

**Definition 3.4.4.** For any  $\eta \geq 0$  and  $\bar{x} \in X^{\text{ad}}$ , we define the *approximate critical cone*

$$\begin{aligned} C_\eta(\bar{x}) &= \{d \in X \mid G'(\bar{x})d \in T(K, G(\bar{x})), f'(\bar{x})d \leq \eta \|d\|_X\} \\ &= \{d \in T_l(K, G, \bar{x}) \mid f'(\bar{x})d \leq \eta \|d\|_X\}. \end{aligned}$$

### 3. Optimality Conditions

As in finite dimensions, we need to assert that the boundary of the feasible set is not too 'wild', and can reasonably be approximated by the linearizing cone. This will be asserted by the condition that for  $X^{\text{ad}} \ni x \rightarrow \bar{x}$  it is

$$\text{dist}(x - \bar{x}, T_l(K, G, \bar{x})) = o(\|x - \bar{x}\|_X).$$

This is indeed the case, given that (RCQ) holds.

**Lemma 3.4.5.** *Let (RCQ) hold at  $\bar{x} \in X^{\text{ad}}$ , then there exists a map  $h: X^{\text{ad}} \rightarrow T_l(K, G, \bar{x})$  such that*

$$\|h(x) - (x - \bar{x})\|_X = o(\|x - \bar{x}\|_X) \quad \text{for } X^{\text{ad}} \ni x \rightarrow \bar{x}.$$

*Beweis.* For any  $x \in X^{\text{ad}}$  the Fréchet-differentiability of  $G$  asserts

$$G(x) = G(\bar{x}) + G'(\bar{x})(x - \bar{x}) + \omega(x)$$

with  $\omega(x) = o(\|x - \bar{x}\|_X)$ . By (3.5) there is  $\delta > 0$  such that

$$\text{cl} B_\delta^Z(0) \subset G'(\bar{x})B_1^X(0) - (T(K, G(\bar{x})) \cap B_1^Z(0)).$$

Hence, we can write

$$\delta \frac{\omega(x)}{\|\omega(x)\|_Z} = G'(\bar{x})\hat{s}(x) - \hat{k}(x)$$

with suitable  $\hat{s}(x) \in B_1^X(0) \subset X$  and  $\hat{k}(x) \in T(K, G(\bar{x})) \cap B_1^Z(0)$ . Rescaling with  $\frac{\|\omega(x)\|_Z}{\delta}$  asserts that

$$\omega(x) = G'(\bar{x})s(x) - k(x)$$

with  $s(x) \in X$  and  $k(x) \in T(K, G(\bar{x}))$  and

$$\|s(x)\|_X \leq \frac{\|\omega(x)\|_Z}{\delta} = o(\|x - \bar{x}\|_X), \quad \|k(x)\|_Z \leq \frac{\|\omega(x)\|_Z}{\delta} = o(\|x - \bar{x}\|_X).$$

Now, we set  $h(x) = x - \bar{x} + s(x)$ , and obtain the desired approximation property

$$\|h(x) - (x - \bar{x})\|_X = \|s(x)\|_X = o(\|x - \bar{x}\|_X)$$

and the feasibility  $G'(\bar{x})h(x) \in T(K, G(\bar{x}))$  from

$$\begin{aligned} G'(\bar{x})h(x) &= G'(\bar{x})(x - \bar{x}) + G'(\bar{x})s(x) \\ &= G'(\bar{x})(x - \bar{x}) + \omega(x) + k(x) \\ &= G(x) - G(\bar{x}) + k(x) \\ &\in K - G(\bar{x}) + T(K, G(\bar{x})) \\ &\subset T(K, G(\bar{x})). \end{aligned}$$

□

With this we can now show the sufficient conditions

### 3.4. Sufficient optimality conditions

**Theorem 3.4.6** (Second order sufficient conditions). *Let  $X, Z$  be Banach-spaces,  $K \subset Z$  closed and convex,  $f : X \rightarrow \mathbb{R}$  and  $G : X \rightarrow Z$  be twice Fréchet-differentiable. Further, let  $\bar{x} \in X^{\text{ad}}$  satisfy (RCQ) and KKT-conditions 3.3.2 with multiplier  $\bar{\lambda}$ , i.e.,*

$$\begin{aligned} f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} &= 0, \\ G(\bar{x}) &\in K, \\ \bar{\lambda} &\in T(K, G(\bar{x}))^\circ. \end{aligned}$$

*If in addition, the second order condition*

$$\mathcal{L}''_{xx}(\bar{x}, \bar{\lambda})(d, d) \geq \gamma \|d\|_X^2 \quad \forall d \in C_\eta(\bar{x})$$

*holds for some  $\gamma > 0$  and  $\eta > 0$ . Then  $\bar{x}$  is an isolated local solution and there exists  $\kappa, \delta > 0$  such that*

$$f(x) - f(\bar{x}) \geq \kappa \|x - \bar{x}\|_X^2 \quad \forall x \in X^{\text{ad}} \cap B_\delta^X(\bar{x}).$$

*Beweis.* Let  $\delta > 0$  be sufficiently small, to be fixed later. For any  $x \in X^{\text{ad}} \cap B_\delta^X(\bar{x})$ , set  $d = x - \bar{x}$ . By Lemma 3.4.5 we can write

$$d = h(x) + \omega(x), \quad \omega(x) = o(\|d\|_X)$$

with  $h(x) \in T_l(K, G, \bar{x})$ . In particular, it holds

$$\|h(x)\|_X = \|d\|_X + o(\|d\|_X).$$

We start assuming,  $h(x) \in C_\eta(\bar{x})$ . Since  $G(x) - G(\bar{x}) \in T(K, G(\bar{x}))$  it is  $\langle \lambda, G(x) - G(\bar{x}) \rangle_{Z^*, Z} \leq 0$  and thus

$$\begin{aligned} f(x) - f(\bar{x}) &\geq \mathcal{L}(x, \bar{\lambda}) - \mathcal{L}(\bar{x}, \bar{\lambda}) \\ &= \mathcal{L}'_x(\bar{x}, \bar{\lambda})d + \frac{1}{2} \mathcal{L}''_{xx}(\bar{x}, \bar{\lambda})(d, d) + o(\|d\|_X^2) \\ &= \frac{1}{2} \mathcal{L}''_{xx}(\bar{x}, \bar{\lambda})(h(x) + \omega(x), h(x) + \omega(x)) + o(\|d\|_X^2) \\ &\geq \frac{1}{2} \mathcal{L}''_{xx}(\bar{x}, \bar{\lambda})(h(x), h(x)) - c \|h(x)\|_X \|\omega(x)\|_X - c \|\omega(x)\|_X^2 + o(\|d\|_X^2) \\ &\geq \frac{\gamma}{2} \|h(x)\|_X^2 - c \|h(x)\|_X \|\omega(x)\|_X + o(\|d\|_X^2) \\ &\geq \frac{\gamma}{2} \|h(x)\|_X^2 - \frac{\gamma}{4} \|h(x)\|^2 - 2c^2 \|\omega(x)\|_X^2 + o(\|d\|_X^2) \\ &\geq \frac{\gamma}{4} \|h(x)\|^2 + o(\|d\|_X^2) \\ &\geq \frac{\gamma}{4} \|d\|^2 + o(\|d\|_X^2) \\ &\geq \frac{\gamma}{8} \|d\|^2 \end{aligned}$$

### 3. Optimality Conditions

once  $\delta$  is small enough to assert  $|o(\|d\|_X^2)| \leq \frac{\gamma}{8}\|d\|_X^2$ .

Otherwise,  $h(x) \notin C_\eta(\bar{x})$ , and thus

$$\begin{aligned} f(x) - f(\bar{x}) &= f'(\bar{x})d + o(\|d\|_X) \\ &= f'(\bar{x})h(x) + f'(\bar{x})\omega(x) + o(\|d\|_X) \\ &\geq \eta\|h(x)\|_X + o(\|d\|_X) \\ &= \eta\|d\|_X + o(\|d\|_X) \\ &\geq \eta\|d\|_X^2 \end{aligned}$$

once  $\delta$  is sufficiently small.

□

## 4. Adjoint Approach

Many optimization problems are of the form  $X = Q \times U$  with

$$\min_{(q,u)} J(q,u) \quad \text{s.t.} \quad E(q,u) = 0, (q,u) \in X^{\text{ad}} \quad (4.1)$$

where  $J : Q \times U \rightarrow \mathbb{R}$  is the cost functional,  $E : Q \times U \rightarrow Z$  is the so called state equation, and  $X^{\text{ad}}$  collects additional constraints on *control*  $q$  and *state*  $u$ . Assuming, that for each  $q \in Q$  the state equation  $E(q,u) = 0$  admits a unique solution  $u = u(q)$ , with the thus defined *solution operator*  $q \mapsto u(q)$ , one sees, that (4.1) is equivalent to

$$\min_q j(q) := J(q, u(q)) \quad \text{s.t.} \quad q \in Q^{\text{ad}} := \{q \in Q \mid (q, u(q)) \in X^{\text{ad}}\}. \quad (4.2)$$

When deriving optimality conditions for this problem, one needs to calculate, e.g., the derivative  $j'$ . By the chain rule, this can be calculated as

$$j'(q) = J'_q(q, u(q)) + J'_u(q, u(q)) \circ u'(q).$$

However,  $u'$  is only implicitly defined, and moreover for computations it may be inconvenient to evaluate this formula for multiple directions. Hence we will discuss how to obtain an alternative so called *adjoint* formula for the derivative which is very efficient in computations.

### 4.1. Adjoint Approach for the First Derivative

To derive a formula for the involved derivatives, we define the Lagrangian  $\mathcal{L} : Q \times U \times Z^* \rightarrow \mathbb{R}$  by

$$\mathcal{L}(q, u, z) = J(q, u) - \langle E(q, u), z \rangle.$$

By definition of  $u(q)$ , it holds

$$j(q) = \mathcal{L}(q, u(q), z)$$

We now let  $z = z(q)$  depend on  $q$  and calculate

$$\begin{aligned} j'(q) &= \mathcal{L}'_q(q, u(q), z(q)) + \mathcal{L}'_u(q, u(q), z(q)) \circ u'(q) + \mathcal{L}'_z(q, u(q), z(q)) \circ z'(q) \\ &= \mathcal{L}'_q(q, u(q), z(q)) + \mathcal{L}'_u(q, u(q), z(q)) \circ u'(q) \end{aligned}$$

since due to the choice  $u = u(q)$  it holds  $\mathcal{L}'_z(q, u(q), z(q)) = E(q, u) = 0$ . Following this idea, we define  $z = z(q)$  to solve

$$\mathcal{L}'_u(q, u(q), z(q)) = 0$$

#### 4. Adjoint Approach

the so called *adjoint equation*. In more detail,  $z = z(q)$  is the solution of the equation

$$E'_u(q, u(q))^* z = J'_u(q, u(q)).$$

We have thus shown, that for this choice of  $z(q)$  it is

$$j'(q) = \mathcal{L}'_q(q, u(q), z(q)) = J'_q(q, u(q)) - E'_q(q, u(q))^* z(q).$$

**Remark 4.1.1.** We notice, that the first order optimality conditions for (4.1) with  $X^{\text{ad}} = X$  (given (RCQ)) are (notice that the sign of  $z$  can be chosen arbitrary!)

$$\begin{aligned}\mathcal{L}'_u(q, u, z) &= J'_u(q, u) - E'_u(q, u)^* z = 0, \\ \mathcal{L}'_q(q, u, z) &= J'_q(q, u) - E'_q(q, u)^* z = 0, \\ E(q, u) &= 0.\end{aligned}$$

This is exactly the optimality condition

$$j'(q) = 0$$

for the reduced problem (4.2) noting that unique solvability of the adjoint equation implies (RCQ), compare Example 3.3.9.

In summary, to obtain the derivative  $j'(q)$  one has to proceed as follows

1. Solve the state equation

$$E(q, u) = 0$$

to get  $u = u(q) \in U$ .

2. Solve the adjoint equation

$$E'_u(q, u)^* z = J'_u(q, u)$$

to get  $z = z(q) \in Z^*$ .

3. Evaluate

$$j'(q) = J'_q(q, u) - E'_q(q, u)^* z.$$

**Remark 4.1.2.** Notice, that the sign of the multiplier  $z$  is arbitrary, i.e., we could also have used the Lagrangian

$$J(q, u) + \langle E(q, u), z \rangle.$$

to obtain  $z$  with the opposite sign.

**Example 4.1.3.** To clarify this, let us consider the example

$$\begin{aligned} \min J(q, u) &= \frac{1}{2} \|u - u^d\|^2 + \frac{\alpha}{2} \|q\|^2 \\ \text{s.t. } \begin{cases} (\nabla u, \nabla \varphi) = (q, \varphi) & \forall \varphi \in H_0^1(\Omega) \\ u \in H_0^1(\Omega) & q \in L^2(\Omega) \end{cases} \end{aligned}$$

on a domain  $\Omega \subset \mathbb{R}^n$  and  $\alpha > 0$ . It follows from Riesz-representation theorem [A.3.7](#) (or Lax-Milgram [A.3.8](#)), that the equation admits for any given  $q \in L^2(\Omega)$  a unique solution  $u = u(q) \in H_0^1(\Omega) \subset L^2(\Omega)$ . This defines the control-to-state map  $S: L^2(\Omega) \rightarrow L^2(\Omega)$ , and the reduced problem

$$\min_{q \in L^2(\Omega)} j(q) := J(q, Sq).$$

Then, by direct calculation, we get the derivative in any direction  $\delta q \in L^2(\Omega)$  as

$$j'(q)\delta q = (Sq - u^d, S\delta q) + \alpha(q, \delta q).$$

Hence to calculate the  $L^2$ -gradient  $\nabla_{L^2} j(q)$  we would need to compute for a Schauder-basis  $\delta q_i$ ,  $i \in \mathbb{N}$  the coefficients

$$(\nabla_{L^2} j(q), \delta q_i) = (Sq - u^d, S\delta q_i) + \alpha(q, \delta q_i).$$

where for each index  $i$  an application of the operator  $S$  (one PDE solve) is needed. With a simple trick this can be avoided, noting that

$$(\nabla_{L^2} j(q), \delta q_i) = (S^*(Sq - u^d), \delta q_i) + \alpha(q, \delta q_i).$$

can easily be solved to get  $\nabla_{L^2} j(q) = S^*(Sq - u^d) + \alpha q$  only needing two PDE solves, one for  $S$  and one for  $S^*$ . Indeed in the above formalism, we have the state equation

$$u = Sq \quad \text{solving} \quad 0 = \mathcal{L}'_z(q, u, z) = -(\nabla u, \nabla \cdot) + (q, \cdot) \in H_0^1(\Omega)^*,$$

the adjoint equation

$$z = S^*(u - u^d) \quad \text{solving} \quad 0 = \mathcal{L}'_u(q, u, z) = -(\nabla \cdot, \nabla z) + (u - u^d, \cdot) \in H_0^1(\Omega)^*,$$

and the gradient representation

$$\nabla_{L^2} j(q) = \nabla_{q, L^2} \mathcal{L}(q, u, z) = z + \alpha q.$$

#### 4. Adjoint Approach

### 4.2. Adjoint Formulas for the Second Derivatives

Similarly, adjoint formulas for the second derivatives can be obtained (assuming that all involved mappings are twice differentiable). In this case one is not interested in obtaining a formula for the Hessian  $j''(q) \in \mathcal{L}(Q, \mathcal{L}(Q; \mathbb{R}))$  but rather in the product ‘Hessian times direction’  $j''(q)\delta q \in \mathcal{L}(Q; \mathbb{R})$  which is needed, e.g., in the iterative solution of the ‘Newton equation’  $j''(q)\delta q = -j'(q)$ .

To this end, we look at the directional derivatives. We obtain, for two directions  $\delta q, \tau q$  and the abbreviations  $\delta u = u'(q)\delta q$ ,  $\tau u = u'(q)\tau q$ ,  $\delta \tau u = u''(q)(\tau q, \delta q)$  and the same for  $\delta z = z'(q)\delta q$ ,  $\tau z = z'(q)\tau q$ , and  $\delta \tau z = z''(q)(\tau q, \delta q)$ , the following representation for the hessian

$$\begin{aligned} j''(q)(\delta q, \tau q) &= \mathcal{L}_{qq}''(x)(\delta q, \tau q) + \mathcal{L}_{qu}''(x)(\delta q, \tau u) + \mathcal{L}_{qz}''(x)(\delta q, \tau z) \\ &\quad + \mathcal{L}_{uq}''(x)(\delta u, \tau q) + \mathcal{L}_{uu}''(x)(\delta u, \tau u) + \mathcal{L}_{uz}''(x)(\delta u, \tau z) \\ &\quad + \mathcal{L}_{zq}''(x)(\delta z, \tau q) + \mathcal{L}_{zu}''(x)(\delta z, \tau u) + \mathcal{L}_{zz}''(x)(\delta z, \tau z) \\ &\quad + \mathcal{L}'_u(x)(\delta \tau u) + \mathcal{L}'_z(x)(\delta \tau z). \end{aligned}$$

By definition of  $u(q)$  and  $z(q)$  the last two terms vanish. Further the term  $\mathcal{L}_{zz}'' = 0$  since  $\mathcal{L}$  is linear in  $z$ . Now, we can group these terms in order to determine a *tangent equation* to define  $\delta u \in U$  such that

$$\mathcal{L}_{qz}''(x)(\delta q, \varphi) + \mathcal{L}_{uz}''(x)(\delta u, \varphi) = 0 \quad \forall \varphi \in U.$$

as well as a *dual-for-hessian* problem to determine  $\delta z \in U$  such that

$$\mathcal{L}_{qu}''(x)(\delta q, \varphi) + \mathcal{L}_{uu}''(x)(\delta u, \varphi) + \mathcal{L}_{zu}''(x)(\delta z, \varphi) = 0 \quad \forall \varphi \in U.$$

In operator notation, we can calculate the directional derivative  $j''(q)\delta q$  of  $j'(q)$  as follows

1. Solve the state equation

$$E(q, u) = 0$$

to get  $u = u(q) \in U$ .

2. Solve the adjoint equation

$$E'_u(q, u)^* z = J'_u(q, u)$$

to get  $z = z(q) \in Z^*$ .

3. Solve the tangent equation

$$E'_u(q, u)\delta u = -E'_q(q, u)\delta q$$

to get  $\delta u \in U$ .



#### 4.2. Adjoint Formulas for the Second Derivatives

4. Solve the auxiliary adjoint (dual-for-hessian) equation

$$E'_u(q, u)^* \delta z = \mathcal{L}''_{uu}(q, u, z) \delta u + \mathcal{L}_{qu}(q, u, z) \delta q$$

to get  $\delta z \in Z^*$ .

5. Evaluate

$$j''(q) \delta q = \mathcal{L}''_{qq}(q, u, z) \delta q + \mathcal{L}''_{uq}(q, u, z) \delta u - E'_q(q, u)^* \delta z.$$

#### *4. Adjoint Approach*

## 5. Methods

### 5.1. Unconstrained Case

We start our discussion, with the unconstrained problem

$$\min_{x \in X} f(x) \tag{5.1}$$

on a Hilbert spaces  $X$  with continuously differentiable  $f : X \rightarrow \mathbb{R}$ .

We will start with a few definitions

**Definition 5.1.1.** We say that a method converges globally, if for any initial value  $x^0$  the sequence  $x^k$  converges to a stationary point  $\bar{x}$  of (5.1), i.e., a point such that  $f'(\bar{x}) = 0$ .

**Definition 5.1.2.** We say that the convergence  $x^k \rightarrow \bar{x}$  is *linear* with rate  $c$  if

$$\|x^{k+1} - \bar{x}\|_X \leq c \|x^k - \bar{x}\|_X.$$

It converges *super-linear* if

$$\|x^{k+1} - \bar{x}\|_X \leq c_k \|x^k - \bar{x}\|_X.$$

with  $c_k \rightarrow 0$ . It converges with rate  $\alpha$  if for some constant  $c$

$$\|x^{k+1} - \bar{x}\|_X \leq c \|x^k - \bar{x}\|_X^\alpha.$$

There is now a variety of algorithm available, we will only consider the simplest case – the generalization of the well-known steepest descent method.

#### 5.1.1. Gradient descent

As a first (and most simple) choice for the direction one can consider the method of steepest descent. To this end we recall that we can define the gradient of  $f$  at  $x^k$ , using Riesz representation Theorem A.3.7, as follows. Find  $\nabla f(x^k) \in X$  that satisfies

$$(\nabla f(x^k), d)_X = f'(x^k)d \quad \forall d \in X. \tag{5.2}$$

## 5. Methods

It is then immediately clear that  $d^k = -\nabla f(x^k)$  is a descent direction

$$f'(x^k)d^k = (\nabla f(x^k), d^k)_X = -\|\nabla f(x^k)\|_X^2 \leq 0.$$

Further if  $f'(x^k)d^k = 0$  then  $\nabla f(x^k) = 0$  and thus the necessary conditions for a minimizer are fulfilled.

As in the finite dimensional case, a suitable step-length selection is needed, and we choose Armijo-Backtracking, i.e., for  $\beta \in (0, 1)$  and  $\gamma \in (0, 1/2)$ , we let  $l \in \mathbb{N}_0$  be minimal such that for  $t_k := \beta^l$  the inequality

$$f(x^k + t_k d^k) \leq f(x^k) + t_k \gamma f'(x^k)(d^k) \quad (.A)$$

is satisfied. Indeed such an index  $l$ , and thus  $t_k$  exists for any continuously differentiable  $f$  (proof as in finite dimensions).

With this, we obtain the gradient descent Algorithm 5.1.

**Algorithm 5.1** (Gradient descend). Choose  $x^0 \in X$  and let  $k = 0$ .  
 Calculate  $d_0 = -\nabla f(x^0)$   
**while**  $\nabla f(x^k) \neq 0$  **do**  
   Determine  $t_k$  via (.A) with Backtracking.  
    $x^{k+1} = x^k + t_k d^k$   
    $d^{k+1} = -\nabla f(x^{k+1})$   
    $k \leftarrow k + 1$   
**end while**

**Remark 5.1.3.** 1. Clearly, in practical computations, the Algorithm 5.1 needs to be run with a more suitable stopping criterion. Moreover, as an evaluation of  $\nabla f$  is typically impossible inexact variants of the algorithm need to be considered.

We can now show the global convergence of the gradient descent Algorithm 5.1 similar to the finite dimensional case.

**Theorem 5.1.4.** For any initial value  $x^0$  Algorithm 5.1 either terminates finitely with  $\nabla f(x^k) = 0$  or it creates an infinite sequence  $x^k$  such that  $f(x^{k+1}) < f(x^k)$  and each accumulation point  $\bar{x}$  is a stationary point of  $f$ , i.e.,  $\nabla f(\bar{x}) = 0$ .

*Beweis.* If the algorithm terminates after finitely many iterations there is nothing to show.

Otherwise, we have an infinite sequence  $x^k \in X$ . Let  $(x^k)_{k \in \mathcal{K}}$  be a convergent subsequence with limit  $\bar{x}$ . By construction  $f(x^k)$  is monotone decreasing, see, (.A) together with the definition (5.2) and  $d^k = -\nabla f(x^k)$ . We obtain

$$\lim_{\mathcal{K} \ni k \rightarrow \infty} f(x^k) = \lim_{k \rightarrow \infty} f(x^k) = f(\bar{x}).$$

Utilizing telescope sums, we obtain from the Armijo condition (A)

$$\begin{aligned}
-\infty &< f(\bar{x}) - f(x^0) \\
&= \lim_{k \rightarrow \infty} f(x^k) - f(x^0) \\
&= \sum_{k=0}^{\infty} (f(x^{k+1}) - f(x^k)) \\
&= \sum_{k=0}^{\infty} (f(x^k + t_k d^k) - f(x^k)) \\
&\leq \sum_{k=0}^{\infty} \gamma t_k f'(x^k)(d^k) \\
&= -\gamma \sum_{k=0}^{\infty} t_k \|\nabla f(x^k)\|_X^2 \\
&\leq 0
\end{aligned}$$

and consequently

$$t_k \|\nabla f(x^k)\|_X^2 \rightarrow 0.$$

Now, if  $\nabla f(x^k) \rightarrow 0$  on  $k \in \mathcal{K}$  nothing is to show. Otherwise, possibly taking again a subsequence,  $\lim_{\mathcal{K} \ni k \rightarrow \infty} t_k = 0$ , w.l.o.g  $t_k < 1$  for all  $k \in \mathcal{K}$ . Since  $l$  is minimal such that  $t_k = \beta^l$  satisfies (A), we know

$$f(x^k + \beta^{-1} t_k x^k) - f(x^k) > \gamma \beta^{-1} t_k f'(x^k)(d^k)$$

and since  $\beta^{-1} t_k \rightarrow 0$  as well, we conclude by the one-dimensional mean-value theorem the existence of  $\theta_k \in (0, 1)$  such that

$$\begin{aligned}
-\|\nabla f(\bar{x})\|_X^2 &= \lim_{\mathcal{K} \ni k \rightarrow \infty} -\|\nabla f(x^k)\|_X^2 \\
&< -\gamma \lim_{\mathcal{K} \ni k \rightarrow \infty} \|\nabla f(x^k)\|_X^2 \\
&= \gamma \lim_{\mathcal{K} \ni k \rightarrow \infty} f'(x^k)(d^k) \\
&\leq \lim_{\mathcal{K} \ni k \rightarrow \infty} \frac{f(x^k + \beta^{-1} t_k d^k) - f(x^k)}{\beta^{-1} t_k} \\
&= \lim_{\mathcal{K} \ni k \rightarrow \infty} f'(x^k + \theta_k \beta^{-1} t_k d^k)(d^k) \\
&= \lim_{\mathcal{K} \ni k \rightarrow \infty} -f'(x^k + \theta_k \beta^{-1} t_k d^k)(\nabla f(x^k)) \\
&= -f'(\bar{x})(\nabla f(\bar{x})) \\
&= -\|\nabla f(\bar{x})\|^2.
\end{aligned}$$

By the strict inequality, we get a contradiction, implying that  $\|\nabla f(\bar{x})\| = 0$ . □

## 5. Methods

**Remark 5.1.5.** 1. Notice, that although the theorem looks identical to the finite dimensional case it is much weaker in that it is less clear when convergent subsequences occur. This is due to the fact that for an infinite dimensional space bounded and closed sets are no longer compact. Hence although  $\alpha\|x^k\|^2 \leq f(x^k) < f(x^0)$  implies boundedness of  $x^k$  this sequence does not necessarily contain a convergent subsequence.

2. However, strong convergence can be asserted if the derivative is monotone in the sense that there exists an  $\varepsilon > 0$  such that

$$(\nabla f(x) - \nabla f(p), x - p)_X \geq \varepsilon \|x - p\|_X^2 \quad \forall x, p \in X.$$

Such conditions are fulfilled, e.g., if  $\widehat{S}$  is linear, i.e.,  $\widehat{S} = S' = S$ , or near stationary points that satisfy the second order sufficient conditions  $f''(\bar{x})(d, d) > \varepsilon \|d\|_X^2$  (different  $\varepsilon$ ).

We will now see that, indeed, the last remark is true.

**Theorem 5.1.6.** Let  $f : X \rightarrow \mathbb{R}$  be continuously differentiable and assume that the set

$$N := \{x \in X \mid f(x) < f(x^0)\}$$

is bounded and nonempty. Further let  $\nabla f(x)$  be uniformly continuous on  $\text{conv}(N)$  and monotone in the sense, that there exists  $\varepsilon > 0$  such that

$$(\nabla f(x) - \nabla f(p), x - p)_X \geq \varepsilon \|x - p\|_X^2$$

is true for all  $x, p \in X$ . Then, Algorithm 5.1 either terminates after finitely many iterations, or it generates a sequence  $x^k$  converging to the only stationary point  $\bar{x} \in N \subset X$

*Beweis.* If the algorithm terminates after finitely many iterations there is nothing to show. Otherwise, we first show that there is at most one stationary point on  $N$ . To see this, let  $\bar{x}_1, \bar{x}_2 \in X$  be stationary point, i.e.,  $\nabla f(\bar{x}_i) = 0$ . Then by monotonicity it holds

$$0 = (\nabla f(\bar{x}_1) - \nabla f(\bar{x}_2), \bar{x}_1 - \bar{x}_2)_X \geq \varepsilon \|\bar{x}_1 - \bar{x}_2\|^2,$$

and thus  $\bar{x}_1 = \bar{x}_2$ .

Now, we need to show, that the sequence generated by the algorithm converges. We note, that since  $\nabla f(x)$  is uniformly continuous on  $N$  there exists a constant  $C$  such that  $\|\nabla f(x)\|_X \leq C$  on the bounded set  $\text{conv}(N)$ ; this is particularly true for the directions  $d^k = -\nabla f(x^k)$ . By the mean-value Theorem 3.1.6  $f$  is bounded on  $\text{conv}(N)$ . Now, analogously to the proof of

### 5.1. Unconstrained Case

Theorem 5.1.4, we assert that the sequence of function values converges (to  $\bar{f}$ ) and thus

$$-\infty < \bar{f} - f(x^0) = -\gamma \sum_{k=0}^{\infty} t_k \|\nabla f(x^k)\|_X^2 \leq 0.$$

and consequently  $t_k \|\nabla f(x^k)\|_X^2 \rightarrow 0$ . If we would know that  $x^k$  converges we would be done now, but we don't! Now, we show  $\lim_{k \rightarrow \infty} \|d^k\|_X = 0$ . To this end, assume for contradiction that this is not the case, i.e., there exists a subsequence  $\mathcal{K}$  and  $\varepsilon > 0$  such that  $\|d^k\|_X \geq \varepsilon$  for all  $k \in \mathcal{K}$ . By the above convergence  $t_k \rightarrow 0$  for  $\mathcal{K} \ni k \rightarrow \infty$  and w.l.o.g  $t_k < 1$ . Now, for  $k \in \mathcal{K}$ , by our Armijo-Backtracking, we know that  $\beta^{-1} t_k \leq 1$  does not satisfy (A), i.e.,

$$f(x^k + \beta^{-1} t_k d^k) - f(x^k) > \beta^{-1} t_k \gamma f'(x^k) d^k$$

and thus by the mean-value theorem for the mapping  $t \mapsto x^k + t_k d^k$  there exists a sequence  $\theta_k \in (0, 1)$  such that

$$\begin{aligned} -\|d^k\|_X^2 &= f'(x^k) d^k \\ &< \gamma f'(x^k) d^k \\ &\leq \frac{f(x^k + \beta^{-1} t_k d^k) - f(x^k)}{\beta^{-1} t_k} \\ &= f'(x^k + \theta_k \beta^{-1} t_k d^k) d^k \\ &= f'(x^k) d^k + (f'(x^k + \theta_k \beta^{-1} t_k d^k) - f'(x^k)) d^k. \end{aligned}$$

Since  $\theta_k \beta^{-1} t_k d^k \rightarrow 0$  uniform-continuity and boundedness of  $d^k$  imply that the second summand converges to zero and we obtain

$$-\liminf_{\mathcal{K} \ni k} \|d^k\|_X^2 = \liminf_{\mathcal{K} \ni k} f'(x^k) d^k < \gamma \liminf_{\mathcal{K} \ni k} f'(x^k) d^k \leq \liminf_{\mathcal{K} \ni k} f'(x^k) d^k + 0 = -\liminf_{\mathcal{K} \ni k} \|d^k\|_X^2.$$

Since  $\liminf_{\mathcal{K} \ni k} \|d^k\|_X \neq 0$  this is a contradiction and thus  $\|d^k\|_X \rightarrow 0$ .

Now, monotonicity gives for any  $k, l$  that

$$\varepsilon \|x^k - x^l\|_X^2 \leq (\nabla f(x^k) - \nabla f(x^l), x^k - x^l)_X \leq \|\nabla f(x^k) - \nabla f(x^l)\|_X \|x^k - x^l\|_X$$

and thus

$$\varepsilon \|x^k - x^l\|_X \leq 2 \max(\|\nabla f(x^k)\|_X, \|\nabla f(x^l)\|_X) \rightarrow 0 \quad (k, l \rightarrow \infty).$$

Consequently, the sequence  $x^k$  is converging to some limit  $\bar{x}$  (with  $\nabla f(\bar{x}) = 0$ ), and it is

$$f(\bar{x}) = \lim_{k \rightarrow \infty} f(x^k) \leq f(x^1) < f(x^0)$$

and thus  $\bar{x} \in N$ . □

Analogous to the finite dimensional case, the above algorithm is not very fast. To see this, we recall the following result:

## 5. Methods

**Lemma 5.1.7.** Let  $f(x) = (Hx, x)_X + (b, x)_X$ , for a symmetric positive definite, and bounded operator  $H \in \mathcal{L}(X, X)$  and  $x \in X$ . Assume that the spectrum  $\sigma(H)$  satisfies  $\sigma(H) \subset [\lambda, \Lambda]$  for some constants  $0 < \lambda \leq \Lambda < \infty$ . We define the condition number

$$\kappa = \frac{\Lambda}{\lambda}.$$

Then the gradient descend method with exact line-search converges linearly in the norm induced by  $H$  with rate  $\frac{(\kappa-1)}{(\kappa+1)}$ , e.g.,

$$\|x^{k+1} - \bar{x}\|_H \leq \frac{(\kappa-1)}{(\kappa+1)} \|x^k - \bar{x}\|_H$$

where

$$\|x\|_H^2 = (Hx, x)_X.$$

*Beweis.* The proof in finite dimensions can be found in any introductory script to numerical analysis. For the infinite dimensional case, see, e.g., [Daniel \[1967\]](#).  $\square$

### 5.1.2. Newton Methods

In a next step, we consider the second derivative of our functional; of course assuming that  $f$  is twice-continuously differentiable.

Analogously to the gradient, we can define a Hessian-operator  $H = H(x) = \nabla^2 f(x) \in \mathcal{L}(X, X)$  by the relation

$$(Hd_1, d_2)_X := f''(x)[d_1, d_2].$$

Near a point  $\bar{x}$  satisfying the second order sufficient condition

$$\nabla f(\bar{x}) = 0, \quad (H(\bar{x})d, d)_X \geq \varepsilon \|d\|_X^2 \quad \forall d \in X$$

one can assume that  $H$  remains positive definite. This is analogous to the finite dimensional case, however the gap between the necessary condition

$$(H(\bar{x})d, d)_X \geq 0 \quad \forall d \in X$$

is larger in the infinite-dimensional case.

As in the finite-dimensional case, the equation  $f'(\bar{x}) = 0$  and Taylor-expansion

$$-\nabla f(x) = \nabla f(\bar{x}) - \nabla f(x) \approx \nabla^2 f(x)(\bar{x} - x)$$

motivates the local Newton-iteration



**Algorithm 5.2** (Local Newton method). Choose  $x^0 \in X$  and let  $k = 0$ .  
**while**  $\nabla f(x^k) \neq 0$  **do**  
    Calculate  $d_k$  satisfying  $\nabla^2 f(x^k)d_k = -\nabla f(x^k)$   
     $x^{k+1} = x^k + d^k$   
     $k \leftarrow k + 1$   
**end while**

### Convergence of the local Newton-method

Assuming that the Functional  $f$  is twice-continuously differentiable, we obtain the usual local convergence theorem:

**Theorem 5.1.8.** Let  $f : X \rightarrow \mathbb{R}$  be twice continuously Fréchet-differentiable and  $\bar{x} \in X$  be a local minimizer satisfying the second order sufficient condition

$$f''(\bar{x})[d, d] \geq \mu \|d\|_X^2 \quad \forall d \in X.$$

Then there exists  $\delta > 0$  such that

1. For any  $x \in B_\delta(\bar{x})$  it is

$$f''(x)[d, d] \geq \frac{\mu}{2} \|d\|_X^2 \quad \forall d \in X.$$

2.  $\bar{x}$  is the only stationary point on  $B_\delta(\bar{x}) := \{x \in X \mid \|x - \bar{x}\|_X < \delta\}$ .
3. For any  $x^0 \in B_\delta(\bar{x})$  the local Newton-method (with  $\text{TOL} = 0$ ) is well-defined and either terminates with  $x^k = \bar{x}$  or creates a sequence that converges  $q$ -superlinearly towards  $\bar{x}$ .

*Beweis.*

1. This is clear since  $f'' \in \mathcal{L}(X; \mathcal{L}(X; \mathbb{R}))$  is continuous.
2. Assume that there is a second stationary point  $\hat{x} \in B_\delta(\bar{x})$ . Then with  $d = \hat{x} - \bar{x}$  it holds

$$f'(\hat{x})d - f'(\bar{x})d = f''(\bar{x} + \xi d)[d, d]$$

for some  $\xi \in (0, 1)$  by the mean value theorem for the mapping  $t \mapsto f'(\bar{x} + td)d$  (note that  $\bar{x} + td \in B_\delta(\bar{x})$  by convexity of the norm). By the lower bound from 1. it is thus

$$0 = f'(\hat{x})d - f'(\bar{x})d \geq \frac{\mu}{2} \|d\|_X^2$$

and thus  $\hat{x} = \bar{x}$ .

## 5. Methods

3. Now, we need to see that the iteration is well-defined. To this end, let  $x^k \in B_\delta(\bar{x})$  be given and let  $H_k = H(x^k)$  be the Hessian-operator. Then the Newton-step  $H_k d^k = -\nabla f(x^k)$  has a unique solution  $d^k$ . This follows directly from the fact, that due to 1. the functional

$$d \mapsto f_k(d) := \frac{1}{2}(H_k d, d)_X + (\nabla f(x^k), d)_X$$

is monotone in the sense of Theorem 5.1.6 (and strictly convex) and thus there is a unique minimizer  $d^k$  equivalently given by the Newton-equation.

It remains to show that the sequence  $x^k$  remains in the ball  $B_\delta(\bar{x})$  and converges superlinearly towards  $\bar{x}$ . To this end, we notice that

$$\begin{aligned} \frac{\mu}{2} \|x^{k+1} - \bar{x}\|_X^2 &\leq f''(x^k)[x^{k+1} - \bar{x}, x^{k+1} - \bar{x}] \\ &= (H_k(x^{k+1} - \bar{x}), x^{k+1} - \bar{x})_X \\ &= (H_k(x^k + d^k - \bar{x}), x^{k+1} - \bar{x})_X \\ &= (H_k(x^k - \bar{x}), x^{k+1} - \bar{x})_X - (\nabla f(x^k), x^{k+1} - \bar{x})_X \\ &= (H_k(x^k - \bar{x}), x^{k+1} - \bar{x})_X + (\nabla f(\bar{x}) - \nabla f(x^k), x^{k+1} - \bar{x})_X \\ &\leq \|\nabla f(\bar{x}) - \nabla f(x^k) + H_k(x^k - \bar{x})\|_X \|x^{k+1} - \bar{x}\|_X. \end{aligned}$$

Consequently, noting that  $H_k = f''(x^k)$ , we get

$$\|x^{k+1} - \bar{x}\|_X \leq \frac{2}{\mu} \sup_{x \in B_{\|x^k - \bar{x}\|_X}(\bar{x})} \|f''(x) - H_k\|_{\mathcal{L}(X; \mathcal{L}(X; \mathbb{R}))} \|x^k - \bar{x}\|_X$$

by the mean-value Theorem 3.1.6 applied to the mapping  $x \mapsto f'(x) - H_k x$ , see also [Ciarlet, 2013, Theorem 7.2.2]. By continuity of  $f''$  (and the definition of  $H_k$ ) there exists some  $\delta$ , such that

$$\sup_{x \in B_\delta(\bar{x})} \|f''(x) - H_k\|_{\mathcal{L}(X; \mathcal{L}(X; \mathbb{R}))} < \frac{\mu}{2}$$

and thus linear-convergence of  $x^k \rightarrow \bar{x}$  follows if  $x^0 \in B_\delta(\bar{x})$  for sufficiently small  $\delta$ . Superlinear convergence is immediate, since indeed, from linear convergence  $\|x^k - \bar{x}\|_X \rightarrow 0$  it follows that

$$\sup_{x \in B_{\|x^k - \bar{x}\|_X}(\bar{x})} \|f''(x) - H_k\|_{\mathcal{L}(X; \mathcal{L}(X; \mathbb{R}))} \rightarrow 0.$$

□

### Remark 5.1.9.

1. As in finite dimensions for Lipschitz-continuous second derivatives, quadratic convergence follows.

2. Indeed, as in finite dimensions, inexact solution of the linear-system are ok, as long as the inexactness tends to zero sufficiently fast. (Exercise!)

## 5.2. Constrained Case

In the constrained case, one can generalize first order methods such as the Gradient descent algorithm 5.1 to work for constrained problems. This can be done for instance by replacing the search directions  $x^{k+1} = x^k + t_k d^k$  with the curve

$$x^{k+1} = P_{X^{\text{ad}}}(x^k + t_k d^k)$$

where  $P_{X^{\text{ad}}}$  denotes a suitable projection onto the feasible set. However such methods will not be any faster than the Gradient descent algorithm 5.1 and moreover requires the evaluation of a projection in each step. Further, any such problem can be approximated by suitable penalty or barrier methods, see, e.g., [Schiela and Wollner \[2011\]](#), [Schiela \[2006\]](#) for barrier methods or [Hintermüller et al. \[2014\]](#), [Ulbrich \[2011\]](#) for penalty approaches.

Instead of these approaches, we will discuss a generalization of Newton's method which is capable to give at least superlinear convergence for many constrained optimization problems.

### 5.2.1. Reformulation of the KKT-conditions as a nonsmooth equation

**Definition 5.2.1.** Let  $K \subset Z$  be a closed convex set. A function  $\Psi: Z \times Z^* \rightarrow \mathbb{R}$  satisfying

$$\Psi(z, \lambda) = 0 \quad \Leftrightarrow \quad z \in K, \lambda \in T(K, z)^\circ$$

is called a *complementarity function*.

It is clear, that any such function  $\Psi$  would allow us to reformulate the KKT-conditions as an equation. That in fact such a complementarity function exists can be seen by the following

**Lemma 5.2.2.** Let  $Z \simeq Z^*$  be a Hilbert space,  $K \subset Z$  be closed and convex. Then for any  $\sigma > 0$  the following is equivalent for  $z, \lambda \in Z$

1.  $z \in K, \lambda \in T(K, z)^\circ$
2.  $z = P_K(z + \sigma \lambda)$

where  $P_K: Z \rightarrow Z$  is the projection onto  $K$  given by

$$P_K(v) = \operatorname{argmin}_{k \in K} \|k - v\|_Z.$$

## 5. Methods

*Beweis.* By definition, 2. means

$$\min_{k \in K} \frac{1}{2} \|k - (z + \sigma \lambda)\|_Z^2 := \min_{k \in K} \varphi(k) = \varphi(z).$$

Now,  $\varphi'(k) = k - (z + \sigma \lambda)$ , and thus  $\varphi'(z) = -\sigma \lambda$  and the first order necessary conditions from Theorem 3.2.3 show

$$z \in K, \quad (-\sigma \lambda, d) \geq 0 \quad \forall d \in T(K, z)$$

which is 1. Since  $\varphi$  is convex the conditions are also sufficient and hence 1. implies 2., too.  $\square$

We have thus seen that to find KKT-pairs  $(\bar{x}, \bar{\lambda}) \in X \times Z$ , where  $Z$  is a Hilbert space, we can equivalently search for solutions of the system

$$\begin{pmatrix} f'(x) + G'(x)^* \lambda \\ G(x) - P_K(G(x) + \sigma \lambda) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

where  $\sigma > 0$  is arbitrary.

To this end, we would like to use a Newton method. Unfortunately, while we may hope that  $f, G$  are twice Fréchet-differentiable, and thus the first line has one derivative, differentiability of the second line can not be assumed in general. This is shown by the following example

**Example 5.2.3.** Consider  $X = Z = L^2(\Omega)$ ,  $G(x) = x$  and

$$K = \{x \in L^2(\Omega) \mid a \leq x \leq b\}$$

for some  $a \leq b \in L^\infty(\Omega)$ .

As we have seen in the exercises, for any  $x \in K$  it holds

$$\begin{aligned} T(K, G(x))^\circ &= T(K, x)^\circ \\ &= \left\{ \lambda \in L^2(\Omega) \mid \lambda \begin{cases} \leq 0 & \text{on } \{x = a\}, \\ \geq 0 & \text{on } \{x = b\}, \\ = 0 & \text{on } \{a < x < b\}. \end{cases} \right\} \end{aligned}$$

The projection onto the admissible set is given by

$$P_K(x) = \min(\max(a, x), b),$$

the pointwise projection onto the interval given by  $a, b$ . Now, since  $P_K$  is induced by a non-differentiable mapping  $\mathbb{R} \rightarrow \mathbb{R}$  it is clear, that it is not differentiable.

Hence, we need to see whether we can generalize the Newton-method to less regular mappings.

## 5.2.2. Generalized Newton Methods

**Definition 5.2.4.** Let  $F: \text{dom}(F) \subset X \rightarrow Y$  be a given mapping between two Banach spaces. We call  $F$  *slant differentiable* at  $x \in X$  if there exists an open neighborhood  $\mathcal{N}(x)$  and a mapping  $G: \mathcal{N}(x) \rightarrow \mathcal{L}(X, Y)$  such that

$$\lim_{\|d\|_X \downarrow 0} \frac{\|F(x+d) - F(x) - G(x+d)d\|_Y}{\|d\|_X} = 0.$$

Before we continue with our analysis we will state the first result which establishes the local fast convergence of the generalized Newton method.

**Theorem 5.2.5.** Let  $\bar{x}$  be a solution to  $F(x) = 0$ . Assume that  $F$  is slant differentiable with derivative  $G$ . If  $G$  is non singular and  $\|G(x)^{-1}\|_{\mathcal{L}(Y, X)} \leq M$  on  $\mathcal{N}(\bar{x})$  then the Newton iteration

$$x_{k+1} = x_k - G(x_k)^{-1}F(x_k)$$

converges  $q$ -superlinearly to  $\bar{x}$  if the initial value  $x_0$  is close enough to  $\bar{x}$ .

*Beweis.* Let  $r > 0$  be such that  $B_r(\bar{x}) \subset \mathcal{N}(\bar{x})$  then by definition of slant differentiability it holds

$$\|F(x_k) - F(\bar{x}) - G(x_k)(x_k - \bar{x})\|_Y \leq c_r \|x_k - \bar{x}\|_X$$

for all  $x_k \in B_r(\bar{x})$  and a constant  $c_r \rightarrow 0$  as  $r \rightarrow 0$ .

We initialize the iteration such that  $x_0 \in B_{r_0}(\bar{x})$  such that  $c_{r_0} < 1/M$ . Now, by definition of the iteration it holds

$$\begin{aligned} \|x_{k+1} - \bar{x}\|_X &= \|x_k - \bar{x} - G(x_k)^{-1}F(x_k)\|_X \\ &\leq \|G(x_k)^{-1}\|_{\mathcal{L}(Y, X)} \|G(x_k)(x_k - \bar{x}) - F(x_k) + F(\bar{x})\|_Y \\ &\leq c_{\|x_k - \bar{x}\|_X} M \|x_k - \bar{x}\|_X. \end{aligned}$$

By induction it follows that  $c_{\|x_k - \bar{x}\|_X} M \leq c_{r_0} M < 1$ . Thus we have

$$\|x_{k+1} - \bar{x}\|_X \leq M^{k+1} \prod_{l=0}^k c_{\|x_l - \bar{x}\|_X} \leq (M c_{r_0})^{k+1} \rightarrow 0.$$

The super-linear convergence then follows as  $c_{\|x_k - \bar{x}\|_X} \rightarrow 0$  with  $k \rightarrow \infty$ .  $\square$

Before we continue to apply this to our model problem we will start with some examples

## 5. Methods

**Example 5.2.6.** 1. Any bounded linear operator  $A$  is slant differentiable with  $G(x) = A$ .

2. In a Hilbert space  $H$  the mapping  $F(u) = \|u\|_H$  is slant differentiable (exercise).

3. The map  $\max(0, \cdot)$  defined as

$$\max(0, \cdot): L^r(\Omega) \rightarrow L^p(\Omega), \quad u \mapsto \max(0, u)$$

is slant differentiable for all  $1 \leq p < r \leq \infty$ . Its slant derivative  $G(u) \in \mathcal{L}(L^r(\Omega), L^p(\Omega))$  is defined by pointwise multiplication with the function

$$g_u(x) = \begin{cases} 1 & u(x) > 0 \\ \delta & u(x) = 0 \\ 0 & \text{otherwise} \end{cases}$$

where  $\delta \in \mathbb{R}$  is arbitrary. This means

$$(G(u)\delta u)(x) = g_u(x)\delta u(x)$$

We will only show that this is true for the choice  $\delta = 0$ . Then, we can write this as

$$G(u)\delta u = \chi_{\mathcal{A}(u)}\delta u$$

where  $\mathcal{A}(u) = \{x \in \Omega | u(x) > 0\}$ .  $\chi_M$  is called characteristic function of the set  $M$  and is defined as

$$\chi_M = \begin{cases} 1 & x \in M \\ 0 & x \notin M \end{cases}$$

To see slant differentiability let  $u, \delta u \in L^r(\Omega)$  be given. Then we denote the remainder by

$$\omega_{u, \delta u}(x) = \max(0, u(x) + \delta u(x)) - \max(0, u(x)) - g_{u+\delta u}(x)\delta u(x).$$

A simple calculation shows

$$|\omega_{u, \delta u}(x)| \begin{cases} \leq |u(x)| & (u(x) + \delta u(x))u(x) < 0, \\ \leq |u(x)| & (u(x) + \delta u(x)) = 0, u(x) \neq 0 \\ = 0 & \text{otherwise.} \end{cases} \quad (5.3)$$

From (5.3), we obtain that we only need to consider the set

$$\Omega_0^{\delta u} = \{x \in \Omega | u(x) \neq 0, u(x)(u(x) + \delta u(x)) \leq 0\}.$$

Now, there are two 'subsets' to consider. One  $\Omega_\varepsilon^{\delta u}$  where  $u$  is large, compared to  $\varepsilon$ , and thus a sign change can only occur if  $\delta u$  is also large; hence we expect this measure of this set to vanish as  $\delta u \rightarrow 0$ . And another set  $\Omega_0^{\delta u} \setminus \Omega_\varepsilon^{\delta u}$ , where  $u$  is small. On this set we have trouble since it depends on both  $\varepsilon$  and  $\delta u$ , and thus we put it into a fixed set  $\tilde{\Omega}_\varepsilon^u \supset \Omega_0^{\delta u} \setminus \Omega_\varepsilon^{\delta u}$  where we need to show that as  $\varepsilon \rightarrow 0$  this set vanishes.

In more detail, we define the following subsets  $\Omega_\varepsilon \subset \Omega_0$  as follows

$$\Omega_\varepsilon^{\delta u} = \{x \in \Omega \mid |u| \geq \varepsilon, u(x)(u(x) + \delta u(x)) \leq 0\}.$$

Now,  $|u(x)| \geq \varepsilon$  on  $\Omega_\varepsilon^{\delta u}$  and thus to get the second defining inequality we have in addition  $|\delta u(x)| \geq \varepsilon$  on  $\Omega_\varepsilon^{\delta u}$ . Hence it holds ( $r < \infty$ )

$$\|\delta u\|_{L^r(\Omega)} \geq \|\delta u\|_{L^r(\Omega_\varepsilon^{\delta u})} \geq \varepsilon |\Omega_\varepsilon^{\delta u}|^{1/r}.$$

In particular, for any fixed  $\varepsilon > 0$  it holds

$$\lim_{\|\delta u\|_{L^r(\Omega)} \rightarrow 0} |\Omega_\varepsilon^{\delta u}| = 0. \quad (5.4)$$

Which is true for  $r = \infty$  as well.

We continue by defining an other set

$$\tilde{\Omega}_\varepsilon^u = \{x \in \Omega \mid 0 < |u(x)| \leq \varepsilon\} \supset \Omega_0^{\delta u} \setminus \Omega_\varepsilon^{\delta u}.$$

For this it holds

$$\tilde{\Omega}_\varepsilon^u \subset \tilde{\Omega}_{\varepsilon'}^u, \quad 0 < \varepsilon \leq \varepsilon', \quad \text{and} \quad \bigcap_{\varepsilon > 0} \tilde{\Omega}_\varepsilon^u = \emptyset.$$

This implies

$$\lim_{\varepsilon \rightarrow 0} |\tilde{\Omega}_\varepsilon^u| = 0. \quad (5.5)$$

By Hölder's inequality, we get for any function  $f \in L^r(\Omega)$ , and  $1 \leq p < r \leq \infty$ , it holds

$$\|f\|_p \leq \|1\|_s \|f\|_r = |\Omega|^{1/s} \|f\|_r, \quad s = \begin{cases} \frac{pr}{r-p} & r < \infty \\ p & r = \infty. \end{cases}$$

## 5. Methods

Then, we use (5.3) to conclude

$$\begin{aligned}
\frac{\|\omega_{u,\delta u}\|_{L^p}}{\|\delta u\|_{L^r}} &\leq \frac{1}{\|\delta u\|_{L^r}} \|u(x)\|_{L^p(\Omega_0^{\delta u})} \\
&\leq \frac{1}{\|\delta u\|_{L^r}} \left\{ \|u(x)\|_{L^p(\Omega_\varepsilon^{\delta u})} + \|u(x)\|_{L^p(\Omega_0^{\delta u} \setminus \Omega_\varepsilon^{\delta u})} \right\} \\
&\leq \frac{1}{\|\delta u\|_{L^r}} \left\{ |\Omega_\varepsilon^{\delta u}|^{1/s} \|u(x)\|_{L^r(\Omega_\varepsilon^{\delta u})} + |\tilde{\Omega}_\varepsilon^u|^{1/s} \|u(x)\|_{L^r(\Omega_0^{\delta u} \setminus \Omega_\varepsilon^{\delta u})} \right\} \\
&\leq \frac{c}{\|\delta u\|_{L^r}} \|u(x)\|_{L^r(\Omega_0^{\delta u})} \left( |\Omega_\varepsilon^{\delta u}|^{1/s} + |\tilde{\Omega}_\varepsilon^u|^{1/s} \right).
\end{aligned}$$

Now, we note that on  $\Omega_0^{\delta u}$  it holds  $|\delta u| \geq |u|$  and hence

$$\begin{aligned}
\frac{\|\omega_{u,\delta u}\|_{L^p}}{\|\delta u\|_{L^r}} &\leq \frac{c}{\|\delta u\|_{L^r}} \|\delta u\|_{L^r} \left( |\Omega_\varepsilon^{\delta u}|^{1/s} + |\tilde{\Omega}_\varepsilon^u|^{1/s} \right) \\
&\leq c \left( |\Omega_\varepsilon^{\delta u}|^{1/s} + |\tilde{\Omega}_\varepsilon^u|^{1/s} \right).
\end{aligned}$$

Now, due to (5.5), the second summand can be made arbitrarily small independent of  $\delta u$ . Further, due to (5.4), the first summand goes to zero which shows the slant differentiability, e.g.

$$\frac{\|\omega_{u,\delta u}\|_{L^p}}{\|\delta u\|_{L^r}} \rightarrow 0, \quad (\|\delta u\|_{L^r} \rightarrow 0).$$

**Remark 5.2.7.** Finally, we like to note that the restriction to consider  $\max(0, u)$  as a mapping from  $L^r$  into  $L^p$  with  $p < r$  is not just for convenience, but in fact the function

$$g_u(x) \begin{cases} 1 & u(x) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

is not a slant derivative if  $r = p$  (exercise).

Now, we would like to apply our knowledge to the KKT-conditions given at the beginning of this section. But before we can do so we first need the following chain rule

**Theorem 5.2.8.** Let  $f : \text{dom}(f) \subset X \rightarrow Y$  be continuously Fréchet differentiable and  $g : Y \rightarrow Z$  be slant differentiable at  $f(x)$  with derivative  $G$ . Assume that  $\|G\|_{\mathcal{L}(Y,Z)} \leq c$  in  $\mathcal{N}(f(x))$ . Then  $F = g \circ f$  is slant differentiable in  $x$  and the slant derivative is given by  $G(f(x+d))f'(x+d) \in \mathcal{L}(X, Z)$  for  $d$  sufficiently small.



Beweis. □

### 5.2.3. Semi-smoothness

On careful inspection of Theorem 5.2.5, we notice that while a nice convergence property holds if a suitable derivative is chosen, the choice is not local at the moment since we require that  $G$  is a slant derivative at the solution  $\bar{x}$ . To allow for a reasonable selection of  $G(x^k)$ , we need to restrict the set of functions to allow for a selection of  $G(x^k)$  based on local information, only.

This will lead to so called semi-smooth Newton methods, see, e.g., [Ito and Kunisch \[2008\]](#), [Ulbrich \[2011\]](#).

**Definition 5.2.9.** Let  $F : X \rightarrow Y$  be a continuous operator between Banach spaces  $X, Y$ . And let  $\partial F : X \rightrightarrows \mathcal{L}(X, Y)$  be a set-valued mapping with  $\partial F(x) \neq \emptyset$  for all  $x \in X$ . Then  $F$  is called  $\partial F$ -semi-smooth at  $x \in X$ , if

$$\lim_{\|d\|_X \rightarrow 0} \sup_{G \in \partial F(x+d)} \frac{\|F(x+d) - F(x) - Gd\|_Y}{\|d\|_X} = 0.$$

Further,  $F$  is called  $\partial F$ -semi-smooth of order  $\alpha > 0$  at  $x \in X$ , if

$$\limsup_{\|d\|_X \rightarrow 0} \sup_{G \in \partial F(x+d)} \frac{\|F(x+d) - F(x) - Gd\|_Y}{\|d\|_X^{1+\alpha}} < \infty.$$

**Algorithm 5.3** (Local semi-smooth Newton method). Choose  $x^0 \in X$ , and let  $k = 0$ .

```

while  $F(x^k) \neq 0$  do
  Choose  $G_k \in \partial F(x^k)$ 
  Calculate  $d^k$  satisfying  $G_k d^k = -F(x^k)$ 
  Set  $x^{k+1} \leftarrow x^k + d^k$ 
   $k \leftarrow k + 1$ 
end while

```

**Corollary 5.2.10.** Let  $\bar{x}$  be a solution to  $F(x) = 0$ . Assume that  $F$  is  $\partial F$ -semi-smooth, and there exists  $\delta > 0$  and  $M > 0$  such that for any  $x \in B_\delta^X(\bar{x})$  any  $G \in \partial F(x)$  is invertible and satisfies

$$\|G^{-1}\|_{\mathcal{L}(Y, X)} \leq M.$$

Then the semi-smooth Newton method 5.3 converges superlinearly to  $\bar{x}$  if the initial value  $x_0$  is close enough to  $\bar{x}$ .

## 5. Methods

Moreover, if  $F$  is  $\partial F$ -semi-smooth of order  $\alpha > 0$ , then the convergence is of order  $1 + \alpha$ .

*Beweis.* The superlinear convergence is an immediate consequence of Theorem 5.2.5.

The case of smoothness order  $\alpha > 0$  is left as an exercise.  $\square$

Similar to slant-differentiable operators, semi-smooth operators satisfy some useful calculus rules, cf., Ulbrich [2011].

**Theorem 5.2.11.** *Let  $X, Y, Z, X_i$ , and  $Y_i$  be Banach spaces. Then it holds:*

1. *If  $F_i: X \rightarrow Y_i$ ,  $i = 1, 2$ , are  $\partial F_i$ -semi-smooth at  $x$  then  $(F_1, F_2): X \rightarrow Y_1 \times Y_2$  is  $(\partial F_1, \partial F_2)$ -semi-smooth at  $x$ .*
2. *If  $F_i: X \rightarrow Y$ ,  $i = 1, 2$ , are  $\partial F_i$ -semi-smooth at  $x$ , then  $F_1 + F_2$  is  $(\partial F_1 + \partial F_2)$ -semi-smooth at  $x$ .*
3. *If  $F_1: Y \rightarrow Z$  and  $F_2: X \rightarrow Y$  are  $\partial F_i$ -semi-smooth at  $F_2(x)$  and  $x$ , respectively,  $\partial F_1$  is bounded near  $F_2(x)$ , and  $F_2$  is Lipschitz-continuous near  $x$ , then  $F = F_1 \circ F_2$  is  $\partial F$ -semi-smooth at  $x$  with*

$$\partial F(x) = \{G_1 G_2 \mid G_1 \in \partial F_1(F_2(x)), G_2 \in \partial F_2(x)\}.$$

4. *If  $F: X \rightarrow Y$  is continuously differentiable at  $x$ , then  $F$  is  $\{F'\}$ -semi-smooth at  $x$ , where  $\{F'\}$  denotes the set-valued operator  $\{F'\}: X \rightrightarrows \mathcal{L}(X, Y)$  with  $\{F'\}(x) = \{F'(x)\}$ .*

*Beweis.* Exercise.  $\square$

### 5.2.4. Semi-Smoothness in Finite Dimensions

We would like to discuss the application of the results towards typical optimization problems. To this end, we need to define the generalized differential. As we have seen in Example 5.2.3, it happens, that the Projection  $P_K: Z \rightarrow Z$  is given as a Nemytskii operator of a locally Lipschitz-continuous mapping  $\psi: \mathbb{R} \rightarrow \mathbb{R}$ . It is thus useful to construct a generalized derivative of  $P_K$  based on a generalized derivative of  $\psi$ .

**Definition 5.2.12.** For a locally Lipschitz-continuous function  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  we define the Clarke generalized Jacobian as

$$\partial^{\text{cl}} F(x) = \text{conv}\{G \mid \text{there exists } x^k \rightarrow x, F \text{ is differentiable at } x^k, F'(x^k) \rightarrow G\}.$$

Indeed, this definition is well-posed since by Rademacher's theorem any Lipschitz-continuous function is almost everywhere differentiable. This result can not be generalized to Fréchet-differentiability in infinite dimensions, but (for locally convex spaces) to Gâteaux derivatives, see, e.g., [Bogachev and Mayer-Wolf \[1996\]](#).

Clearly, if  $F$  is continuously differentiable in a point  $x$  then  $\partial^{\text{cl}}F(x) = \{F'(x)\}$ .

**Example 5.2.13.** Let us consider the function  $\psi: \mathbb{R} \rightarrow \mathbb{R}$  with  $\psi(x) = P_{[a,b]}(x) = \min(\max(a, x), b)$  as given in [Example 5.2.3](#). The Clarke generalized derivative is

$$\partial^{\text{cl}}\psi(x) = \begin{cases} 0 & x \notin [a, b], \\ 1 & x \in (a, b), \\ \text{conv}\{0, 1\} = [0, 1] & x \in \{a, b\}. \end{cases}$$

Indeed,  $\psi$  is then  $\partial^{\text{cl}}\psi$ -semi-smooth.

### 5.2.5. Semi-Smoothness of Nemyzkii-Operators

A natural way, to construct a generalized derivative for a Nemyzkii-operator  $\Psi$  given by superposition with  $\psi$  is by superposition with the generalized derivative of  $\psi$ . This is analogous to the differentiable case, where our candidate for a derivative is given by the superposition of the pointwise derivatives.

**Definition 5.2.14.** Let  $\psi: \mathbb{R}^m \rightarrow \mathbb{R}$  be Lipschitz-continuous and  $\partial^{\text{cl}}\psi(x)$  semi-smooth. For any  $1 \leq r \leq p \leq \infty$  we define the superposition operator

$$\Psi: L^p(\Omega)^m \rightarrow L^r(\Omega), \quad \Psi(u)(x) = \psi(u(x)) \quad \forall u \in L^p(\Omega), x \in \Omega$$

and the generalized differential

$$\partial\Psi(u) = \{G \mid Gw = g \cdot w, g \in L^\infty(\Omega)^m, g(x) \in \partial^{\text{cl}}\psi(u(x)) \text{ for all } w \in L^p(\Omega)^m \text{ for a.e. } x \in \Omega\}.$$

We note, that  $g(x)$  is uniformly bounded by Lipschitz-continuity of  $\psi$ .

From [Example 5.2.3](#) one would like to use  $p = r = 2$  as then a generalized derivative of  $P_K(G(x) + \sigma\lambda)$  could immediately be obtained by the chain-rule from [Theorem 5.2.11-3](#).

As we have to expect from [Remark 5.2.7](#) this is not possible (unless  $p = r = \infty$ ).

**Lemma 5.2.15.** Let  $\psi: \mathbb{R} \rightarrow \mathbb{R}$  be Lipschitz-continuous and **not** affine linear, and  $\Omega \subset \mathbb{R}^n$

## 5. Methods

a domain. Then, for any  $r \in [1, \infty)$  the operator

$$\Psi: L^r(\Omega) \ni u \mapsto \psi(u(\cdot)) \in L^r(\Omega)$$

is not  $\partial\Psi$ -semi-smooth.

*Beweis.* Let  $b \in \mathbb{R}$  be given and choose  $g_b \in \partial^{\text{cl}}\psi(b)$ . Since  $\psi$  is not affine linear, there is  $a \in \mathbb{R}$  such that

$$\psi(a) \neq \psi(b) + g_b(a - b)$$

and thus

$$\rho := |\psi(b) - \psi(a) - g_b(b - a)| > 0.$$

Now, let  $x_0 \in \Omega$  and  $U_\varepsilon = (x_0 - h_\varepsilon, x_0 + h_\varepsilon)^n$  with  $h_\varepsilon = \varepsilon^{1/n}/2$ . Define

$$u(x) = a, \quad x \in \Omega, \quad d_\varepsilon(x) = \begin{cases} b - a & x \in U_\varepsilon, \\ 0 & \text{otherwise.} \end{cases}$$

Hence

$$\|d_\varepsilon\|_{L^r}^r = \int_{\Omega} |d_\varepsilon(x)|^r dx = \int_{U_\varepsilon} |b - a|^r dx = \varepsilon |b - a|^r.$$

Now, let  $g_a \in \partial\psi(a)$  be arbitrary, and define

$$g_\varepsilon(x) = \begin{cases} g_b & x \in U_\varepsilon, \\ g_a & \text{otherwise.} \end{cases}$$

By this construction  $G: L^r(\Omega) \ni v \mapsto g_\varepsilon \cdot v \in L^r(\Omega)$  is an element of  $\partial\Psi(u + d_\varepsilon)$ . For arbitrary  $x \in \Omega$ , we have the remainder

$$|\psi(u(x) + d_\varepsilon(x)) - \psi(u(x)) - g_\varepsilon(x)d_\varepsilon(x)| = \begin{cases} |\psi(b) - \psi(a) - g_b(b - a)| = \rho > 0 & x \in U_\varepsilon, \\ |\psi(a) - \psi(a) - g_a \cdot 0| = 0 & \text{otherwise.} \end{cases}$$

We obtain

$$\begin{aligned} \|\Psi(u + d_\varepsilon) - \Psi(u) - Gd_\varepsilon\|_{L^r}^r &= \int_{\Omega} |\psi(u(x) + d_\varepsilon(x)) - \psi(u(x)) - g_\varepsilon(x)d_\varepsilon(x)|^r dx \\ &= \int_{U_\varepsilon} \rho^r dx \\ &= \varepsilon \rho^r \\ &= \frac{\rho^r}{|b - a|^r} \|d_\varepsilon\|_{L^r}^r. \end{aligned}$$

This shows

$$\frac{\|\Psi(u + d_\varepsilon) - \Psi(u) - Gd_\varepsilon\|_{L^r}}{\|d_\varepsilon\|_{L^r}} = \frac{\rho}{|b - a|} \neq 0$$

and thus  $\Psi$  is not  $\partial\Psi$ -semi-smooth. □

This shows, that unfortunately we can not simply use the chain-rule for proving semi-smoothness. Instead, we need to utilize the combined mapping, noting that  $G$  is not only semi-smooth, but infact Fréchet-differentiable and moreover in many cases maps to a slightly better space  $L^p$  with  $p > r$ .

**Theorem 5.2.16.** *Let  $\Omega \subset \mathbb{R}^n$  be measurable with  $0 < |\Omega| < \infty$ . Further, let  $\psi : \mathbb{R}^m \rightarrow \mathbb{R}$  be Lipschitz-continuous and  $\partial^{\text{cl}}\psi$ -semi-smooth. Let  $1 \leq r < p \leq \infty$ . Then the Nemyzkii-operator  $\Psi$  is  $\partial\Psi$ -semi-smooth (with the generalized differential given in Definition 5.2.14).*

*Beweis.* Assume that  $\Psi$  is not semi-smooth at  $u \in L^p(\Omega)^m$ . Then there exists  $\varepsilon > 0$ ,  $d^k \rightarrow 0$  in  $L^p(\Omega)^m$ , and a sequence  $G_k \in \partial\Psi(u + d^k)$  with

$$\|\psi(u + d^k) - \psi(u) - G_k d^k\|_{L^r} \geq \varepsilon \|d^k\|_{L^p}.$$

By definition it is (for almost every  $x \in \Omega$ )

$$(\Psi(u + d^k) - \Psi(u) - G_k d^k)(x) = \psi(u(x) + d^k(x)) - \psi(u(x)) - g^k(x) d^k(x)$$

with  $g^k(x) \in \partial^{\text{cl}}\psi(u(x) + d^k(x))$ . If  $L$  is the Lipschitz-constant of  $\psi$  it holds

$$\begin{aligned} |\psi(u(x) + d^k(x)) - \psi(u(x)) - g^k(x) d^k(x)| &\leq |\psi(u(x) + d^k(x)) - \psi(u(x))| + |g^k(x) d^k(x)| \\ &\leq L |d^k(x)| + |g^k(x) d^k(x)| \\ &\leq 2L |d^k(x)|. \end{aligned}$$

By convergence  $d^k \rightarrow 0$  in  $L^p$ , we can select a subsequence such that  $d^k(x) \rightarrow 0$  pointwise almost everywhere. As a consequence of pointwise convergence and the definition  $g^k(x) \in \partial^{\text{cl}}\psi(u(x) + d^k(x))$ , for any given  $\delta > 0$  the sets

$$M_{\delta,k} := \{x \in \Omega \mid |\psi(u(x) + d^k(x)) - \psi(u(x)) - g^k(x) d^k(x)| \geq \delta |d^k(x)|\}$$

satisfy  $|M_{\delta,k}| \rightarrow 0$  as  $k \rightarrow \infty$ . Now,  $p > r$  and thus Hölder inequality with  $\frac{1}{r} = \frac{1}{p} + \frac{1}{s}$  and thus  $\frac{1}{s} = \frac{p-r}{pr} > 0$  (with the suitable interpretation for  $p = \infty$ ) gives

$$\begin{aligned} \|\Psi(u + d^k) - \Psi(u) - G_k d^k\|_{L^r} &\leq \|\Psi(u + d^k) - \Psi(u) - G_k d^k\|_{L^r(M_{\delta,k}^c)} + \|\Psi(u + d^k) - \Psi(u) - G_k d^k\|_{L^r(M_{\delta,k})} \\ &\leq \delta \|d^k\|_{L^r(M_{\delta,k}^c)} + 2L \|d^k\|_{L^r(M_{\delta,k})} \\ &\leq \delta \|1\|_{L^s(M_{\delta,k}^c)} \|d^k\|_{L^p(M_{\delta,k}^c)} + 2L \|1\|_{L^s(M_{\delta,k})} \|d^k\|_{L^p(M_{\delta,k})} \\ &\leq \left( \delta |\Omega|^{1/s} + 2L |M_{\delta,k}|^{1/s} \right) \|d^k\|_{L^p}. \end{aligned}$$

Selecting  $\delta |\Omega|^{1/s} \leq \varepsilon/2$  gives

$$\frac{\|\Psi(u + d^k) - \Psi(u) - G_k d^k\|_{L^r}}{\|d^k\|_{L^p}} \leq \varepsilon/2 + 2L |M_{\delta,k}|^{1/s} \rightarrow \varepsilon/2$$

as  $k \rightarrow \infty$ ; contradicting the assumption. □

## 5. Methods

**Corollary 5.2.17.** *Let  $\psi: \mathbb{R}^m \rightarrow \mathbb{R}$  be Lipschitz-continuous and  $\partial^{\text{cl}}\psi(x)$  semi-smooth. Further, let  $1 \leq r < p \leq \infty$  be given,  $U$  be a Banach space,  $G: U \rightarrow L^p(\Omega)^m$  Fréchet-differentiable, and  $\Omega \subset \mathbb{R}^n$  be measurable with  $0 < |\Omega| < \infty$ . Define*

$$\Psi_G: U \rightarrow L^r(\Omega), \quad \Psi_G(u)(x) = \psi(G(u)(x))$$

*for any  $u \in U$  and  $x \in \Omega$ . Further, we define the generalized differential*

$$\partial \Psi_G: U \rightrightarrows \mathcal{L}(U; L^r(\Omega))$$

$$\partial \Psi_G(u) = \{G \mid Gv = g \cdot (G'(u)v), g \in L^\infty(\Omega)^m, g(x) \in \partial^{\text{cl}}\psi(G(u)(x)) \text{ for a.e. } x \in \Omega\}.$$

*Then  $\Psi_G$  is  $\partial \Psi_G$  semi-smooth.*

*Beweis.* This is an immediate consequence of Theorem 5.2.16 and Theorem 5.2.11 (3. and 4.).  $\square$

### 5.2.6. Application

We now consider an optimization problem to discuss the applicability of the afore discussed theory.

$$\begin{aligned} \min j(q) \\ \text{s.t. } q \in L^2(\Omega), \quad a \leq q \leq b \text{ a.e. on } \Omega. \end{aligned} \tag{5.6}$$

Where  $\Omega \subset \mathbb{R}^n$  is some measurable set (Domain, boundary, ...) with  $0 < |\Omega| < \infty$ . W.l.o.g.  $a < b \in \mathbb{R}$ , otherwise non-constant bounds could be transformed to constant bounds by substituting  $q \mapsto \frac{q-a}{b-a}$ . Clearly, in this example  $G = \text{Id}$  (asserting (RCQ)) and thus the first order optimality conditions can be written as

$$\begin{aligned} \nabla j(q) + \lambda &= 0, \\ q - P_K(q + \sigma \lambda) &= 0 \end{aligned}$$

following Example 5.2.3, with the projection

$$P_K(f)(x) = \min(\max(a, f(x)), b).$$

Clearly, the two lines can equivalently be written as

$$q - P_K(q - \sigma \nabla j(q)) = 0.$$

As we have seen,  $P_K: L^2 \rightarrow L^2$  will not be semi-smooth. Hence, we need to assume additional structure on  $\nabla j(q)$  which can be satisfied in many cases.

**Assumption 5.2.18.** Assume that there exists  $\alpha > 0$  and  $p > 2$  such that

1.  $\nabla j(q) = \alpha q + H(q)$
2.  $H: L^2(\Omega) \rightarrow L^p(\Omega)$  is continuously Fréchet-differentiable.

With this assumptions, we can select  $\sigma = 1/\alpha$  to have the optimality conditions

$$\Phi(q) = q - P_K(q - \sigma \nabla j(q)) = q - P_K\left(\frac{-1}{\alpha} H(q)\right) = 0. \quad (5.7)$$

This operator  $\Phi(q) = q - P_K(\frac{-1}{\alpha} H(q))$  is now matching our conditions of Corollary 5.2.17 and is thus semi-smooth.

Hence we have seen

**Corollary 5.2.19.** Consider the problem (5.6) with Assumption 5.2.18. Then, for  $\sigma = 1/\alpha$  the operator  $\Phi$  in (5.7) is  $\partial\Phi$ -semi-smooth with

$$\partial\Phi: L^2(\Omega) \rightrightarrows \mathcal{L}(L^2(\Omega), L^2(\Omega)),$$

$$\partial\Phi(q) = \{G \mid G = \text{Id} + \frac{g}{\alpha} H'(q), \quad g \in L^\infty(\Omega), \quad g(x) \in \partial^{\text{cl}} P_{[a,b]}(-1/\alpha H(q)(x)) \text{ for a.e. } x \in \Omega\}$$

where

$$\partial^{\text{cl}} P_{[a,b]}(s) = \begin{cases} 0 & x \notin [a, b], \\ 1 & x \in (a, b), \\ [0, 1] & \text{otherwise.} \end{cases}$$

For the application of the semi-smooth Newton method 5.3 one now needs to show, that the inverse of all elements in  $\partial\Phi(q)$  are uniformly bounded close to  $\bar{q}$ , i.e.,

$$\|G^{-1}\|_{\mathcal{L}(L^2, L^2)} \leq C \quad \forall G \in \partial\Phi(q), q \in B_\delta^{L^2}(\bar{q}).$$

### 5.2.7. Control Constrained Optimal Control

We come back to a control constrained setting similar to Example 3.3.9 but with a bit more explicit structure.

$$\begin{aligned} \min J(q, u) &= \frac{1}{2} \|u - u^d\|^2 + \frac{\alpha}{2} \|q\|^2 \\ \text{s.t. } &\begin{cases} -\Delta u = q, & \in H^{-1}(\Omega) = H_0^1(\Omega)^* \\ q \in Q^{\text{ad}} := \{p \in Q = L^2(\Omega) \mid a \leq p \leq b\}, \end{cases} \end{aligned}$$

## 5. Methods

on a bounded domain  $\Omega \subset \mathbb{R}^2$ . Since the operator  $-\Delta$  defines an isomorphism from  $H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$  we can consider the (linear) solution operator  $S: H^{-1}(\Omega) \rightarrow H_0^1(\Omega)$ . Notice, that the embeddings  $H_0^1(\Omega) \rightarrow L^p(\Omega)$  are compact (in 2d for any  $p < \infty$ ). This gives the reduced problem

$$\begin{aligned} \min j(q) &= \frac{1}{2} \|Sq - u^d\|^2 + \frac{\alpha}{2} \|q\|^2 \\ \text{s.t. } q &\in Q^{\text{ad}} := \{p \in Q = L^2(\Omega) \mid a \leq p \leq b\}. \end{aligned}$$

A short calculation shows that w.r.t the  $L^2$ -inner product we have

$$\begin{aligned} \nabla j(q) &= \alpha q + S^*(Sq - u^d) \in L^2(\Omega), \\ \nabla^2 j(q) &= \alpha \text{Id} + S^*S \in \mathcal{L}(L^2, L^2). \end{aligned}$$

Hence we identify Assumption 5.2.18 with  $H(q) = S^*(Sq - u^d)$ . Note that the adjoint  $S^*: H^{-1}(\Omega) \rightarrow H_0^1(\Omega)$  and hence

$$S^*(Sq - u^d) \in H_0^1(\Omega) \subset L^p(\Omega) \subset L^2(\Omega)$$

for some  $p > 2$ , and thus it is differentiable as a mapping  $L^2(\Omega) \rightarrow L^p(\Omega)$ .

Hence our optimality system can be written in terms of the semi-smooth equation

$$\Phi(q) = q - P_{[a,b]} \left( \frac{-1}{\alpha} S^*Sq \right) = 0.$$

To simplify the following calculations, we select the particular generalized derivative given by selecting

$$g(x) = \begin{cases} 1 & \frac{-1}{\alpha} S^*(Sq(x) - u^d(x)) \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

then the generalized derivative of  $\Phi$  is given by

$$\partial \Phi(q) = \left\{ \text{Id} + \frac{1}{\alpha} \chi_{\mathcal{J}} S^*S \right\} \quad (5.8)$$

where  $\chi_{\mathcal{J}}$  is the indicator function of the inactive set of the projection  $P_{[a,b]}$

$$\mathcal{J} = \left\{ x \in \Omega \mid \frac{-1}{\alpha} (S^*Sq)(x) \in [a, b] \right\}$$

**Lemma 5.2.20.** *Let  $G \in \partial \Phi(q)$  given by (5.8). Then  $G$  is invertible and satisfies*

$$\|G^{-1}\|_{\mathcal{L}(L^2(\Omega), L^2(\Omega))} \leq 2 + \frac{1}{\alpha} \|S^*S\|_{\mathcal{L}(L^2(\Omega), L^2(\Omega))}.$$



*Beweis.* To see invertibility of  $G$ , we notice that  $G = \text{id} + \frac{1}{\alpha}\chi_{\mathcal{J}}S^*S$  is a compact perturbation of the identity and thus Fredholm of index zero and invertibility follows if  $\text{kern } G = \{0\}$ .

To this end, let  $d, F \in L^2(\Omega)$  be given satisfying the equation  $Gd = F$ , i.e.,

$$F = Gd = \frac{1}{\alpha}\chi_{\mathcal{J}}S^*Sd + d.$$

On the active set  $\mathcal{A} = \Omega \setminus \mathcal{J}$  this gives

$$\chi_{\mathcal{A}}d = \chi_{\mathcal{A}}F.$$

Now we use this to determine  $d|_{\mathcal{J}}$ . To this end, we consider the relation

$$\frac{1}{\alpha}S^*S(\chi_{\mathcal{J}}d) + \chi_{\mathcal{J}}d = F - \frac{1}{\alpha}S^*S(\chi_{\mathcal{A}}d)$$

on  $\mathcal{J}$ .

Then it holds

$$\begin{aligned} \int_{\Omega} (\chi_{\mathcal{J}}d)^2 dx &\leq \frac{1}{\alpha} \int_{\Omega} (S^*S(\chi_{\mathcal{J}}d)\chi_{\mathcal{J}}d) dx + \int_{\Omega} (\chi_{\mathcal{J}}d)^2 dx \\ &= \int_{\Omega} (F - \frac{1}{\alpha}S^*S(\chi_{\mathcal{A}}d)(\chi_{\mathcal{J}}d)) dx. \end{aligned}$$

By Hölder's inequality, we get

$$\begin{aligned} \|d\|_{L^2(\Omega)} &\leq \|\chi_{\mathcal{J}}d\|_{L^2(\Omega)} + \|\chi_{\mathcal{A}}d\|_{L^2(\Omega)} \\ &\leq (\|F\|_{L^2(\Omega)} + \frac{1}{\alpha}\|S^*S(\chi_{\mathcal{A}}d)\|_{L^2(\Omega)}) + \|\chi_{\mathcal{A}}d\|_{L^2(\Omega)} \\ &\leq 2\|F\|_{L^2(\Omega)} + \frac{1}{\alpha}\|S^*S\|_{\mathcal{L}(L^2(\Omega), L^2(\Omega))}\|F\|_{L^2(\Omega)}. \end{aligned}$$

Hence taking  $F = 0$  this gives  $\text{kern } G = \{0\}$  and thus

$$\|G^{-1}F\|_{L^2(\Omega)} \leq 2\|F\|_{L^2(\Omega)} + \frac{1}{\alpha}\|S^*S\|_{\mathcal{L}(L^2(\Omega), L^2(\Omega))}\|F\|_{L^2(\Omega)}$$

showing the assertion. □

## 5. *Methods*

## 6. Discretization

To sketch the necessary steps in analysing discrete approximations to optimization problems in function spaces, we consider the following model problem

$$\min J(q, u) \quad \text{s.t. (6.1).}$$

Where the constraint is to find  $u \in H_0^1(\Omega)$  such that

$$(a \nabla u, \nabla \varphi) = (f, \varphi) \quad \forall \varphi \in H_0^1(\Omega) \quad (6.1)$$

for some  $a \in W^{1,\infty}(\Omega)$  with  $0 < \underline{a} \leq a(a) \leq \bar{a}$  on  $\Omega$ . Then a generic way to discretize this equation is to choose some finite dimensional space  $V_h$  (for convenience let  $V_h \subset H_0^1(\Omega)$  then  $V_h$  is called  $H_0^1$ -conform). Then one considers the problem to find  $u_h \in V_h$  such that

$$(a \nabla u_h, \nabla \varphi_h) = (f, \varphi_h) \quad \forall \varphi_h \in V_h. \quad (6.2)$$

In particular, one searches for a solution  $u_h$  in the *ansatz space*  $V_h$  such that the variational equation is satisfied for all  $\varphi_h$  in the *test space*  $V_h$  (Although here they are the same.). The equation (6.2) has a unique solution by the same arguments that are used for (6.1).

**Theorem 6.0.1.** *Let  $u$  be the solution to (6.1) and  $u_h$  be the solution to (6.2). Then the following quasi-best-approximation property holds*

$$\|\nabla(u - u_h)\| \leq \frac{\bar{a}}{\underline{a}} \inf_{\varphi_h \in V_h} \|\nabla(u - \varphi_h)\|. \quad (6.3)$$

*Beweis.* The proof follows from the Galerkin orthogonality, e.g.,

$$(a \nabla(u - u_h), \nabla \varphi_h) = 0 \quad \forall \varphi_h \in V_h.$$

□

By choosing a basis  $\varphi_h^{(i)}$ ,  $i = 1, \dots, N$  of  $V_h$  the problem (6.2) can be rewritten as a linear algebraic problem as follows: Find  $\mathbf{u}_h = (u_h^{(i)})_{i=1, \dots, N}$  such that

$$\mathbf{A}_h \mathbf{u}_h = \mathbf{f}_h$$

## 6. Discretization

with the stiffness matrix

$$\mathbf{A}_h = (a_{ij})_{i,j=1,\dots,N}, \quad a_{ij} = (a \nabla \varphi_h^{(j)}, \nabla \varphi_h^{(i)})$$

and the right hand side

$$\mathbf{f}_h = ((f, \varphi_h^{(i)}))_{i=1,\dots,N}.$$

Now, in order to be able to solve this problem efficiently one needs to choose the basis in a clever way, in particular such that it is „easy” to invert  $\mathbf{A}_h$ . The idea to choose an orthogonal basis is usually not computationally feasible, and may not be possible for general equations. One possible way to obtain a reasonably well behaved matrix  $\mathbf{A}_h$  is given by the so called *finite element method*.

### 6.1. Linear Finite Elements for Elliptic Problems

For simplicity let  $\Omega \subset \mathbb{R}^2$  be a bounded polygonal domain. Then we consider a sequence of decompositions (triangulations)  $\mathcal{T}_h = T$  of  $\overline{\Omega}$  into closed polygons  $T$  with  $h := \max_{T \in \mathcal{T}_h} \text{diam}(T) \rightarrow 0$ . For simplicity let us assume that there exists some reference element  $\hat{T}$  such that any element  $T \in \mathcal{T}_h$  is given by an affine-linear transformation of  $\hat{T}$ . Typical elements  $T$  are triangles or quadrilaterals.

We require the following regularity properties for  $\mathcal{T}_h$ .

1. (Structural regularity or feasibility) Any two different elements  $T_1, T_2 \in \mathcal{T}_h$  intersect at most in a vertex or along an entire edge.
2. (Shape regularity) For the radius  $\rho_T$  of the incircle and  $h_T$  of the circumscribed circle of an element  $T \in \mathcal{T}_h$  it holds

$$\max_{T \in \mathcal{T}_h} \frac{h_T}{\rho_T} \leq c$$

uniformly in  $h \rightarrow 0$ .

3. (Size regularity) It holds

$$\max_{T \in \mathcal{T}_h} h_T \leq c \min_{T \in \mathcal{T}_h} h_T$$

uniformly in  $T$ .

We can now describe the construction of a finite dimensional space  $V_h$ . For simplicity let us assume that the decompositions  $\mathcal{T}_h$  consist of triangles only. Then we can define

$$V_h^{(1)} := \{v_h \in C(\overline{\Omega}) \mid v_h|_T \in P_1(T), T \in \mathcal{T}_h, v_h|_{\Gamma} = 0\}.$$

where  $P_1(T)$  denotes the space of polynomials of degree  $\leq 1$  on  $T$ . This defines the space of (piecewise) linear finite elements. It is straightforward to see that  $V_h^{(1)} \subset H_0^1(\Omega)$ .

## 6.1. Linear Finite Elements for Elliptic Problems

Now to define a suitable basis of  $V_h^{(1)}$  we denote the interior vertices of  $\mathcal{T}_h$  by  $x_i$  for  $i = 1, \dots, N$ . Then we define the nodal basis  $\varphi_h^{(i)}$  (or Lagrange basis) as follows:

$$\varphi_h^{(i)}(x_j) = \delta_{ij} = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

**Remark 6.1.1.** Of course higher order polynomials are also possible, e.g., piecewise polynomials of degree two. Which then require additional basis functions, for example midpoints of the edges.

On quadrilateral elements the definition of piecewise polynomials is not satisfactory hence one considers usually functions that are piecewise bilinear (biquadratic, ...). However the definition on general quadrilaterals is more involved, as one needs to account for the transformation from the reference element (which may be non affine!).

### 6.1.1. A Priori Error Estimates

As we have already seen, the error in the natural norm is given by the best approximation property. Hence to obtain error estimates it is natural to ask what is the best possible approximation. To answer this, we consider the Lagrange interpolation

**Definition 6.1.2.** For a function  $u \in C(\bar{\Omega})$  we define the Lagrange interpolation  $I_h u: C(\bar{\Omega}) \rightarrow V_h^{(1)}$  as follows

$$I_h u(x) = \sum_{i=1}^N u(x_i) \varphi_h^{(i)}(x).$$

**Theorem 6.1.3.** Let  $\mathcal{T}_h$  satisfy the regularity assumptions above. Then there exists a constant  $c > 0$  independent of  $h$  such that for  $u \in W^{2,2}(\Omega)$  it holds

$$\|u - I_h u\|_2 + h \|\nabla(u - I_h u)\|_2 \leq ch^2 \|u\|_{2,2}.$$

*Beweis.* See, e.g., [Brenner and Scott, 2008, Theorem 4.4.20] □

**Theorem 6.1.4.** Let  $\mathcal{T}_h$  satisfy the regularity assumptions above, and let  $\Omega$  be convex. Then

## 6. Discretization

the solution  $u$  to (6.1) and  $u_h$  to (6.2) satisfy for  $f \in L^2(\Omega)$

$$\|u - u_h\| + h\|\nabla(u - u_h)\| \leq ch^2\|f\|.$$

*Beweis.* From Theorem 6.0.1 and Theorem 6.1.3 we have that

$$\|\nabla(u - u_h)\| \leq ch\|u\|_{2,2}$$

and by the elliptic regularity, see, e.g., [Gilbarg and Trudinger \[2001\]](#), it follows

$$\|\nabla(u - u_h)\| \leq ch\|f\|.$$

For the error in the  $L^2$ -Norm we employ the Aubin-Nitsche trick. Define  $e_h = u - u_h$  and  $p \in H_0^1(\Omega)$  as solution to

$$(a\nabla\varphi, \nabla p) = 1/\|e_h\|(e_h, \varphi) \quad \forall \varphi \in H_0^1(\Omega).$$

Then, because  $e_h \in L^2(\Omega)$  it is  $p \in H^2(\Omega)$  and it holds

$$\|p\|_{2,2} \leq c\|e_h\|/\|e_h\| = c.$$

Then we obtain

$$\|e_h\| = (a\nabla e_h, \nabla p) = (a\nabla e_h, \nabla(p - I_h p)) \leq \|\nabla e_h\| \|\nabla(p - I_h p)\| \leq ch\|\nabla e_h\|$$

hence the assertion follows.  $\square$

**Remark 6.1.5.** The above theorem states that the operator  $S: L^2(\Omega) \rightarrow L^2(\Omega)$  given by (6.1) and the discrete operator  $S_h: L^2(\Omega) \rightarrow L^2(\Omega)$  given by (6.2) satisfy

$$\|(S - S_h)f\| \leq ch^2\|f\| \quad \forall f \in L^2(\Omega)$$

or equivalently

$$\|S - S_h\|_{\mathcal{L}(L^2(\Omega))} \leq ch^2.$$

## 6.2. Discretization of the Model Problem

We return to the model problem

$$\begin{aligned} \min J(q, u) &= \frac{1}{2}\|u - u^d\|^2 + \frac{\alpha}{2}\|q\|^2 \\ \text{s.t. } u &\in H_0^1(\Omega), \\ (a\nabla u, \nabla \varphi) &= (q, \varphi) \quad \forall \varphi \in H_0^1(\Omega), \\ q_{\min} &\leq q \leq q_{\max} \quad \text{a.e. on } \Omega. \end{aligned} \tag{6.4}$$

Here we assume that  $\Omega \subset \mathbb{R}^2$  is a convex, bounded polygonal domain. Further we assume that we have a sequence of feasible and (shape and size) regular triangulations  $\mathcal{T}_h$  of  $\Omega$ .

### 6.2.1. Variational Discretization

Now we can replace the state equation using the finite element method on a sequence of feasible and (shape and size) regular triangulations. This gives

$$\begin{aligned} \min J(q_h, u_h) &= \frac{1}{2} \|u_h - u^d\|^2 + \frac{\alpha}{2} \|q_h\|^2 \\ \text{s.t. } u_h &\in V_h, \\ (a \nabla u_h, \nabla \varphi_h) &= (q_h, \varphi_h) \quad \forall \varphi_h \in V_h, \\ q_{\min} &\leq q_h \leq q_{\max} \quad \text{a.e. on } \Omega. \end{aligned} \tag{6.5}$$

We note that we did not discretize the control space, and hence this is (at least formally) an infinite dimensional problem.

**Theorem 6.2.1.** *For  $\alpha > 0$  there exists a unique solution  $(\bar{q}_h, \bar{u}_h) = (\bar{q}_h, S_h \bar{q}_h) \in Q^{ad} \times V_h$  to (6.5). This solution is (equivalently) characterized by the existence of  $\bar{z}_h \in V_h$  such that the following holds:*

$$\begin{aligned} (a \nabla \bar{u}_h, \nabla \varphi_h) &= (\bar{q}_h, \varphi_h) & \forall \varphi_h \in V_h, \\ (a \nabla \varphi_h, \nabla \bar{z}_h) &= (\bar{u}_h - u^d, \varphi_h) & \forall \varphi_h \in V_h, \\ (\bar{z}_h + \alpha \bar{q}_h, q - \bar{q}_h) &\geq 0 & \forall q \in Q^{ad}. \end{aligned}$$

*Beweis.* We note that the operator  $S_h : L^2(\Omega) \rightarrow V_h \subset L^2(\Omega)$  is linear and continuous. Hence the reduced cost functional

$$j_h(q_h) = J(q_h, S_h q_h)$$

is convex and continuous and thus w.l.s.c. with  $j_h(q_h) \rightarrow \infty$  as  $\|q_h\| \rightarrow \infty$ . Thus we can apply Theorem 2.2.3 to obtain existence of a solution  $(\bar{q}_h, \bar{u}_h)$  to (6.5). Due to strict convexity of  $j_h$  this solution is unique.

Now as for  $j$  the functional  $j_h$  is continuously Fréchet differentiable and hence by Theorem 3.2.3 the solution  $\bar{q}_h$  to

$$\min_{q_{\min} \leq q_h \leq q_{\max}} j_h(q_h)$$

is equivalently characterized by the variational inequality

$$(S_h^*(S_h \bar{q}_h - u^d) + \alpha \bar{q}_h, q - \bar{q}_h) \geq 0 \quad \forall q \in Q^{ad}.$$

Now, to conclude we need to derive an equation that is defining  $S_h^*$ . To this end, we compute for  $f, g \in L^2(\Omega)$

$$(S_h^* g, f) = (g, S_h f) = (a \nabla S_h g, \nabla S_h f) = (a \nabla S_h f, \nabla S_h g) = (f, S_h g)$$

Thus  $S_h^* = S_h$  and the assertion follows.  $\square$

## 6. Discretization

**Remark 6.2.2.** We remark that the optimality conditions obtained here are the discretized optimality conditions of the continuous case. This means that it does not matter whether we discretize the optimization problem and compute necessary optimality conditions or if we directly discretize the optimality conditions. This is often referred to as ‘optimize-then-discretize = discretize-then-optimize’. This is not always the case, e.g., in finite difference schemes this need not happen.

**Theorem 6.2.3.** Let  $(\bar{q}, \bar{u})$  be the solution to (6.4) and  $(\bar{q}_h, \bar{u}_h)$  be the solution to (6.5). Then there exists a constant (independent of  $h$  and  $\alpha$ ) such that

$$\alpha \|\bar{q} - \bar{q}_h\|^2 + \|\bar{u} - \bar{u}_h\|^2 \leq c \left(1 + \frac{1}{\alpha}\right) h^4.$$

*Beweis.* We recall the necessary optimality condition for  $\bar{q}$  and  $\bar{q}_h$  namely

$$\begin{aligned} (S^*(S\bar{q} - u^d) + \alpha\bar{q}, q - \bar{q}) &\geq 0 \quad \forall q \in Q^{\text{ad}}, \\ (S_h^*(S_h\bar{q}_h - u^d) + \alpha\bar{q}_h, q - \bar{q}_h) &\geq 0 \quad \forall q \in Q^{\text{ad}}. \end{aligned}$$

Now we see that  $\bar{q}_h$  is a feasible test function for the first inequality while  $\bar{q}$  is a feasible test function for the second inequality. Together this gives

$$\begin{aligned} 0 &\leq (S^*(S\bar{q} - u^d) - S_h^*(S_h\bar{q}_h - u^d) + \alpha\bar{q} - \alpha\bar{q}_h, \bar{q}_h - \bar{q}) \\ &= -\alpha \|\bar{q}_h - \bar{q}\|^2 \\ &\quad + (S\bar{q} - u^d, S(\bar{q}_h - \bar{q})) - (S_h\bar{q}_h - u^d, S_h(\bar{q}_h - \bar{q})) \\ &= -\alpha \|\bar{q}_h - \bar{q}\|^2 \\ &\quad + (S\bar{q} - S_h\bar{q}_h, S_h(\bar{q}_h - \bar{q})) + (S\bar{q} - u^d, (S - S_h)(\bar{q}_h - \bar{q})) \\ &= -\alpha \|\bar{q}_h - \bar{q}\|^2 \\ &\quad + (S\bar{q} - S_h\bar{q}_h, S_h\bar{q}_h - S\bar{q}) + (S\bar{q} - S_h\bar{q}_h, S\bar{q} - S_h\bar{q}) + (S\bar{q} - u^d, (S - S_h)(\bar{q}_h - \bar{q})) \\ &= -\alpha \|\bar{q}_h - \bar{q}\|^2 - \|\bar{u} - \bar{u}_h\|^2 \\ &\quad + (S\bar{q} - S_h\bar{q}_h, (S - S_h)\bar{q}) + ((S - S_h)^*(S\bar{q} - u^d), \bar{q}_h - \bar{q}) \end{aligned}$$

We obtain, using  $(S - S_h)^* = S^* - S_h^*$ , Young’s inequality, and a generic constant  $c$  (which may



be different in each line of the inequality),

$$\begin{aligned}
\alpha \|\bar{q}_h - \bar{q}\|^2 + \|\bar{u} - \bar{u}_h\|^2 &\leq \|S\bar{q} - S_h \bar{q}_h\| \| (S - S_h) \bar{q} \| + \| (S^* - S_h^*) (S\bar{q} - u^d) \| \| \bar{q}_h - \bar{q} \| \\
&\leq (\| (S - S_h) \bar{q} \| + \| S_h (\bar{q} - \bar{q}_h) \|) \| (S - S_h) \bar{q} \| \\
&\quad + \| (S^* - S_h^*) (S\bar{q} - u^d) \| \| \bar{q}_h - \bar{q} \| \\
&\leq ch^2 (\| (S - S_h) \bar{q} \| + \| S_h (\bar{q} - \bar{q}_h) \| + \| \bar{q}_h - \bar{q} \|) \\
&\leq ch^4 + ch^2 \| \bar{q}_h - \bar{q} \| \\
&\leq ch^4 + \frac{c}{\alpha} h^4 + \frac{\alpha}{2} \| \bar{q}_h - \bar{q} \|^2
\end{aligned}$$

By subtracting the last term, we get

$$\alpha \|\bar{q}_h - \bar{q}\|^2 + \|\bar{u} - \bar{u}_h\|^2 \leq c(1 + \frac{1}{\alpha})h^4$$

which proofs the assertion.  $\square$

**Remark 6.2.4.** It is clear that if  $Q$  is a finite dimensional space, e.g., we have only finitely many control variables, then this discretization is sufficient to obtain a finite dimensional problem.

However, this technique yields a discrete optimal control problem even if  $Q$  is not finite dimensional. To see this one can again consider the case of pure control constraints, in this case the variational inequality

$$(\bar{z}_h + \alpha \bar{q}_h, q - \bar{q}_h) \geq 0 \quad \forall q \in Q^{\text{ad}}.$$

holds again pointwise almost everywhere, by our usual arguments, yielding that

$$\bar{q}_h = \mathcal{P}_{Q^{\text{ad}}} \left( \frac{-1}{\alpha} \bar{z}_h \right).$$

This can then be used to eliminate the control in the optimality conditions given by Theorem 6.2.1. This means that

$$\left( \mathcal{P}_{Q^{\text{ad}}} \left( \frac{-1}{\alpha} \bar{z}_h \right), \bar{u}_h \right)$$

solves (6.5) if and only if  $(\bar{u}_h, \bar{z}_h)$  solves the following nonlinear (and non smooth) problem

$$\begin{aligned}
(\nabla \bar{u}_h, \nabla \varphi_h) &= (\mathcal{P}_{Q^{\text{ad}}}(-1/\alpha \bar{z}_h), \varphi) & \forall \varphi_h \in V_h, \\
(\nabla \varphi_h, \nabla \bar{z}_h) &= (\bar{u}_h - u^d, \varphi_h) & \forall \varphi_h \in V_h.
\end{aligned}$$

## 6. Discretization

This technique is useful in particular if one uses simple constraints such as box-constraints when the projection can be evaluated easily.

This technique has been advertised in [Hinze \[2005\]](#), see also [Hinze et al. \[2009\]](#)

**Remark 6.2.5.** It was recently shown [Gaspoz et al. \[2020\]](#) that in fact the preasymptotic estimate in Theorem 6.2.3 overestimates the true error. In fact, it was shown, that the quasi-best approximation estimate

$$\|\nabla(u - u_h)\| + \|\nabla(z - z_h)\| \leq c_h \inf_{\varphi_h, \psi_h \in V_h} \left( \|\nabla(u - \varphi_h)\| + \|\nabla(z - \psi_h)\| \right) \quad (6.6)$$

with a constant  $c_h \rightarrow \frac{\bar{a}}{\underline{a}}$  so that asymptotically the same best-approximation property as for the PDE (6.3) holds.

**Remark 6.2.6.** Of course, next to the implicit, variational, discretization one can choose to discretize the control space as well. Analogous estimates to Theorem 6.2.3 can be shown, by more complicated arguments. However, a quasi-best approximation as in 6.6 will no longer hold, see, e.g., [Gaspoz et al. \[2020\]](#).

# A. Functional Analytic Background

## A.1. Normed Linear Spaces

**Definition A.1.1.** A  $\mathbb{R}$ -linear space  $V$  is called normed linear space if there exists a map  $\|\cdot\| = \|\cdot\|_V: V \rightarrow \mathbb{R}_{\geq 0}$  which satisfies:

1.  $\|v\| = 0$  if and only if  $v = 0$ .
2.  $\|v + w\| \leq \|v\| + \|w\|$  for all  $v, w \in V$ .
3.  $\|\lambda v\| = |\lambda| \|v\|$  for all  $\lambda \in \mathbb{R}$  and  $v \in V$ .

**Definition A.1.2.** A sequence  $v_k \in V$  in a normed linear space  $V$  is called *Cauchy sequence* if for any given  $\varepsilon$  there exists some  $n_0 \in \mathbb{N}$  such that

$$\|v_k - v_l\| \leq \varepsilon \quad \forall k, l \geq n_0.$$

**Definition A.1.3.** A normed linear space  $V$  is called *Banach space* if it is complete, i.e., any Cauchy sequence in  $V$  has a limit in  $V$ .

**Definition A.1.4.** A Banach space  $V$  whose norm is induced by a scalar product is called a *Hilbert space*. We recall, a scalar product is a positive definite bilinear form,  $(\cdot, \cdot)_V = (\cdot, \cdot): V \times V \rightarrow \mathbb{R}$ , e.g.,

- $(v, v) \geq 0$  for all  $v \in V$  and  $(v, v) = 0$  if and only if  $v = 0$ .
- $(v, w) = (w, v)$  for all  $v, w \in V$ .
- $(\lambda v_1 + v_2, w) = \lambda(v_1, w) + (v_2, w)$  for all  $\lambda \in \mathbb{R}$ ,  $v_1, v_2, w \in V$ .

The scalar product induces a norm by the following definition

$$\|v\| = \sqrt{(v, v)}.$$

## A. Functional Analytic Background

**Example A.1.5.** An example for an infinite dimensional Banach space is given by the space  $C[0, 1]$  of continuous functions on the interval  $[0, 1]$  equipped with the norm

$$\|f\|_{\infty} = \max_{x \in [0, 1]} |f(x)|.$$

**Example A.1.6.** An example for an infinite dimensional Hilbert space is given by the space  $l_2$  of square summable sequences, e.g.,  $(x_i)_{i=0}^{\infty} \in l_2$  if and only if  $\sum_{i=0}^{\infty} x_i^2 < \infty$ . The scalar product is given by

$$(x, y) = \sum_{i=0}^{\infty} x_i y_i, \quad x, y \in l_2.$$

**Remark A.1.7.** A norm  $\|\cdot\|$  on a normed space  $V$  is induced by a scalar product if and only if the parallelogram law is satisfied, e.g., it holds for  $v, w \in V$

$$2\|v\|^2 + 2\|w\|^2 = \|v + w\|^2 + \|v - w\|^2.$$

Further, if  $\|\cdot\|$  is induced by a scalar product  $(\cdot, \cdot)$  then the Cauchy Schwarz inequality

$$|(v, w)| \leq \|v\| \|w\|$$

holds for any  $v, w \in V$ .

We conclude this section with an important definition.

**Definition A.1.8.** A normed linear space  $V$  is called *separable*, if there exists a countable dense (w.r.t the norm on  $V$ ) subset of  $V$ .

**Example A.1.9.** A typical example of a separable Banach space is again  $C[0, 1]$  with the norm  $\|\cdot\|_{\infty}$ . Since it is known from calculus lessons, that any continuous function can be approximated uniformly by polynomials (Theorem of Stone-Weierstrass).

## A.2. Convexity

**Definition A.2.1.** Let  $V$  be a vector space. A set  $C \subset V$  is called *convex*, if for any  $\lambda \in (0, 1)$  and  $v, w \in C$  it holds

$$\lambda v + (1 - \lambda)w \in C.$$

The *convex hull* of a set  $A \subset V$  is defined as

$$\bigcap_{\substack{C \supset A \\ C \text{ convex}}} C.$$

**Definition A.2.2.** Let  $V$  be a vector space and  $C \subset V$  be convex. A function  $f : C \rightarrow \mathbb{R} \cup \{\infty\}$  is called *convex*, if (with the usual definitions for arithmetic with  $\infty$ ) for any  $\lambda \in (0, 1)$  and  $v, w \in C$  it holds

$$f(\lambda v + (1 - \lambda)w) \leq \lambda f(v) + (1 - \lambda)f(w).$$

**Theorem A.2.3.** Let  $V$  be a Hilbert space and  $C \subset V$  nonempty, closed, and convex. Then there exists a unique mapping  $\mathcal{P}_C : V \rightarrow C$  given by

$$\|v - \mathcal{P}_C(v)\| = \text{dist}(v, C) = \inf_{w \in C} \|v - w\| \quad \forall v \in V.$$

*Beweis.* Let  $v_k \in C$  be a minimizing sequence, e.g.,

$$\|v - v_k\| \rightarrow \inf_{w \in C} \|v - w\| =: d.$$

By the parallelogram law we get for any  $k, l \in \mathbb{N}$

$$\|(v - v_k) - (v - v_l)\|^2 + \|(v - v_k) + (v - v_l)\|^2 = 2(\|v - v_l\|^2 + \|v - v_k\|^2).$$

This yields

$$\|v_k - v_l\|^2 = 2(\|v - v_l\|^2 + \|v - v_k\|^2 - 2\|v - \frac{1}{2}(v_k + v_l)\|^2).$$

Now,  $\frac{1}{2}(v_k + v_l) \in C$  and thus  $\|v - \frac{v_k + v_l}{2}\|^2 \geq d^2$ . Hence, we get

$$\|v_k - v_l\|^2 \leq 2(\|v - v_l\|^2 + \|v - v_k\|^2 - 2d^2) \rightarrow 0 \quad (k, l \rightarrow \infty).$$

Now,  $V$  is complete, and  $C$  is closed, hence there exists  $v_\infty \in C$  with  $v_k \rightarrow v_\infty$ . By continuity of the norm it follows

$$\|v - v_\infty\| = d.$$

To see uniqueness, assume that there is some  $v'_\infty \in A$  satisfying

$$\|v - v'_\infty\| = d$$

### A. Functional Analytic Background

an application of the parallelogram law yields as above

$$\|v_\infty - v'_\infty\| \leq 2(\|v - v_\infty\|^2 + \|v - v_\infty\|^2 - 2d^2) = 0$$

and hence  $\mathcal{P}_C(v) = v_\infty$ . □

**Lemma A.2.4.** Let  $V$  be a Hilbert space and  $C \subset V$  nonempty, closed, and convex. The element  $\mathcal{P}_C$  in Theorem A.2.3 is given equivalently by the variational inequality

$$(v - \mathcal{P}_C(v), c - \mathcal{P}_C(v)) \leq 0 \quad \forall c \in C.$$

*Beweis.* Let  $c \in C$  be given. Then for any  $\lambda \in (0, 1)$  it holds  $(1 - \lambda)\mathcal{P}_C(v) + \lambda c \in C$  and hence

$$\begin{aligned} \|v - \mathcal{P}_C(v)\|^2 &\leq \|v - ((1 - \lambda)\mathcal{P}_C(v) + \lambda c)\|^2 \\ &= \|v - \mathcal{P}_C(v)\|^2 - 2\lambda(v - \mathcal{P}_C(v), c - \mathcal{P}_C(v)) + \lambda^2\|c - \mathcal{P}_C(v)\|^2. \end{aligned}$$

Dividing by  $\lambda$  and taking the limit  $\lambda \rightarrow 0$  yields the desired variational inequality.

To see the converse, if the variational inequality holds we have

$$\begin{aligned} \|v - c\|^2 &= \|v - \mathcal{P}_C(v) + \mathcal{P}_C(v) - c\|^2 \\ &= \|v - \mathcal{P}_C(v)\|^2 + 2(v - \mathcal{P}_C(v), \mathcal{P}_C(v) - c) + \lambda^2\|c - \mathcal{P}_C(v)\|^2 \\ &\geq \|v - \mathcal{P}_C(v)\|^2. \end{aligned}$$

□

## A.3. Linear Operators

In the following, let  $V, W$  be normed linear spaces.

**Definition A.3.1.** A mapping  $A: V \rightarrow W$  is called linear or *linear operator* if for any  $\lambda \in \mathbb{R}$  and  $v_1, v_2 \in V$  it holds

$$A(\lambda v_1 + v_2) = \lambda A v_1 + A v_2.$$

If  $A: V \rightarrow \mathbb{R}$  is linear, then it is called a *linear functional*.

**Definition A.3.2.** A linear operator  $A: V \rightarrow W$  is *bounded* if there exists a constant  $c_A$  independent of  $v$  such that

$$\|A v\|_W \leq c_A \|v\|_V \quad \forall v \in V.$$

The set of all bounded linear operators from  $V$  to  $W$  is denoted by

$$\mathcal{L}(V, W).$$

If  $V = W$  we write  $\mathcal{L}(V) := \mathcal{L}(V, V)$ .

**Theorem A.3.3.**  $\mathcal{L}(V, W)$  is a normed linear space (exercise) with norm

$$\|A\| = \|A\|_{\mathcal{L}(V, W)} = \sup_{\|v\|_V=1} \|Av\|_W.$$

$\mathcal{L}(V, W)$  is a Banach space if  $W$  is a Banach space.

**Theorem A.3.4.** A linear operator  $A: V \rightarrow W$  is bounded if and only if it is continuous.

**Example A.3.5** (Multiplication operator). Consider  $V = C[0, 1]$ . Let  $f \in V$  be an arbitrary function. Then we define a linear operator  $A$  by

$$(Ag)(x) = f(x)g(x) \quad \forall x \in [0, 1], g \in C[0, 1].$$

The operator  $A$  is an element in  $\mathcal{L}(C[0, 1], C[0, 1])$  because

$$\|Ag\|_\infty = \max_{x \in [0, 1]} |f(x)g(x)| \leq \max_{x \in [0, 1]} |f(x)| \max_{x \in [0, 1]} |g(x)| = \|f\|_\infty \|g\|_\infty.$$

Hence  $\|A\| \leq \|f\|_\infty$ . In fact  $\|A\|_\infty = \|f\|_\infty$  (exercise).

**Definition A.3.6.** The set  $\mathcal{L}(V, \mathbb{R})$  is called (topological) dual space to  $V$  and is denoted by  $V^*$ .

**Theorem A.3.7** (Riesz representation). Let  $V$  be a Hilbert space with scalar product  $(\cdot, \cdot)$ . Then for any bounded linear functional  $F \in V^*$  there exists a unique element  $f \in V$  such that

$$(f, v) = F(v) \quad \forall v \in V.$$

*Beweis.* Assume w.l.o.g. that  $F \neq 0$ . The set  $C = \mathcal{N}(F) = \{x \in H \mid F(x) = 0\}$  is clearly convex, closed and nonempty (Exercise). Hence by Theorem A.2.3 the projection  $\mathcal{P}_C$  is well defined.

### A. Functional Analytic Background

We choose some  $w \in V$  with  $F(w) = 1$  and define  $\tilde{f} = w - \mathcal{P}_C(w)$ . We obtain  $F(\tilde{f}) = 1$  and hence  $\tilde{f} \neq 0$ . From Lemma A.2.4 (and the fact that  $C$  is a subspace) we obtain

$$(\tilde{f}, w') = 0 \quad \forall w' \in C = \mathcal{N}(F).$$

Now let  $v \in V$  be arbitrary, then  $v - F(v)\tilde{f} \in \mathcal{N}(F)$  and thus

$$(\tilde{f}, v) = (\tilde{f}, v - F(v)\tilde{f}) + (\tilde{f}, F(v)\tilde{f}) = F(v)\|\tilde{f}\|^2.$$

Hence  $f = \frac{\tilde{f}}{\|\tilde{f}\|^2}$  yields the desired.

To see uniqueness we use the non degeneracy of the scalar product. Let  $\hat{f} \in V$  be a second element with  $(\hat{f}, v) = F(v)$  for all  $v \in V$ . Then it holds

$$0 = F(f - \hat{f}) - F(f - \hat{f}) = (f, f - \hat{f}) - (\hat{f}, f - \hat{f}) = \|f - \hat{f}\|^2.$$

□

For more general cases one may need the following generalization

**Theorem A.3.8 (Lax-Milgram).** *Let  $V$  be a Hilbert space and  $a: V \times V \rightarrow \mathbb{R}$  a bilinear form which satisfies for some constants  $\alpha, \beta > 0$  that for any  $u, v \in V$  it holds*

1.  $|a(u, v)| \leq \beta \|u\| \|v\|$  (Continuity)
2.  $\alpha \|v\|^2 \leq a(v, v)$  (Coercivity).

*Then there exists a unique bijective  $A \in \mathcal{L}(V)$  such that*

$$(Au, v) = a(u, v) \quad \forall u, v \in V.$$

*Further it holds*

$$\|A\| \leq \beta, \quad \|A^{-1}\| \leq \frac{1}{\alpha}$$

*Beweis.* Exercise.

□

For reasons of notational simplicity it is sometimes convenient to define a duality pairing  $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{B, B^*}: B \times B^* \rightarrow \mathbb{R}$  for arbitrary Banach spaces  $B$  by the following definition

$$\langle v, f \rangle = f(v) \quad \forall v \in B, f \in B^*.$$

Before we conclude this section we will define the bidual to a Banach space  $B$  as  $B^{**} = (B^*)^* = \mathcal{L}(B^*, \mathbb{R})$ . It is then clear, that we can define a mapping  $i \in \mathcal{L}(B, B^{**})$  as follows:

$$\langle v', i(v) \rangle_{B^*, B^{**}} = \langle v, v' \rangle_{B, B^*} = v'(v) \quad \forall v \in B, v' \in B^*.$$

In fact the mapping  $i$  is an isometry, e.g.,  $\|i(v)\| = \|v\|$ .



**Definition A.3.9.** A Banach space  $B$  is called reflexive if  $i$  is surjective, e.g.,  $B$  and  $B^{**}$  are isometrically isomorphic.

**Example A.3.10.** 1. Every Hilbert space is reflexive.

2. The spaces  $l_p$  and  $L^p(\Omega)$  ( $\Omega \subset \mathbb{R}^n$ ) are reflexive if  $1 < p < \infty$ .

3. The space  $C(\overline{\Omega})$ ,  $L^1(\Omega)$ , and  $L^\infty(\Omega)$  are not reflexive if  $\Omega \subset \mathbb{R}^n$  is open and nonempty.

## A.4. Adjoint

Let  $U, V$  be Banach spaces. Let  $A \in \mathcal{L}(U, V)$  be given. Then we can define for any  $f \in V^*$  an element  $g \in U^*$  by setting

$$g(u) = f(Au) \quad \forall u \in U$$

clearly  $g$  is linear and because of

$$|g(u)| \leq \|f\|_{V^*} \|A\|_{\mathcal{L}(U, V)} \|u\|_U$$

it holds  $g \in U^*$  with

$$\|g\|_{U^*} \leq \|f\|_{V^*} \|A\|_{\mathcal{L}(U, V)}.$$

**Definition A.4.1.** For any  $A \in \mathcal{L}(U, V)$  the above mapping  $f \mapsto g = f \circ A$  defines a linear operator  $A^*: V^* \rightarrow U^*$ ,  $f \mapsto f \circ A$  the *adjoint* or *dual* operator to  $A$ .

It is given by the relation

$$\langle Au, f \rangle_{V \times V^*} = \langle u, A^* f \rangle_{U \times U^*} \quad \forall u \in U, f \in V^*.$$

Moreover it holds  $\|A^*\|_{\mathcal{L}(V^*, U^*)} = \|A\|_{\mathcal{L}(U, V)}$ .

For Hilbert spaces there exists a natural choice for the duality pairing given by Theorem A.3.7. This gives the following

**Definition A.4.2.** Let  $U$  and  $V$  be Hilbert spaces and  $A \in \mathcal{L}(U, V)$ . Then we define the

## A. Functional Analytic Background

Hilbert space adjoint  $A^* \in \mathcal{L}(V, U)$  to  $A$  by

$$(Au, v)_V = (u, A^*v)_U \quad \forall u \in U, v \in V.$$

Note that the notation is sloppy as there is a difference between the adjoint and the Hilbert space adjoint (Which one?). We will ignore this because it will be clear from the context which one is meant.

### A.5. Weak Convergence

We have seen in some examples that the usual finite dimensional proof for existence of minimizers fails, because in infinite dimensional spaces the closed unit ball is not compact. To circumvent this we will derive a new notion of convergence for these spaces.

**Definition A.5.1** (Weak convergence). Let  $B$  be a Banach space. A sequence  $v_k \in B$  converges weakly to some  $v \in B$  if

$$\langle v_k, f \rangle \rightarrow \langle v, f \rangle \quad \forall f \in B^*.$$

This is written as  $v_k \rightharpoonup v$ .

**Definition A.5.2** (Weak\* convergence). Let  $B$  be a Banach space. A sequence  $f_k \in B^*$  converges weakly\* to some  $f \in B^*$  if

$$\langle v, f_k \rangle \rightarrow \langle v, f \rangle \quad \forall v \in B.$$

This is written as  $f_k \rightharpoonup^* f$ .

**Example A.5.3.** Consider the Hilbert space  $L^2(0, 2\pi)$  with the sequence

$$v_k = \sin(kx)$$

Then we have for any  $f \in L^2(0, 2\pi)$  that

$$(f, v_k) = \int_0^{2\pi} f(x) \sin(kx) dx \rightarrow 0$$

by Parseval's identity (The scalar product is up to a constant the  $k$ -th Fourier coefficient.

This means  $v_k \rightharpoonup 0$ . Moreover we have that

$$\|v_k\|^2 = \pi.$$

And hence we immediately obtain that the norm is not weakly continuous (with the obvious definition) because

$$0 = \|0\| \neq \lim_{k \rightarrow \infty} \|v_k\| = \sqrt{\pi}.$$

The following important properties hold:

**Lemma A.5.4.** *Let  $B$  be a Banach and  $H$  be a Hilbert space.*

1. *If  $B$  is reflexive then weak and weak\* convergence coincide on  $B^*$ .*
2. *If  $v_k \rightarrow v$  in  $B$  (strong convergence) then  $v_k \rightharpoonup v$  (weak convergence).*
3. *If  $v_k \rightharpoonup v$  in  $B$  then  $\|v_k\| \leq C < \infty$ .*
4. *If  $u_k \rightarrow u$  and  $v_k \rightharpoonup v$  in  $H$  then*

$$(u_k, v_k) \rightarrow (u, v).$$

*(The analog property holds for Banach spaces)*

5. *If  $v_k \rightharpoonup v$  in  $H$  and  $\|v_k\| \rightarrow \|v\|$  then  $v_k \rightarrow v$ .*

**Definition A.5.5.** A set  $M \subset B$  of a Banach space  $B$  is *weak sequentially closed* if for any sequence  $v_k \in M$  with weak limit  $v \in B$  it holds  $v \in M$ .

A set  $M \subset B$  of a Banach space  $B$  is *weak sequentially compact* if for any sequence  $v_k \in M$  there exists a subsequence  $v_{k_l}$  and an element  $v \in M$  with  $v_{k_l} \rightharpoonup v$ .

We have the following important properties:

**Theorem A.5.6.** *Let  $B$  be a Banach space.*

1. *If  $M \subset B$  is closed and convex then  $M$  is weak sequentially closed.*
2. *If  $B$  is reflexive and  $M \subset B$  bounded, closed, and convex then  $M$  is weak sequential compact.*
3. *If  $B$  is separable then the ball  $\overline{B_1(0)}$  in  $B^*$  is weakly\* sequentially compact.*

**Theorem A.5.7.** Let  $B$  be a Banach space and  $J : B \rightarrow \mathbb{R}$  be a continuous and convex functional. Then  $J$  is weakly lower semicontinuous, e.g., if  $v_k \rightharpoonup v$  then

$$\liminf_{k \rightarrow \infty} J(v_k) \geq J(v).$$

**Remark A.5.8.** In particular the norm  $\|\cdot\| : B \rightarrow \mathbb{R}$  is convex and continuous and thus weakly lower semicontinuous.

**Definition A.5.9.** Let  $V, W$  be two Banach spaces. We call a map  $A \in \mathcal{L}(V, W)$  compact, if for any bounded sequence  $v_k \in V$  there exists a subsequence of  $Av_k$  which converges strongly.

In particular, if  $A \in \mathcal{L}(V, W)$  is compact then if  $v_k \rightarrow v$  in  $V$  it holds  $Av_k \rightarrow Av$ . (Exercise)

## A.6. Lebesgue and Sobolev spaces

### A.6.1. Domain regularity

Throughout this lecture, we consider  $\Omega \subset \mathbb{R}^n$  to be a bounded domain, i.e.,  $\Omega$  is open, connected, and nonempty.

**Definition A.6.1.** We say that a bounded domain  $\Omega$  with boundary  $\Gamma$  is of class  $C^{k,1}$  if the boundary of  $\Omega$  is locally the graph of a  $C^{k,1}$  function such that  $\Omega$  is locally on one side of the boundary.

For the technical detail see, e.g., [Adams and Fournier \[2003\]](#), [Wloka \[1982\]](#). For our purposes the most important regularity will be Lipschitz domains ( $C^{0,1}$ ) domains.

An other important class is that of polygonally bounded domains, e.g., those domains whose boundary is a polygon, in 2d or polyhedral boundaries in 3d.

### A.6.2. Lebesgue spaces

**Definition A.6.2.** On a domain  $\Omega$  we define the spaces  $L^p(\Omega)$  as follows

$$L^p(\Omega) = \{f \mid f \text{ is Lebesgue-measurable, } \|f\|_p^p = \int_{\Omega} |f(x)|^p dx < \infty\}, \quad 1 \leq p < \infty$$

$$L^\infty(\Omega) = \{f \mid f \text{ is Lebesgue-measurable, } \|f\|_\infty = \text{ess sup}_{x \in \Omega} |f(x)| < \infty\}$$

The space  $L^1_{\text{loc}}(\Omega)$  is defined as those functions that are in  $L^1(C)$  for any compact set  $C \subset \Omega$ .

**Remark A.6.3.** Note that from now on if we write down a norm or scalar product for functions without any indices this will be the  $L^2(\Omega)$  norm or scalar product.

**Theorem A.6.4.** The dual space to  $L^p(\Omega)$  for  $1 \leq p < \infty$  is isometrically isomorphic to  $L^{p'}(\Omega)$  where  $\frac{1}{p} + \frac{1}{p'} = 1$ .

Further they are reflexive for  $1 < p < \infty$ .

**Theorem A.6.5** (Lebesgue differentiation theorem). Let  $f \in L^1(\Omega)$  for some bounded domain  $\Omega \subset \mathbb{R}^n$ . Then for almost all  $x \in \Omega$  it holds

$$f(x) = \lim_{r \rightarrow 0} \frac{1}{|B_r(x)|} \int_{B_r(x)} f(y) dy.$$

*Beweis.* For a proof see [Giaquinta et al. \[1998\]](#). □

### A.6.3. Sobolev spaces

We recall that  $C_0^\infty(\Omega)$  denotes the space of all arbitrarily times differentiable functions with compact support in  $\Omega$ .

Motivated by the formula for partial integration of smooth functions we define the following.

**Definition A.6.6** (Weak derivatives). Let  $u \in L^1_{\text{loc}}(\Omega)$  and  $\alpha$  a multi-index. If there exists

### A. Functional Analytic Background

a function  $w \in L^1_{\text{loc}}(\Omega)$  such that

$$\int_{\Omega} u(x) \partial^{\alpha} \varphi(x) dx = (-1)^{|\alpha|} \int_{\Omega} \varphi(x) w(x) dx \quad \forall \varphi \in C_0^{\infty}(\Omega)$$

then we call  $w$  the *weak derivative* of order  $\alpha$  of  $u$  and write  $w = \partial^{\alpha} u$ .

As usual we denoted by  $\partial^{\alpha}$  the term

$$\partial^{\alpha} = \frac{\partial^{\alpha_1}}{\partial x_1} \cdots \frac{\partial^{\alpha_n}}{\partial x_n}.$$

**Example A.6.7.** The function  $u(x) = |x|$  on  $\Omega = (-1, 1)$  has the weak derivative

$$\partial^1 u(x) = u'(x) = \begin{cases} -1 & x \in (-1, 0), \\ 1 & x \in (0, 1). \end{cases}$$

However,  $u$  has no higher weak derivatives (Exercise).

We can now define the Sobolev spaces

**Definition A.6.8.** On a bounded domain  $\Omega$  with given  $k \in \mathbb{N}_0$  and  $1 \leq p \leq \infty$ , we define

$$W^{k,p}(\Omega) = \{u \mid \partial^{\alpha} u \in L^p(\Omega)\}.$$

Together with the norm

$$\|u\|_{k,p} = \left( \sum_{|\alpha| \leq k} \|\partial^{\alpha} u\|_p^p \right)^{1/p}$$

for  $1 \leq p < \infty$  and

$$\|u\|_{k,\infty} = \max_{|\alpha| \leq k} \|\partial^{\alpha} u\|_{\infty}$$

they are Banach spaces. For  $1 < p < \infty$  they are reflexive.

**Definition A.6.9.** On a bounded domain  $\Omega$ , we define for  $k \in \mathbb{N}_0$  and  $1 \leq p < \infty$  we define  $W_0^{k,p}(\Omega)$  as the closure of  $C_0^{\infty}(\Omega)$  with respect to the norm  $\|\cdot\|_{k,p}$ .

In the special case that  $p = 2$  these spaces are Hilbert spaces, we write  $H^k(\Omega) = W^{k,2}(\Omega)$  and  $H_0^k(\Omega) = W_0^{k,2}(\Omega)$ .

**Example A.6.10.** In the special case of  $H^1(\Omega)$ , which we will frequently use, the norm has the following form

$$\|u\|_{1,2}^2 = \int_{\Omega} u^2 + |\nabla u|^2 dx$$

The corresponding scalar product is given as

$$(u, v)_{1,2} = \int_{\Omega} uv + \nabla u \cdot \nabla v dx$$

where  $\cdot$  denotes the euclidean scalar product in  $\mathbb{R}^n$ .

**Theorem A.6.11** (Boundary Traces). *Let  $\Omega$  be a bounded Lipschitz domain and  $1 \leq p \leq \infty$ . Then there exists a unique operator  $\tau \in \mathcal{L}(W^{1,p}(\Omega), L^p(\partial\Omega))$  such that*

$$\tau u = u|_{\partial\Omega} \quad \forall u \in W^{1,p}(\Omega) \cap C(\bar{\Omega}).$$

In slightly sloppy notation, we will usually neglect to explicitly write down the operator  $\tau$  when it is clear that the trace is meant.

With these preparations we can give an equivalent definition for the space  $H_0^1(\Omega)$  by the following

$$H_0^1(\Omega) = \{u \in H^1(\Omega) \mid \tau u = 0\}. \quad (\text{A.1})$$

**Theorem A.6.12** (Poincaré inequality). *Let  $\Omega$  be a bounded domain. Then there exists a constant  $c_{\Omega}$  depending on  $\Omega$  only such that*

$$\|u\|^2 = \int_{\Omega} u^2 dx \leq c_{\Omega} \|\nabla u\|^2 = c_{\Omega} \int_{\Omega} |\nabla u|^2 dx \quad \forall u \in H_0^1(\Omega).$$

**Corollary A.6.13.** *The norms  $\|\cdot\|_{1,2}$  and  $\|\nabla \cdot\|$  are equivalent on  $H_0^1(\Omega)$ .*

#### A.6.4. Embedding theorems

Now, in order to define constraints onto our solution variables it is nice to know which Sobolev space are contained in spaces of continuous functions. To see this we will recall several embedding theorems.

## A. Functional Analytic Background

**Theorem A.6.14.** Let  $\Omega \subset \mathbb{R}^n$  be a Lipschitz domain. Let  $k_1, k_2 \in \mathbb{N}_0$  and  $1 \leq p_1, p_2 < \infty$ .  
a) If it holds

$$k_1 - \frac{n}{p_1} \geq k_2 - \frac{n}{p_2}, \quad k_1 \geq k_2$$

then the embedding  $\text{Id}: W^{k_1, p_1}(\Omega) \rightarrow W^{k_2, p_2}(\Omega)$  exists and is continuous (but not compact). This means in particular that there exists a constant  $c$  depending on  $n, \Omega, k_1, k_2, p_1, p_2$  such that

$$\|u\|_{k_2, p_2} \leq c \|u\|_{k_1, p_1}.$$

b) If it holds

$$k_1 - \frac{n}{p_1} > k_2 - \frac{n}{p_2}, \quad k_1 > k_2$$

then the embedding  $\text{Id}: W^{k_1, p_1}(\Omega) \rightarrow W^{k_2, p_2}(\Omega)$  is in addition compact.

**Theorem A.6.15.** Let  $\Omega \subset \mathbb{R}^n$  be a Lipschitz domain. Let  $k_1, k_2 \in \mathbb{N}_0$  with  $k_1 \geq 1$  and  $1 \leq p < \infty$  and  $0 \leq \alpha \leq 1$ . If it holds

$$k_1 - \frac{n}{p} > k_2 + \alpha,$$

then the embedding  $\text{Id}: W^{k_1, p}(\Omega) \rightarrow C^{k_2, \alpha}(\overline{\Omega})$  exists and is continuous and compact.

**Theorem A.6.16.** Let  $\Omega \subset \mathbb{R}^n$  be a Lipschitz domain. Let  $k \in \mathbb{N}$ , then the embedding

$$C^{k, 1}(\overline{\Omega}) \rightarrow W^{k+1, \infty}(\Omega)$$

exists and is an isomorphism.

**Example A.6.17.** We consider again the case of the space  $H^1(\Omega)$ . Now we consider the question whether or not these functions are continuous. To this end we note that the relevant inequality gets

$$1 - \frac{n}{2} > 0.$$

Hence if  $n = 1$  we obtain that all  $H^1(\Omega)$  functions are continuous. However the inequality fails for any dimension  $n \geq 2$ .

In fact for these dimensions functions in  $H^1(\Omega)$  are not necessarily continuous. To see this



let  $\Omega = B_1(0) \subset \mathbb{R}^2$  and consider the function  $u(x) = \ln(\ln(r^{-1}) + 1)$ . A straightforward computation shows that  $u \in H^1(\Omega)$  but clearly  $u \notin C^0(\overline{\Omega})$ .

### *A. Functional Analytic Background*

# Literaturverzeichnis

- R. A. Adams and J. J. F. Fournier. *Sobolev Spaces*, volume 140 of *Pure and Applied Mathematics (Amsterdam)*. Elsevier/Academic Press, Amsterdam, second edition, 2003.
- Hans W. Alt. *Lineare Funktionalanalysis: Eine anwendungsorientierte Einführung*. Springer, Berlin – Heidelberg – New York, 2002.
- Jürgen Appell and Petr P. Zabrejko. *Nonlinear Superposition Operators*. Cambridge University Press, 1990.
- V. I. Bogachev and E. Mayer-Wolf. Some remarks on Rademacher’s theorem in infinite dimensions. *Potential Analysis*, 5:23–30, 1996.
- J. F. Bonnans and A. Shapiro. Optimization problems with perturbations: A guided tour. *SIAM Rev.*, 40(2):228–264, 1998. doi: 10.1137/S0036144596302644.
- J. F. Bonnans and A. Shapiro. *Perturbation Analysis of Optimization Problems*. Springer Series in Operations Research. Springer, 2000.
- S. C. Brenner and L. R. Scott. *The Mathematical Theory of Finite Element Methods*. Springer Verlag, New York, 3. edition, 2008.
- P. G. Ciarlet. *Linear and nonlinear functional analysis with applications*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2013.
- B. Dacorogna. *Direct Methods in the Calculus of Variations*, volume 78 of *Applied Mathematical Sciences*. Springer, second edition, 2008.
- J. W. Daniel. The conjugate gradient method for linear and nonlinear operator equations. *SIAM J. Numer. Anal.*, 4:10–26, 1967.
- J. C. De los Reyes. *Numerical PDE-constrained optimization*. SpringerBriefs in Optimization. Springer, Cham, 2015.
- F. Gaspoz, C. Kreuzer, A. Veiser, and W. Wollner. Quasi-best approximation in optimization with PDE constraints. *Inverse Problems*, 36(1), 2020. doi: 10.1088/1361-6420/ab47f3.
- M. Giaquinta, G. Modica, and J. Souček. *Cartesian Currents in the Calculus of Variations I*. *Ergebnisse der Mathematik und ihrer Grenzgebiete*. Springer, Berlin – Heidelberg – New York, 1. edition, 1998.
- D. Gilbarg and N. S. Trudinger. *Elliptic Partial Differential Equations of Second Order*, volume 224 of *Grundlehren der mathematischen Wissenschaften*. Springer, revised 3. edition, 2001.

- M. Hintermüller, A. Schiela, and W. Wollner. The length of the primal-dual path in Moreau-Yosida-based path-following for state constrained optimal control. *SIAM J. Optim.*, 24(1): 108–126, 2014.
- M. Hinze. A variational discretization concept in control constrained optimization: The linear-quadratic case. *Comp. Optim. Appl.*, 30(1):45–61, 2005.
- M. Hinze, R. Pinnau, M. Ulbrich, and S. Ulbrich. *Optimization with PDE Constraints*, volume 23 of *Mathematical Modelling: Theory and Applications*. Springer, 2009.
- K. Ito and K. Kunisch. *Lagrange Multiplier Approach to Variational Problems and Applications*, volume 15 of *Advances in Design and Control*. Society for Industrial and Applied Mathematics (SIAM), 2008.
- M. Lavrentieff. Sur quelques problèmes du calcul des variations. *Ann. Mat. Pura Appl.*, 4(1): 7–28, 1927.
- Stephen M. Robinson. Stability theory for systems of inequalities. II. Differentiable nonlinear systems. *SIAM J. Numer. Anal.*, 13(4):497–513, 1976.
- W. Rudin. *Functional analysis*. International Series in Pure and Applied Mathematics. McGraw-Hill Inc., 1991.
- A. Schiela and W. Wollner. Barrier methods for optimal control problems with convex nonlinear gradient state constraints. *SIAM J. Optim.*, 21(1):269–286, 2011.
- Anton Schiela. *The Control Reduced Interior Point Method*. Verlag Dr. Hut, München, 2006.
- F. Tröltzsch. *Optimale Steuerung partieller Differentialgleichungen*. Vieweg, 1. edition, 2005.
- M. Ulbrich. *Semismooth Newton Methods for Variational Inequalities and Constrained Optimization Problems*. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics (SIAM), 2011.
- D. Werner. *Funktionalanalysis*. Springer, Berlin – Heidelberg – New York, 2005.
- J. Wloka. *Partielle Differentialgleichungen. Sobolewräume und Randwertaufgabe*. Teubner, Leipzig, 1. edition, 1982.
- Jochem Zowe and Stanislaw Kurcyusz. Regularity and stability for the mathematical programming problem in Banach spaces. *Appl. Math. Optim.*, 5:49–62, 1979.