

# OPTIMIZATION WITH PARTIAL DIFFERENTIAL EQUATIONS

Lecture notes  
Winter term 2022/23

HANNES MEINLSCHMIDT

[meinlschmidt@math.fau.de](mailto:meinlschmidt@math.fau.de)

Professor for Applied Analysis  
Chair for Dynamics, Control and Numerics  
Department of Data Science

Preliminary version  
February 10, 2023



---

# Preface

These lecture notes cover topics suitable for an introductory 5 ECTS (2+1 weekly hours) course about the theory of optimization in function spaces with applications to partial differential equations. Students are assumed to be familiar with general principles of functional analysis and nonlinear optimization.

The goal of the course is to give an abstract treatment of optimization problems posed in infinite-dimensional spaces, which is—as far as time allows it—self-contained in its presentation, giving students a theoretical background for further studies in the field of infinite-dimensional optimization, in particular in connection with optimal control or generally PDE-constrained problems. After completion of this course, students should be able to identify the underlying machinery for practical problems and follow their optimization setup or execute it themselves.

We refer to the excellent textbooks of Brezis [Br10] and Nocedal and Wright [NW06] as well as Ulbrich and Ulbrich [UU12] for references to functional analysis and basic PDE theory, and theory of nonlinear optimization, respectively. Topics of this course are (partially) covered in, for example, the books of Bonnans and Shapiro [BS00], Schirotzek [Sc07] or Hinze, Pinnau, Ulbrich and Ulbrich [HPUU09]. Cross-reading is recommended also with Tröltzsch’s book on optimal control of PDEs [Tr10].

The notes are based on lecture notes of Michael Ulbrich (TU München), Stefan Ulbrich (TU Darmstadt), and Christian Meyer (TU Dortmund). Comments, suggestions, and notification of errors are welcome by eMail to [meinlschmidt@math.fau.de](mailto:meinlschmidt@math.fau.de).

# Contents

<b>1</b>	<b>Motivation and problem setting</b>	<b>1</b>
1.1	Problem setting . . . . .	2
<b>2</b>	<b>Existence of Solutions</b>	<b>5</b>
<b>3</b>	<b>Optimality conditions</b>	<b>12</b>
3.1	Robinson’s constraint qualification . . . . .	20
3.2	The Karush-Kuhn-Tucker conditions . . . . .	25
3.3	Sufficient optimality conditions . . . . .	45



# 1 Motivation and problem setting

Modeling, optimization and numerical simulation of complex systems plays an important role in physics, engineering, mechanics, chemistry, biology, medicine, finance, and in many other disciplines. These models often use ordinary or partial differential equations to describe the system, which means that the resulting optimization problems involve *generically infinite-dimensional* objects such as functions, or in general quantities that live in Banach spaces. (Thus, most often, *function spaces* are involved.) There is a rising interest to solve optimization problems involving such models and there are many fascinating applications. We give a (totally incomplete) selection of examples:

- Optimal space mission design (trajectories),
- optimal control of robot movements,
- identification of model parameters by means of measurements [ $\rightarrow$  *inverse problems*], such as:
  - geological material properties from seismic measurements,
  - data assimilation: initial conditions for weather and climate models from scattered observations of all kinds (weather stations, satellite data, etc.),
  - tomography
  - derivation of parameters (e.g., volatility) in models for option pricing (*Black Scholes* PDE) from market prices ( $\rightarrow$  *mathematical finance*),
- optimal control of actuators, built on the surface of an aircraft, to avoid noise generation, material fatigue, or vortex generation,
- optimal radiation therapy planning,
- optimal design of the body of a ship,
- optimization of the shape of a wing, a turbine blade, etc.,
- optimal control of laser hardening of steel,
- optimal control of crystal pulling (heating, pulling speed),
- and many more.

The aim of this course is to develop a rigorous theory and establish methods for infinite dimensional optimization problems. This provides the appropriate framework for investigating and solving the above problem classes. Of course, the problems above are of very different natures with respect to their dynamics and the rigidity with which they need be treated. We will thus aim for a quite general problem setting in order to encapsulate as many problems as possible, but for which we can also establish a satisfying theory.

An exemplary PDE-constrained problem could look like this:

**Example.** The following is an *optimal boundary control problem* for a *semilinear elliptic* equation:

$$\begin{aligned} \min_{\substack{y: \Omega \rightarrow \mathbb{R} \\ u: \partial\Omega \rightarrow \mathbb{R}}} & \quad \frac{1}{2} \int_{\Omega} (y(x) - y_d(x))^2 dx + \frac{\alpha}{2} \int_{\partial\Omega} u(x)^2 d\mathcal{H}_{n-1}(x) \\ \text{s.t.} & \quad \begin{cases} -\Delta y + y^3 = 0 & \text{on } \Omega, \\ \frac{\partial y}{\partial \nu} + y = u & \text{on } \partial\Omega, \\ a \leq u \leq b & \text{on } \partial\Omega, \end{cases} \end{aligned}$$

for example arising from (stationary) nonlinear heat transfer. The goal is to find an ambient temperature  $u$  such that the temperature profile  $y$  in  $\Omega \subset \mathbb{R}^n$  (e.g., a work piece) is as close as possible to a given desired temperature profile  $y_d$  in an average sense. Thereby, the ambient temperature must stay between bound functions  $a$  and  $b$  which represent e.g. physical or process-related limitations, and a trade-off between optimization accuracy and heating cost is possible by the parameter  $\alpha$ .

We will come back to this example later. The fundamental (theoretical) questions which we consider in this lecture with respect to infinite-dimensional optimization problems such as the one above are as follows:

- Is the problem well posed to begin with, that is: does the problem admit an optimal solution at all, and in which function class?
- If it does admit an optimal solution, can we characterize these? This means, under which condition are there *necessary* and *sufficient* optimality conditions of e.g. first and second order?

The latter optimality conditions are also a very useful starting point for numerical algorithms to find an optimal solution.

## 1.1 Problem setting

The fundamental general problem class is given by

$$\min_{x \in X} f(x) \quad \text{s.t.} \quad G(x) \in K, \tag{P}$$

with the *feasible set* (zulässige Menge)

$$\mathcal{F} := \left\{ x \in X : G(x) \in K \right\} := G^{-1}[K].$$

We pose the following fundamental assumptions on the data in (P):

**Assumption 1.1** (Problem data). The following properties of the data in (P) are valid:

- (a)  $G: X \rightarrow Z$  is a continuous mapping between real Banach spaces,
- (b)  $K \subseteq Z$  is a closed convex set,
- (c)  $f: X \rightarrow \mathbb{R}$  is a lower semicontinuous objective function,
- (d) the feasible set  $\mathcal{F}$  is nonempty.

In (P) we thus minimize with respect to the linear space  $X$ ; additional constraints on the variable  $x$  are realized by the constraint  $G(x) \in K$ . Note that  $\mathcal{F}$  is closed in  $X$  under [Assumption 1.1](#) since  $G$  is continuous and  $K$  is closed. Let us moreover briefly recall the notion of lower semicontinuity.

**Reminder:** We say that the function  $f: X \rightarrow \mathbb{R}$  is *lower semicontinuous* (*unterhalbstetig*) if for all  $x \in X$  and all sequences  $(x_k) \subset X$  there holds

$$x_k \rightarrow x \implies \liminf_{k \rightarrow \infty} f(x_k) \geq f(x), \quad (1.1)$$

or, equivalently, if the level set  $N_f(\alpha) := \{x \in X: f(x) \leq \alpha\}$  (*Niveaumenge*) is sequentially closed for all  $\alpha \in \mathbb{R}$ .

Of course, a continuous function is also lower semicontinuous. Moreover, lower semicontinuity can also be defined for functions  $f: S \rightarrow \mathbb{R}$  defined on a proper subset  $S \subsetneq X$ . One then considers sequences  $S \ni x_k \rightarrow x \in S$ . The equivalent definition via level sets however requires that  $S$  is closed.

The compact notation (P) includes essentially all relevant situations which occur in particular problems. The two most common variations are the following:

1. The *total problem*

$$\min_{x \in X} f(x) \quad \text{s.t.} \quad E(x) = 0, \quad F(x) \in C, \quad x \in M,$$

where the continuous mappings  $E: X \rightarrow Z_1$  and  $F: X \rightarrow Z_2$  take their values in Banach spaces  $Z_1$  and  $Z_2$ , the set  $C \subseteq Z_2$  is a closed convex cone, and  $M \subseteq X$  is a closed convex set. This is a problem of type (P) for  $Z = Z_1 \times Z_2 \times X$  and  $K = \{0_{Z_1}\} \times C \times M$  with  $G(x) := (E(x), F(x), x)$ .

**Reminder:** Recall that a set  $C$  is called *cone* if from  $x \in C$  it follows that  $\lambda x \in C$  for all  $\lambda > 0$ .

2. In the total problem as above, the variable  $x$  often has two components  $x = (y, u) \in X = Y \times U$ , with  $u$  being a *control* (**Steuerung**), design or parameter within the system, and  $y$  being the associated *state* (**Zustand**) of the system. They are linked by the system dynamics (usually: the differential equations) described abstractly by  $E(y, u) = 0$ . The set  $M$  is then also decoupled into  $M = M_y \times M_u$  with  $M_y \subseteq Y$  and  $M_u \subseteq U$ . Problems of this particular structure are called *optimal control* problems.

In this context, one often encounters optimal control problems in *reduced form* (as opposed to the *total* problem above), that is,

$$\min_{u \in U} j(u) \quad \text{s.t.} \quad H(u) \in C, \quad u \in U_{\text{ad}}.$$

These are obtained from the total optimal control problem by eliminating the state equation  $E(y, u) = 0$ . This is done by showing that for each  $u$  there exists a unique state  $y = y(u)$  such that  $E(y(u), u) = 0$ , and inserting the so-called *control-to-state operator* (**Steuerungs-Zustands-Operator**)  $y(u)$  for  $y$ , that is,

$$j(u) := f(y(u), u), \quad H(u) := F(y(u), u), \quad U_{\text{ad}} := M_u \cap y^{-1}[M_y]$$

where  $M = M_y \times M_u \subseteq Y \times U$ . Of course, the state equation  $E(y(u), u) = 0$  is then superfluous because it is always satisfied by construction. Such a reduced problem is of formally easier nature, but with (much) more involved functions  $j, H$  due to the control-to-state operator  $u \mapsto y(u)$ .

**Remark 1.2.** The conditions  $F(x) \in C$  or  $H(u) \in C$  with  $C$  denoting a closed convex cone can be interpreted to describe *abstract inequality constraints*. If for instance  $Z_2 = \mathbb{R}^m$  and  $C = (-\infty, 0]^m$ , then the constraint  $F(x) \in C$  is the same as  $F(x) \leq 0$  (component-wise), which corresponds to the standard inequality constraints in nonlinear programming. Note that there is an order structure on  $C$  in this case. This additional structure is quite helpful in several situations. We will come back to this later.

**Remark 1.3.** (a) As in [Assumption 1.1](#), we will *always* assume all occurring Banach spaces to be real in the following, even if it is not mentioned explicitly.

(b) Over the course of this lecture, we will also consider first order necessary optimality conditions for problems of type **(P)**. Differentiability assumptions on  $f$  and  $G$  will then be stated as required.



## 2 Existence of Solutions

We first discuss the question of existence of solutions to the problem (P). The easiest case is that of a compact feasible set.

**Theorem 2.1.** Consider (P) and let [Assumption 1.1](#) be satisfied. Assume further that there exist a sequentially precompact set  $C \subseteq \mathcal{F}$  and  $x_0 \in C$  such that  $f(x) \geq f(x_0)$  for all  $x \in \mathcal{F} \setminus C$ . Then (P) possesses a global solution.

*Proof.* By assumption, no point in  $\mathcal{F} \setminus C$  provides a lower objective function value than  $x_0 \in C$ . Hence, we can restrict our search for the minimum on the set  $C$ , which is sequentially precompact and nonempty. Now, consider a minimizing sequence  $(x_k) \subset C$  satisfying  $f(x_k) \rightarrow f^* := \inf_{x \in \mathcal{F}} f(x)$ . Due to sequential precompactness of  $C$ , there exists a subsequence, again denoted by  $(x_k)$ , that converges to a limit  $\bar{x}$ . The feasible set  $\mathcal{F}$  is closed, hence there holds  $\bar{x} \in \mathcal{F}$ . Now, since  $f$  is lower semicontinuous, we have

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

and thus  $\bar{x} \in \mathcal{F}$  is a solution of (P). □

**Remark 2.2.** Note that if  $\mathcal{F}$  is already compact itself, then the choice  $C = \mathcal{F}$  is perfectly valid in the foregoing theorem.

In infinite-dimensional spaces, compactness is a strong requirement, often: *too* strong. Recall that the closed unit ball  $B_X(0)$  in a Banach space  $X$  is compact *if and only if*  $X$  is finite-dimensional! On the other hand, *weak* sequential compactness is often verifiable.

**Reminder:** We recall that  $(x_k)$  converges weakly to  $x$  in  $X$ , written  $x_k \rightharpoonup x$ , iff

$$x'(x_k) =: \langle x', x_k \rangle \rightarrow \langle x', x \rangle := x'(x) \quad \text{for all } x' \in X^*,$$

where  $X^* = \mathcal{L}(X; \mathbb{R})$  is the space of continuous linear functionals on  $X$ . A set  $M$  is weakly sequentially compact (**schwach folgenkompakt**) if every sequence  $(x_k) \subset M$  admits a weakly convergent subsequence whose limit is again in  $M$ , and a function  $f: X \rightarrow \mathbb{R}$  is weakly lower semicontinuous if  $x_k \rightharpoonup x$  implies  $\liminf_{k \rightarrow \infty} f(x_k) \geq f(x)$ .

We now want to replace the compactness requirement in [Theorem 2.1](#) by weak sequential compactness. This works particularly well in *reflexive* Banach spaces, since these are precisely the ones for which the closed unit ball  $B_X(0)$  is weakly sequentially compact.

**Reminder:** A Banach space  $X$  is said to be *reflexive* if the canonical injection

$$J: X \rightarrow (X^*)^* := X^{**}, \quad \langle Jx, x' \rangle := \langle x', x \rangle \quad \text{for all } x' \in X^*$$

is surjective (and thus a bijection).

We need the following facts from functional analysis:

**Proposition 2.3** (See [Br10, Ch. 3]). *Let  $X$  be a Banach space.*

- (a) *Every closed convex subset of  $X$  is weakly sequentially closed.*
- (b) *The space  $X$  is reflexive if and only if every bounded sequence contains a weakly convergent subsequence.*
- (c) *If  $f: X \rightarrow \mathbb{R}$  is convex and lower semicontinuous, then  $f$  is also weakly lower semicontinuous.*
- (d) *Any weakly convergent sequence in  $X$  is bounded.*

Together, the first and second statement in [Proposition 2.3](#) in particular imply that the closed unit ball  $\overline{B_X(0)}$  in a reflexive Banach space  $X$  is weakly sequentially compact.

Note moreover that the third assertion is derived from the first since every level set  $N_f(\alpha)$  of  $f$  is convex (convexity of  $f$ ) and closed (lower semicontinuity) and thus also weakly closed, hence weakly sequentially closed. Therefore,  $f$  is weakly lower semicontinuous. A prime example of a weakly lower semicontinuous function is the norm function  $x \mapsto \|x\|_X$ .

The most commonly encountered reflexive Banach space class are Hilbert spaces. Also important, but for  $p \neq 2$  non-Hilbert, are the Lebesgue spaces  $L^p$  for  $1 < p < \infty$ — $L^1$  and  $L^\infty$  are *not* reflexive on nontrivial measure spaces—and the associated Sobolev spaces  $W^{k,p}$ , where  $k \in \mathbb{N}$ . The fact that these spaces are reflexive plays a huge role in their prevalence in functional analytic methods for PDEs and will also become extremely helpful in the context of optimization theory as we develop in this lecture.

We obtain the following existence result for [\(P\)](#):

**Theorem 2.4.** *Consider [\(P\)](#) and let the following assumptions be satisfied:*

- (a)  *$G: X \rightarrow Z$  is weakly sequentially continuous from the reflexive Banach space  $X$  to the Banach space  $Z$ , i.e., if  $x_k \rightharpoonup x$  in  $X$ , then  $G(x_k) \rightarrow G(x)$  in  $Z$ ,*
- (b)  *$f: X \rightarrow \mathbb{R}$  is weakly lower semicontinuous,*
- (c)  *$K \subseteq Z$  is closed and convex,*

(d) there exist a bounded set  $C \subseteq \mathcal{F}$  and  $x_0 \in C$  such that  $f(x) \geq f(x_0)$  for all  $x \in \mathcal{F} \setminus C$ .

Then the problem (P) possesses a global solution.

*Proof.* By assumption, analogously to the proof of [Theorem 2.1](#), there exists a minimizing sequence  $(x_k) \subset C$  with  $f(x_k) \rightarrow f^* := \inf_{x \in \mathcal{F}} f(x)$ . Since  $C$  is bounded and  $X$  is reflexive, there exists a subsequence, again denoted by  $(x_k)$ , that converges weakly to a limit  $\bar{x} \in X$ . Therefore, by weak sequential continuity,  $G(x_k) \rightharpoonup G(\bar{x})$  in  $Z$ . The closed convex set  $K$  is weakly closed in  $Z$ , and thus  $K \ni G(x_k) \rightharpoonup G(\bar{x}) \in K$ . We conclude  $\bar{x} \in \mathcal{F}$ . Now, since  $f: X \rightarrow \mathbb{R}$  is weakly lower semicontinuous, we obtain

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

and hence  $\bar{x} \in \mathcal{F}$  is a solution of (P). □

**Remark 2.5.** Inspecting the proof of [Theorem 2.4](#), one observes that we could have replaced the assumptions on  $G$  and  $K$  by the general assumption that  $\mathcal{F}$  be weakly closed. As seen in the proof, the assumptions posed are in fact sufficient for  $\mathcal{F}$  to be weakly closed. Sometimes it is more convenient to show this more general condition.

Of course, [Remark 2.2](#) applies again for [Theorem 2.4](#). We note that weak continuity, as assumed for  $G$  in [Theorem 2.4](#), is a delicate topic when nonlinear functions are involved. Generally, it is only to be expected to hold if the mapping under consideration is in fact (affine-) linear and continuous (see the exercises), or if there is *compactness* involved.

**Reminder:** A continuous linear operator  $A \in \mathcal{L}(X; Y)$  between Banach spaces  $X$  and  $Y$  is said to be *compact* if it maps bounded sets in  $X$  to precompact sets in  $Y$ .

Compactness is a very useful technique to transfer properties for the weak topology to properties for the norm topology, as the following lemma demonstrates.

**Lemma 2.6.** *Let  $A \in \mathcal{L}(X; Y)$  be a compact linear operator between the Banach spaces  $X$  and  $Y$ . If  $(x_k) \subset X$  converges weakly to  $x$  in  $X$ , then  $(Ax_k) \subset Y$  converges strongly to  $Ax$  in  $Y$ .*

*Proof.* Note first that weak convergence  $x_k \rightharpoonup x$  in  $X$  also implies  $Ax_k \rightharpoonup Ax$  in  $Y$  (why?).

Assume that  $(Ax_k)$  does not converge in norm to  $Ax$  in  $Y$ . Then there is a sufficiently small neighborhood  $\mathcal{U}$  of  $Ax$  in  $Y$  and a subsequence  $(Ax_{k_\ell})_\ell$  such that  $Ax_{k_\ell} \notin \mathcal{U}$  for all  $\ell \in \mathbb{N}$ .

On the other hand, the weakly convergent sequence  $(x_k)$  is bounded in  $X$ , hence  $(Ax_k)$  is a precompact set in  $Y$ . In particular, the sequence  $(Ax_{k_\ell})_\ell$  must admit a norm convergent subsequence, still denoted by  $(Ax_{k_\ell})_\ell$ . Denoting the limit of that subsequence by  $y$ , we must have  $y \neq Ax$  by construction. But then the sequence  $(Ax_{k_\ell})_\ell$  also converges weakly to  $y \neq Ax$ . Since  $(Ax_{k_\ell})_\ell$  is a subsequence of  $(Ax_k)$ , this is a contradiction to  $Ax_k \rightharpoonup Ax$ .  $\square$

The ansatz to transfer properties for the weak topology to the norm topology is most often used for embeddings.

**Reminder:** A Banach space  $X$  is said to be *embedded* (*eingebettet*) into another Banach space  $X_0$  if there exists a continuous linear injective operator  $i: X \rightarrow X_0$ , the *embedding*, and we write  $X \hookrightarrow X_0$  in this case.

Since an embedding  $i: X \rightarrow X_0$  is by definition injective, it is often most useful to identify  $X$  with its image  $i(X) \subset X_0$ . In most cases, the embedding operator  $i$  is in fact given by the identity mapping  $i = \text{id}: X \rightarrow X_0$ . In this case, we usually do not refer to  $i$  explicitly. (It is sometimes reasonable to do so, however.)

**Definition 2.7.** A Banach space  $X$  is *compactly embedded* in the Banach space  $X_0$  if  $X \hookrightarrow X_0$  and every bounded sequence in  $X$  contains a (strongly) convergent subsequence in  $X_0$ . Equivalently, the embedding operator  $i$  is a compact linear operator from  $X$  to  $X_0$ .

**Corollary 2.8.** *If the Banach space  $X$  is compactly embedded in the Banach space  $X_0$ , then every weakly convergent sequence  $(x_k)$  with limit  $x$  in  $X$  is norm convergent to the same limit in  $X_0$  up to identification of  $X$  with  $i(X) \subset X_0$ .*

*Proof.* Apply [Lemma 2.6](#) to the embedding operator  $i: X \rightarrow X_0$ .  $\square$

Using the foregoing corollary, we observe that by combining compactness and continuity properties for the norm topology, we can generate continuity properties for the weak topology. In the context of [Theorem 2.4](#), we have in particular:

- If the embedding  $X \hookrightarrow X_0$  is compact and  $f_0: X_0 \rightarrow \mathbb{R}$  is lower semicontinuous, then  $f := (f_0 \circ i): X \rightarrow \mathbb{R}$  is weakly lower semicontinuous.

In fact,  $x_k \rightharpoonup x$  in  $X$  implies  $i(x_k) \rightarrow i(x)$  in  $X_0$  and  $\liminf_{k \rightarrow \infty} f_0(i(x_k)) \geq f_0(i(x))$ .

- If the embedding  $X \hookrightarrow X_0$  is compact and  $G_0: X_0 \rightarrow Z$  is sequentially continuous from the norm topology of  $X_0$  to the weak topology of  $Z$  (*strong-weak continuous*), then  $G := (G_0 \circ i)$  is weakly sequentially continuous from  $X$  to  $Z$ .

This follows from

$$x_k \rightharpoonup x \text{ in } X \quad \implies \quad i(x_k) \rightarrow i(x) \text{ in } X_0 \quad \implies \quad G_0(i(x_k)) \rightharpoonup G_0(i(x)) \text{ in } Z.$$

October 26, 2022

← 2

## Example

We pick up the example for an optimal control problem from the introduction, to which we wish to apply the above existence theory.

**Example 2.9.** Consider the following optimal boundary control problem for a semilinear elliptic equation:

$$\begin{aligned} \min_{\substack{y: \Omega \rightarrow \mathbb{R} \\ u: \partial\Omega \rightarrow \mathbb{R}}} J(y, u) &:= \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L^2(\partial\Omega)}^2 \\ \text{s.t.} \quad &\begin{cases} -\Delta y + y^3 = 0 & \text{on } \Omega, \\ \frac{\partial y}{\partial \nu} + y = u & \text{on } \partial\Omega, \\ a \leq u \leq b & \text{on } \partial\Omega, \end{cases} \end{aligned}$$

where  $\Omega \subset \mathbb{R}^n$ , for  $n = 2$  or  $n = 3$ , is an open and bounded Lipschitz domain,  $y_d \in L^2(\Omega)$ ,  $\alpha \geq 0$ , and, with the surface measure  $\omega$  on  $\partial\Omega$ ,

$$a, b \in L^2(\partial\Omega) \quad \text{with} \quad a \leq b \quad \omega\text{-a.e. in } \partial\Omega.$$

Let  $U = L^2(\partial\Omega)$  and  $Y = H^1(\Omega)$ , and

$$U_{\text{ad}} = \left\{ u \in U : a \leq u \leq b \quad \omega\text{-a.e. in } \partial\Omega \right\}.$$

We consider weak solutions of the state equation encoded in  $E(y, u) = 0$ , i.e.,

$$E: Y \times U \rightarrow Y^*,$$

$$\langle E(y, u), v \rangle_{Y^*, Y} := (\nabla y, \nabla v)_{L^2(\Omega)^n} + (y^3, v)_{L^2(\Omega)} + (y - u, v)_{L^2(\partial\Omega)} \quad \text{for all } v \in Y.$$

This formulation is obtained by testing (that is, multiplying and integrating) the PDE with  $v \in H^1(\Omega) = Y$  and integrating by parts. Note that, by Sobolev embeddings, we have  $H^1(\Omega) \hookrightarrow L^6(\Omega)$  for  $n \leq 3$  and thus  $y^3 \in L^2(\Omega)$ . Moreover, for  $w \in H^1(\Omega)$

there is a well defined notion of the restriction  $w|_{\partial\Omega} \in L^2(\partial\Omega)$  and  $w \mapsto w|_{\partial\Omega}$  is continuous between these spaces. (Why is this mapping not an embedding?)

One can show that for each  $u \in U$ , the state equation possesses a unique solution  $y(u) \in Y$ . This defines the control-to-state operator  $U \ni u \mapsto y(u) \in Y$ . Furthermore, there holds

$$\|y(u)\|_{H^1(\Omega)} = \|y(u)\|_Y \leq C\|u\|_U = \|u\|_{L^2(\partial\Omega)}. \quad (2.1)$$

with a constant  $C > 0$  independent of  $u$ . The estimate (2.1) follows by choosing  $v = y(u)$  in the weak formulation, which gives:

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

By the generalized Friedrichs' inequality, there holds

$$\|y\|_{L^2(\Omega)}^2 \leq C(\Omega) \left[ \|\nabla y\|_{L^2(\Omega)^n}^2 + \|y\|_{L^2(\partial\Omega)}^2 \right] \quad (2.2)$$

with a constant  $C(\Omega)$  independent of  $y$ , from which now (2.1) follows. (Work this out!)

In the terminology of (P), we put

$$X = Y \times U, \quad Z = Y^* \times U, \quad K = \{0_{Y^*}\} \times U_{\text{ad}},$$

with  $G(y, u) := (E(y, u), u)$  and  $f := J$ . We verify the requirements of [Theorem 2.4](#) to show that there exists an optimal solution pair  $(y, u) \in X = H^1(\Omega) \times L^2(\partial\Omega)$ . There holds:

(a) The mapping  $E: Y \times U \rightarrow Y^*$  is weakly sequentially continuous:

In fact, the operator  $(y, u) \mapsto E(y, u) - (y^3, \cdot)_{L^2(\Omega)}$  is linear and continuous from  $Y \times U$  to  $Y^*$ , hence weakly sequentially continuous. (Work this out!) It thus remains to show that  $Y \ni y \mapsto (y^3, \cdot)_{L^2(\Omega)} \in Y^*$  is weakly sequentially continuous.

Firstly,  $Y = H^1(\Omega)$  is compactly embedded into  $L^5(\Omega)$  for  $n = 2, 3$ . Further, the function  $\varphi(t) := t^3$  clearly satisfies

$$|\varphi(t)| = |t|^3 = |t|^{\frac{5}{\frac{5}{3}}},$$

such that the superposition operator  $y \mapsto y^3$  is continuous from  $L^5(\Omega)$  to  $L^{\frac{5}{3}}(\Omega)$ . (Exercises!) Further, since  $Y = H^1(\Omega) \hookrightarrow L^{\frac{5}{2}}(\Omega)$  densely, we also have the adjoint embedding  $(L^{\frac{5}{2}}(\Omega))^* = L^{\frac{5}{3}}(\Omega) \hookrightarrow Y^*$ . By [Corollary 2.8](#) and the following factorization diagram, this means that the mapping  $Y \ni y \mapsto$

$(y^3, \cdot)_{L^2(\Omega)} \in Y^*$  is sequentially continuous from the weak topology on  $Y$  to the norm topology on  $Y^*$ :

$$Y \xrightarrow[\text{compact}]{\text{id}} L^5(\Omega) \xrightarrow[\text{continuous}]{y \mapsto y^3} L^{\frac{5}{3}}(\Omega) \xrightarrow[\text{continuous}]{\hookrightarrow} Y^*$$

The second component of  $G$  is the identity  $\text{id}$  on  $U$  and as such obviously weakly sequentially continuous. Overall,  $G$  is weakly sequentially continuous from  $Y \times U$  to  $Y^* \times U$ .

- (b) The function  $J$  is convex and continuous on  $Y \times U$ , hence weakly lower semicontinuous by [Proposition 2.3 \(c\)](#). (Here note that  $Y = H^1(\Omega) \hookrightarrow L^2(\Omega)$  and that the norm on a Banach space is convex and by definition continuous.)
- (c)  $U_{\text{ad}} \subset U$  is bounded, closed and convex.
- (d) The feasible set  $\mathcal{F} = \{(y, u) \in Y \times U : E(y, u) = 0, u \in U_{\text{ad}}\}$  is nonempty and bounded:

Since  $U_{\text{ad}}$  is bounded, the set  $\mathcal{F}$  is also bounded by [\(2.1\)](#). Further,  $(y(a), a) \in \mathcal{F}$ .

Thus, the assumptions of [Theorem 2.4](#) are verified and there exists a globally optimal solution to the semilinear optimal boundary control problem.

From the foregoing example it should become clear that the functional-analytic setup for [\(P\)](#) plays a decisive role in whether the problem has a globally optimal solution or not. (See also the exercises.) This is a general paradigm in the theory of infinite-dimensional optimization problems and in particular in the context of PDE-constrained optimization where this paradigm connects in a natural way with the concept of weak or strong solutions to PDEs.

## The elephant in the room

As mentioned before, [\(P\)](#) is in general an infinite-dimensional optimization problem. In particular, the searched-for variable  $x$  is an object of infinite dimension. As such, we can generally not represent it exactly on a computer, which will at some point prove to be a problem when a practical problem instance of [\(P\)](#) need be solved numerically. One could thus reasonably pose the question why one has to suffer through the complicated reasoning for infinite-dimensional optimization problems at all.

In fact, what would probably happen for a numerical implementation is that one would employ finite-dimensional approximations  $X_n$  and  $Z_n$  of dimension  $n$  of the occurring Banach spaces and the associated discretizations  $G_n$  and  $f_n$  of the constraint function (in particular, the differential equations) and the objective function. Then one would

solve the resulting finite-dimensional optimization problems

$$\min_{x \in X_n} f_n(x) \quad \text{s.t.} \quad G_n(x) \in K_n, \quad (\mathbf{P}_n)$$

numerically with ones favorite algorithm from classical Nonlinear Optimization, repeatedly with increasing discretization accuracy, so  $n \rightarrow \infty$ . Since the underlying problems are then posed in finite-dimensional space, one could avoid the delicate mathematical issues arising for infinite-dimensional problems in order to deal with  $(\mathbf{P}_n)$ . This ansatz is nowadays summarized by “first-discretize-then-optimize”.

However, it turns out that in general it is not a good idea to ignore the infinite-dimensional original problem  $(\mathbf{P})$ . Indeed, already in order to argue that the solutions obtained for  $(\mathbf{P}_n)$  are really a good approximation for a “true” solution to  $(\mathbf{P})$ , one at least needs to know that such a solution exists, and in what function space norm sense we can hope for an approximation at all. (It can *very well* happen that the family  $(\mathbf{P}_n)$  admits optimal solutions but  $(\mathbf{P})$  will not!)

Another problem is that with increasing discretization accuracy, and thus increasing space dimension  $n$ , the numerical algorithms used to solve  $(\mathbf{P}_n)$  can become very inefficient and slow. It is usually preferable to determine an optimization algorithm in the infinite-dimensional setting and then pass to an appropriate discretization of this *in each step*. This would be “first-optimize-then-discretize”. For example, for some algorithms derived in such a way, one can show that they will perform in a certain sense independent of the discretization accuracy and are thus extremely effective. Such algorithms are usually based on a so-called first-order optimality system for  $(\mathbf{P})$ , the derivation of which is the subject of the next section which will occupy us until nearly the end of the course.

In this sense, although it may not seem so from the start, it is in fact of actual practical importance to deal with the full infinite-dimensional optimization problem  $(\mathbf{P})$  in order to derive efficient algorithms and appropriate approximation properties for numerical implementations. We also refer to the course “Numerik optimaler Steuerungen” which is planned for the coming summer term 2023.

November 2, 2022

3 →

### 3 Optimality conditions

We now turn to optimality conditions for the general problem  $(\mathbf{P})$ , both of first order necessary– and second order sufficient type. These will be conditions which characterize optimal solutions of  $(\mathbf{P})$ . Since  $(\mathbf{P})$  is in general a nonconvex problem, there will be possibly multiple local solutions which are not necessarily also global ones. (The techniques in the previous section yield existence of a *globally* optimal solution, though!) We can in general only expect characterizations of *local* solutions of  $(\mathbf{P})$ .



**Reminder:** A *local* solution  $\bar{x}$  to an optimization problem  $\min f(x)$  such that  $x \in M$  is characterized by the existence of an  $\varepsilon > 0$  such that  $\bar{x}$  is a global solution to the localized problem  $\min f(x)$  with the constraint  $x \in M \cap B(\bar{x}, \varepsilon)$ .

Before we begin with the actual definitions and statements leading to optimality conditions, we recall and define some concepts which will be needed in the following. Firstly, from the name *first* or *second order* optimality conditions, it is clear that derivatives will play a decisive role from now on.

**Definition 3.1** (Differentiability). Let  $X, Y$  be Banach spaces and let  $F: X \supseteq U \rightarrow Y$  be an operator between  $X$  and  $Y$  defined on an open subset  $U \neq \emptyset$  of  $X$ .

- (a) We say that  $F$  is *directionally differentiable* (**richtungsdifferenzierbar**) at  $x \in U$  if the limit

$$dF(x, h) = \lim_{t \searrow 0} \frac{F(x + th) - F(x)}{t} \in Y$$

exists for all  $h \in X$ . In this case,  $dF(x, h)$  is called *directional derivative* of  $F$  at  $x$  in the direction  $h$ .

- (b) Moreover,  $F$  is called *Gâteaux differentiable* (G-differentiable) at  $x \in U$  if  $F$  is directionally differentiable at  $x$  and the directional derivative as a function in  $h$ , so  $X \ni h \mapsto dF(x, h) \in Y$ , is bounded and linear, i.e., it is given by a linear operator  $A \in \mathcal{L}(X; Y)$  for which we set  $F'(x) := A$ .
- (c) Finally, we say that  $F$  is *Fréchet differentiable* (F-differentiable) at  $x \in U$  if  $F$  is Gâteaux differentiable at  $x$  and the following approximation condition holds:

$$\|F(x + h) - F(x) - F'(x)h\|_Y = o(\|h\|_X) \quad \text{for } h \rightarrow 0 \text{ in } X.$$

- (d) If  $F$  is directionally-/G-/F-differentiable at every  $x \in V$  for  $V \subseteq U$  open, then  $F$  is called *directionally-/G-/F-differentiable on  $V$* .
- (e) If  $F$  is G- or F-differentiable on a neighborhood of  $x \in U$  and the derivative mapping  $X \ni x \mapsto F'(x) \in \mathcal{L}(X; Y)$  is continuous, then we say that  $F$  is *continuously G-/F-differentiable*.

**Reminder:** The  $o$  notation in the definition of F-differentiability means that

$$\lim_{h \rightarrow 0} \frac{\|F(x+h) - F(x) - F'(x)h\|_Y}{\|h\|_X} = 0,$$

so, when  $h \rightarrow 0$  in  $X$ , then the linear approximation error  $\|F(x+h) - F(x) - F'(x)h\|_Y$  goes to zero *much faster* than  $\|h\|_X$ . In this context, it is also useful to observe that the approximation condition determines the derivative uniquely; that is, if some linear operator  $A \in \mathcal{L}(X; Y)$  satisfies

$$\|F(x+h) - F(x) - Ah\|_Y = o(\|h\|_X) \quad \text{for } h \rightarrow 0 \text{ in } X,$$

then it already follows that  $F'(x) = A$ .

**Example 3.2.** Classical examples and counterexamples for differentiable mappings in Banach spaces include (see the exercises):

- (a) Every continuous linear function  $F(x) = Ax$  defined by a bounded linear operator  $A \in \mathcal{L}(X; Y)$  is continuously F-differentiable and its derivative in every point  $x \in X$  is given by  $F'(x) = A$ ; in particular,  $F'(x)h = F(h)$ .
- (b) The quadratic form  $X \ni u \mapsto \frac{1}{2}a(u, u) \in \mathbb{R}$  induced by a (symmetric) continuous bilinear form  $a: X \times X \rightarrow \mathbb{R}$  is continuously F-differentiable and its derivative in  $u \in X$  is given by  $h \mapsto a(u, h)$ .
- (c) The superposition operator given by  $\sin: \mathbb{R} \rightarrow \mathbb{R}$  is continuously F-differentiable as a mapping from  $L^\infty(0, 1)$  into itself with the derivative  $L^\infty(0, 1) \ni h \mapsto \cos(f)h$  in  $f \in L^\infty(0, 1)$ . The operator is *not* F-differentiable as a mapping from  $L^p(0, 1)$  into itself for any  $1 \leq p < \infty$ .
- (d) The superposition operator given by the real function  $t \mapsto |t|^3$  is continuously F-differentiable as a mapping from  $L^6(\Omega)$  to  $L^2(\Omega)$  with the derivative  $L^6(\Omega) \ni h \mapsto 3|f|^2 h$  in  $y \in L^6(\Omega)$ . (Here  $\Omega \subseteq \mathbb{R}^n$  is a nonempty measurable set.)

By imitating the proof for the finite-dimensional case, many classical formulas as the following are easily established. (Work this out!)

**Proposition 3.3.** *Let  $X, Y$  and  $Z$  be Banach spaces.*

- (a) *If  $F, G: X \rightarrow Y$  are (continuously) F-differentiable, then so are  $\alpha F + \beta G$  for every  $\alpha, \beta \in \mathbb{R}$  and  $(\alpha F + \beta G)' = \alpha F' + \beta G'$ .*
- (b) *If  $F: X \rightarrow Y$  and  $G: Y \rightarrow Z$  are (continuously) F-differentiable in  $x \in X$  and  $F(x) \in Y$ , respectively, then  $G \circ F$  is (continuously) F-differentiable in  $x$  and*

its derivative is given by

$$(G \circ F)'(x) = G'(F(x)) \circ F'(x) \in \mathcal{L}(X; Z),$$

where the latter means that  $(G \circ F)'(x)h$  is given by the application of the operator  $G'(F(x)) \in \mathcal{L}(Y; Z)$  to  $F'(x)h \in Y$  for every  $h \in X$ .

**Example 3.4** (PDE Operator). Consider the mapping  $H^1(\Omega) \times L^2(\partial\Omega) \ni (y, u) \mapsto E(y, u) \in H^1(\Omega)^*$  corresponding to the weak formulation of the semilinear problem as posed in [Example 2.9](#), that is,

$$\langle E(y, u), v \rangle := (\nabla y, \nabla v)_{L^2(\Omega)^n} + (y^3, v)_{L^2(\Omega)} + (y - u, v)_{L^2(\partial\Omega)} \text{ for all } v \in H^1(\Omega).$$

As observed in [Example 2.9](#),  $(y, u) = x \mapsto F(x) := E(y, u) - (y^3, \cdot)_{L^2(\Omega)}$  is a linear and continuous mapping  $H^1(\Omega) \times L^2(\partial\Omega) \rightarrow H^1(\Omega)^*$ . According to [Example 3.2](#), it is thus continuously F-differentiable with the derivative  $F'(x)h = F(h)$ . Moreover, the superposition operator  $y \mapsto y^3$  is continuously F-differentiable from  $L^6(\Omega)$  to  $L^2(\Omega)$  as noted in [Example 3.2](#). Accordingly,  $y \mapsto (y^3, \cdot)_{L^2(\Omega)}$  is continuously F-differentiable  $H^1(\Omega) \rightarrow H^1(\Omega)^*$  since it factors into

$$H^1(\Omega) \xrightarrow[\text{cont. linear}]{\text{id}} L^6(\Omega) \xrightarrow[\text{continuous}]{y \mapsto y^3} L^2(\Omega) \xrightarrow[\text{cont. linear}]{z \mapsto (z, \cdot)_{L^2(\Omega)}} H^1(\Omega)^*.$$

Here we use a chain rule as in [Proposition 3.3](#). Altogether this yields that  $E$  is continuously F-differentiable with  $E'(y, u) \in \mathcal{L}(H^1(\Omega) \times L^2(\partial\Omega); H^1(\Omega)^*)$  given by

$$E'(y, u)(z, h): \left[ v \mapsto (\nabla z, \nabla v)_{L^2(\Omega)^n} + (z - h, v)_{L^2(\partial\Omega)} + 3(y^2 z, v)_{L^2(\Omega)} \right] \in H^1(\Omega)^*.$$

The partial derivatives of  $E$  are given by  $E'_y(y, u) \in \mathcal{L}(H^1(\Omega); H^1(\Omega)^*)$ ,

$$E'_y(y, u)z: \left[ v \mapsto (\nabla z, \nabla v)_{L^2(\Omega)^n} + (z, v)_{L^2(\partial\Omega)} + 3(y^2 z, v)_{L^2(\Omega)} \right] \in H^1(\Omega)^*$$

and  $E'_u(y, u) \in \mathcal{L}(L^2(\partial\Omega); H^1(\Omega)^*)$ ,

$$E'_u(y, u)h: \left[ v \mapsto -(h, v)_{L^2(\partial\Omega)} \right] \in H^1(\Omega)^*$$

In particular, the linearized operator  $E'_y(y, u) \in \mathcal{L}(H^1(\Omega); H^1(\Omega)^*)$  and the equation  $E'_y(y, u)z = w$  give rise to the weak formulation of the *linear* PDE

$$\begin{aligned} -\Delta z + 3y^2 z &= f && \text{on } \Omega, \\ \frac{\partial z}{\partial \nu} + z &= g && \text{on } \partial\Omega \end{aligned}$$

for the functional  $w \in H^1(\Omega)^*$  defined by

$$\langle w, v \rangle := \int_{\Omega} f(x)v(x) \, dx + \int_{\partial\Omega} g(x)v(x) \, d\omega(x).$$

A main goal of this chapter is to ultimately derive conditions under which, given a local solution  $\bar{x}$  of (P), a so-called *Lagrange multiplier* (associated to  $\bar{x}$ )  $\bar{\lambda} \in Z^*$  exists. The generic tools to prove *existence* of functionals in a dual space such as  $Z^*$  are *separation theorems* which we will use quite frequently. (Such arguments are also used in the finite-dimensional case for Nonlinear Optimization, or even already in Linear Programming, usually in so-called *Farkas lemmas*.) We thus recall the geometric, or separation, versions of the fundamental *Hahn-Banach theorem* [Br10, Ch. 1.2].

**Reminder:** For two subset  $A, B \subset X$  of a Banach space  $X$ , we say that the hyperplane

$$H = [f = \alpha] := \left\{ x \in X : \langle f, x \rangle = \alpha \right\}$$

induced by  $0 \neq f \in X^*$  and  $\alpha \in \mathbb{R}$  *separates* (trennt)  $A$  and  $B$  if  $\langle f, x \rangle \leq \alpha$  for all  $x \in A$  and  $\langle f, x \rangle \geq \alpha$  for all  $x \in B$ . We say the hyperplane  $H$  *strictly separates*  $A$  and  $B$  if there exists  $\varepsilon > 0$  such that  $\langle f, x \rangle \leq \alpha - \varepsilon$  for all  $x \in A$  and  $\langle f, x \rangle \geq \alpha + \varepsilon$  for all  $x \in B$ .

**Proposition 3.5** (Hahn-Banach, geometric form). *Let  $X$  be a Banach space and let  $A, B \subset X$  be two nonempty convex subsets such that  $A \cap B = \emptyset$ .*

- (a) *Assume that  $A$  or  $B$  is open. Then there exists a hyperplane separating  $A$  and  $B$ .*
- (b) *Assume that  $A$  is closed and that  $B$  is compact. Then there exists a hyperplane strictly separating  $A$  and  $B$ .*

A general class of first order necessary optimality conditions is obtained by observing that the Fréchet derivative at a local solution  $\bar{x}$  is nonnegative along all directions that are tangential to the feasible set or point into the feasible set. These directions are characterized by the *tangent cone* (Tangentialkegel) of the feasible set at the solution.

**Definition 3.6** (Tangent cone). The *Bouligand tangent cone* (or *contingent cone*) of a set  $M \subseteq X$ , where  $X$  is a Banach space, at  $x \in M$  is defined by

$$T(M, x) = \left\{ d \in X : \exists (x_k) \subseteq M, x_k \rightarrow x, (\eta_k) > 0 : \eta_k(x_k - x) \rightarrow d \right\}.$$

The set  $T(M, x)$  is a closed cone for every  $x \in M$ .

See the exercises for the proof that  $T(M, x)$  is closed. In the case when  $M$  is *convex*, a—possibly—more intuitive description of  $T(M, x)$  can be derived using the conical hull. This description relies on the idea that we are interested in all directions which point into the feasible set. We need two more notions, which are however of independent interest and will be used frequently in the following.

**Reminder:** Recall that the *Minkowski sum* of two sets  $A, B \subset X$  is given by

$$A + B := \{a + b : a \in A, b \in B\}$$

with the convention  $a + B := \{a\} + B$ , allowing to write the *Minkowski difference* as

$$A - B := \{c \in X : c + B \subseteq A\}.$$

**Definition 3.7** (Conical hull, cone of radial directions). Let  $X$  be a vector space and let  $\emptyset \neq A \subseteq X$  be convex. Then the *conical hull* (*konische Hülle*) of  $A$  is given by

$$\text{cone}(A) := \{\lambda y : y \in A, \lambda > 0\}.$$

It is the smallest cone which includes the set  $A$ . More generally, the *cone of radial directions* of  $A$  in a point  $x \in X$  is given by

$$\text{cone}(A, x) := \text{cone}(A - x) = \{z \in X : z = \lambda(y - x), y \in A, \lambda > 0\}.$$

It is the smallest cone  $C$  such that  $A \subseteq x + C$ . See also [Figure 1](#).

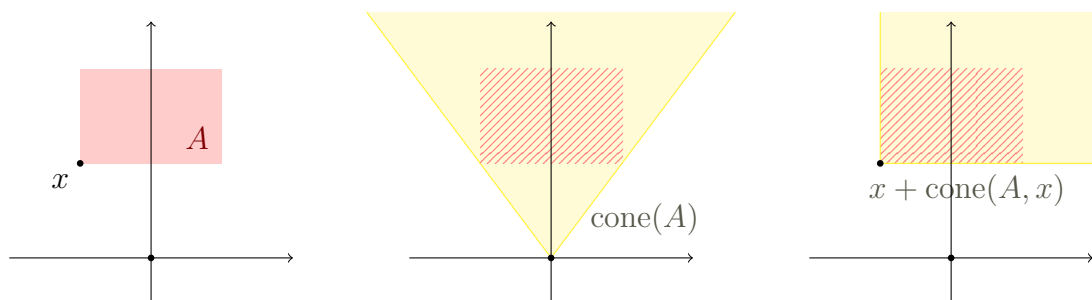


Figure 1: Conical hull of  $A$  and cone of radial directions of  $A$  at  $x$

**Lemma 3.8.** *Let  $M \subseteq X$  be convex and  $x \in M$ , then*

$$T(M, x) = \overline{\text{cone}(M, x)}.$$

*Moreover,  $\text{cone}(M, x)$  and thus  $T(M, x)$  are convex.*

*Proof.* Let  $d \in T(M, x)$ . Then there exist  $(x_k) \subset M$  with  $x_k \rightarrow x$  and  $(\eta_k) > 0$  with  $d_k := \eta_k(x_k - x) \rightarrow d$ . By definition,  $(d_k) \subset \text{cone}(M, x)$ . (Set  $y = x_k$  and  $\lambda = \eta_k$  for each  $k$ .) Hence,  $\lim_{k \rightarrow \infty} d_k = d \in \overline{\text{cone}(M, x)}$ .

Conversely, to prove  $\overline{\text{cone}(M, x)} \subset T(M, x)$  it is sufficient to prove  $\text{cone}(M, x) \subset T(M, x)$  since  $T(M, x)$  is closed. To this end, consider any  $d \in \text{cone}(M, x)$ . Then there exist  $\lambda > 0$  and  $y \in M$  with  $d = \lambda(y - x)$ . Now set  $\eta_k := k\lambda$  and  $x_k := x + (y - x)/k$  for  $k \geq 1$ . Then  $x_k \in M$  by convexity of  $M$  and  $x_k \rightarrow x$ . Further,  $\eta_k(x_k - x) = \lambda(y - x) = d$ . This shows  $d \in T(M, x)$ .

We lastly show that if  $M$  is convex, then so is  $\text{cone}(M, x)$ . Since the closure of convex sets is convex (why?), this implies that  $T(M, x)$  is convex, too. Let  $d_1, d_2 \in \text{cone}(M, x)$ , i.e., for  $i = 1, 2$ , there are  $\lambda_i > 0$  and  $y_i \in M$  such that  $d_i = \lambda_i(y_i - x)$ . Then, for  $\alpha \in (0, 1)$ , we have

$$(1 - \alpha)d_1 + \alpha d_2 = ((1 - \alpha)\lambda_1 + \alpha\lambda_2) \left( \frac{(1 - \alpha)\lambda_1}{(1 - \alpha)\lambda_1 + \alpha\lambda_2} y_1 + \frac{\alpha\lambda_2}{(1 - \alpha)\lambda_1 + \alpha\lambda_2} y_2 - x \right).$$

Since  $M$  is assumed to be convex,  $(1 - \alpha)d_1 + \alpha d_2 \in \text{cone}(M, x)$  follows.  $\square$

As apparent from the proof, the inclusion  $T(M, x) \subseteq \overline{\text{cone}(M, x)}$  is always true, independent of convexity of  $M$ . (Think of nonconvex set  $M$  where  $\text{cone}(M, x) \not\subseteq T(M, x)$ !)

We now can state a first order optimality condition.

**Theorem 3.9.** *Let  $f: U \rightarrow \mathbb{R}$  be defined on an open neighborhood  $U$  of the set  $M \subseteq X$ , where  $X$  is a Banach space. Let  $\bar{x} \in M$  be a local minimizer of  $f$  on  $M$ , i.e., a local solution of  $\min_{x \in M} f(x)$  and assume that  $f$  is  $F$ -differentiable at  $\bar{x}$ . Then there holds*

$$\langle f'(\bar{x}), d \rangle_{X^*, X} \geq 0 \quad \text{for all } d \in T(M, \bar{x}). \quad (3.1)$$

*Proof.* For all  $d \in T(M, \bar{x})$  there exist sequences  $(x_k) \subseteq M$  and  $(\eta_k) > 0$  such that  $x_k \rightarrow \bar{x}$  and  $\eta_k(x_k - \bar{x}) \rightarrow d$ .

We then have  $\eta_k o(\|x_k - \bar{x}\|_X) \rightarrow 0$ . Now, for sufficiently large  $k$ , there holds  $f(x_k) \geq f(\bar{x})$  since  $\bar{x}$  was a local minimum, and thus

$$\begin{aligned} \langle f'(\bar{x}), d \rangle &= \lim_{k \rightarrow \infty} \eta_k \langle f'(\bar{x}), x_k - \bar{x} \rangle = \lim_{k \rightarrow \infty} \left[ \eta_k (f(x_k) - f(\bar{x})) + \eta_k o(\|x_k - \bar{x}\|_X) \right] \\ &\geq \lim_{k \rightarrow \infty} \eta_k o(\|x_k - \bar{x}\|_X) = 0. \end{aligned} \quad \square$$

**Remark 3.10.** Note that if  $x$  is an *interior point* of  $M$ —in particular, if  $M = X$ —, then  $T(M, \bar{x}) = X$ , and the foregoing optimality condition becomes  $f'(\bar{x}) = 0$  (work this out!). This is the well known classical stationarity condition which one already learns in high school.

We are interested in applying this result to the problem (P), i.e., with  $M = \mathcal{F}$ . However, the cone  $T(\mathcal{F}, \bar{x})$  is difficult and often impossible to compute in practice. Hence, we approximate it by linearization and give conditions (so-called *constraint qualifications*) under which the linearizing cone and the contingent cone coincide. We recall that the set  $K$  is convex as in the basic assumptions in [Assumption 1.1](#).

**Definition 3.11** (Linearizing cone). Let  $G$  be F-differentiable at  $x \in \mathcal{F} = G^{-1}[K]$ . The *linearizing cone* at  $x$  is given by

$$\begin{aligned} T_\ell(G, K, x) &= \left\{ d \in X : G'(x)d \in T(K, G(x)) \right\} \\ &= \left\{ d \in X : G'(x)d \in \overline{\text{cone}(K, G(x))} \right\}. \end{aligned}$$

**Remark 3.12.** Note that due to convexity of  $K$ , the linearizing cone  $T_\ell(G, K, x)$  is always convex. (See [Lemma 3.8](#).) It is moreover closed.

**Remark 3.13.** For the classical NLP

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

with  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$  continuously differentiable, given in our setting for (P) by  $X = \mathbb{R}^n$ ,  $Z = \mathbb{R}^m \times \mathbb{R}^p$ ,  $G(x) = \begin{pmatrix} g(x) \\ h(x) \end{pmatrix}$  and  $K = (-\infty, 0]^m \times \{0\}^p$ , we indeed have (see the exercises)

$$T_\ell(G, K, x) = \left\{ d \in \mathbb{R}^n : \nabla h(x)^T d = 0, \nabla g_i(x)^T d \leq 0 \text{ for } i \in \mathcal{A}(x) \right\},$$

where  $\mathcal{A}(x) = \{i: g_i(x) = 0\}$  is the set of active inequality constraints. This is the standard linearizing cone from Nonlinear Optimization.

We would like to infer from (3.1) that the optimality condition holds not only for directions  $d$  from the cone  $T(\mathcal{F}, \bar{x})$ , using  $M = \mathcal{F}$ , but also for directions from the cone  $T_\ell(G, K, \bar{x})$ . This, however, is in general only true if  $T_\ell(G, K, \bar{x}) \subseteq T(\mathcal{F}, \bar{x})$ .

**Definition 3.14** (Abadie Constraint Qualification). The *Abadie Constraint Qualification* (ACQ) holds true at a point  $\bar{x} \in \mathcal{F}$  if the condition

$$T_\ell(G, K, \bar{x}) \subseteq T(\mathcal{F}, \bar{x})$$

is satisfied. We call any condition which implies (ACQ) a *constraint qualification* (CQ).

Note that there are sets  $\mathcal{F}$  where (ACQ) is *never* satisfied. For example, this is always the case when  $T(\mathcal{F}, x)$  is not convex, since we have noted in Remark 3.12 that  $T_\ell(G, K, x)$  is always convex, and the reverse inclusion is always true:

**Lemma 3.15.** *Let  $G$  be  $F$ -differentiable at  $x \in \mathcal{F} = G^{-1}[K]$ . Then  $T(\mathcal{F}, x) \subseteq T_\ell(G, K, x)$ .*

*Proof.* For  $d \in T(\mathcal{F}, x)$  there exist sequences  $(x_k) \subseteq \mathcal{F}$  and  $(\eta_k) > 0$  such that  $x_k \rightarrow x$  and  $\eta_k(x_k - x) \rightarrow d$ . Without loss of generality, we can assume  $\eta_k \rightarrow \infty$  (why?). We need to show that  $G'(x)d \in T(K, G(x))$ . As in the proof of Theorem 3.9, we find  $\eta_k o(\|x_k - \bar{x}\|_X) \rightarrow 0$  and thus

$$\eta_k [G(x_k) - G(x)] = G'(x) [\eta_k(x_k - x)] + \eta_k o(\|x_k - x\|_X) \rightarrow G'(x)d.$$

Due to  $\eta_k \rightarrow \infty$ , or, alternatively, continuity of  $G$ , we have  $G(x_k) \rightarrow G(x)$  and, of course,  $G(x_k) \in K$  for every  $k$ . We conclude  $G'(x)d \in T(K, G(x))$ .  $\square$

### 3.1 Robinson's constraint qualification

The first CQ we consider is an algebraic-topological condition due to Robinson [Ro76] and will be of great relevance. It will be used to prove the *Karush-Kuhn-Tucker conditions* for (P), which are the most common form of first order optimality conditions.

**Definition 3.16** (Robinson's constraint qualification, regularity). We say that *Robinson's constraint qualification* (RCQ) is satisfied for the problem (P) at  $\bar{x} \in \mathcal{F}$  if there holds

$$0 \in \text{int}(G(\bar{x}) + G'(\bar{x})X - K). \quad (3.2)$$

In this case, we also say that  $\bar{x} \in \mathcal{F}$  is *regular*.

In this section we will show that (3.2) implies the ACQ. (So it *is* a CQ.) First, we establish the connection to a well known CQ in the finite-dimensional case.



**Example 3.17.** Consider the case of an NLP

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

as in [Remark 3.13](#). We will show that in this case (3.2) is equivalent to the *Mangasarian Fromovitz constraint qualification* (MFCQ):

$$\text{rank } \nabla h(\bar{x}) = p, \quad \exists d \in \mathbb{R}^n: \quad \nabla h(\bar{x})^T d = 0, \quad \nabla g_i(\bar{x})^T d < 0 \quad \text{for } i \in \mathcal{A}(\bar{x}). \quad (3.3)$$

Let RCQ (3.2) be true for  $\bar{x}$  satisfying  $g(\bar{x}) \leq 0$  and  $h(\bar{x}) = 0$ , i.e.,

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

The lower block requires that  $\nabla h(\bar{x})^T$  is surjective, which implies  $\text{rank } \nabla h(\bar{x}) = p$ . Now, let  $\delta > 0$  and set  $w \in \mathbb{R}^m$  by  $w_i = -\delta$  if  $g_i(\bar{x}) = 0$  and  $w_i = 0$  if  $g_i(\bar{x}) < 0$ . Then

$$\begin{pmatrix} w \\ 0 \end{pmatrix} \in \left\{ \begin{pmatrix} g(\bar{x}) + \nabla g(\bar{x})^T s - v \\ \nabla h(\bar{x})^T s \end{pmatrix} : s \in \mathbb{R}^n, v \in (-\infty, 0]^m \right\}$$

if we choose  $\delta$  sufficiently small. This means that there exist  $s \in \mathbb{R}^n$  and  $v \in (-\infty, 0]^m$  with

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

Hence,  $\nabla h(\bar{x})^T s = 0$  and, for all  $i$  with  $g_i(\bar{x}) = 0$ :

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

Thus, the MFCQ (3.3) is satisfied for  $d := s$ .

Conversely, let the MFCQ (3.3) hold true and let  $d \in \mathbb{R}^n$  be the corresponding vector. We show that there exists  $\varepsilon > 0$  such that

$$B_{\varepsilon, \mathbb{R}^{m+p}} \subseteq \left\{ \begin{pmatrix} g(\bar{x}) + \nabla g(\bar{x})^T s - v \\ \nabla h(\bar{x})^T s \end{pmatrix} : s \in \mathbb{R}^n, v \in (-\infty, 0]^m \right\}.$$

Firstly, from  $g_i(\bar{x}) < 0$  for  $i \notin \mathcal{A}(\bar{x})$ , we can find numbers  $\delta_1, t > 0$  such that

$$g_i(\bar{x}) + \nabla g_i(\bar{x})^T t d < -2\delta_1 \quad \text{for } i \notin \mathcal{A}(\bar{x}).$$

Secondly, for  $i \in \mathcal{A}(\bar{x})$  we know that  $\nabla g_i(\bar{x})^T d < 0$ , thus we can find another number  $\delta_2 > 0$  such that  $\nabla g_i(\bar{x})^T t d < -2\delta_2$  for all  $i \in \mathcal{A}(\bar{x})$ . But then we have

$$g_i(\bar{x}) + \nabla g_i(\bar{x})^T t d < -2\delta \quad \text{for } i = 1, \dots, m,$$

where  $\delta := \min(\delta_1, \delta_2)$ , and there exists  $\rho > 0$  such that

$$g_i(\bar{x}) + \nabla g_i(\bar{x})^T(td + s_0) < -\delta \quad \text{for } i = 1, \dots, m,$$

for all  $s_0 \in B_{\rho, \mathbb{R}^n}(0)$ . From the rank assumption on  $\nabla h(\bar{x})$  in (3.3), we finally obtain a number  $\varepsilon_1$  such that  $B_{\varepsilon_1, \mathbb{R}^p}(0) \subseteq \nabla h(\bar{x})^T B_{\rho, \mathbb{R}^n}(0)$  (why?) and set  $\varepsilon := \min(\delta, \varepsilon_1)$ .

Now consider an arbitrary vector  $w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \in \mathbb{R}^{m+p}$  with  $w \in B_{\varepsilon, \mathbb{R}^{m+p}}(0)$ . Then we also have  $\|w_i\|_2 < \varepsilon$  for  $i = 1, 2$ , and by construction of  $\rho$  and  $\varepsilon$  there exists a vector  $s_0 \in B_{\rho, \mathbb{R}^n}(0) \subset \mathbb{R}^n$  with  $w_2 = \nabla h(\bar{x})^T s_0 = \nabla h(\bar{x})^T(td + s_0)$ . Choosing  $s := td + s_0$  and

$$v := -w_1 + g(\bar{x}) + \nabla g(\bar{x})^T s,$$

we find  $v_i < \varepsilon - \delta \leq 0$  and thus  $v \in (-\infty, 0]^m$ . But this means that

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

and, since  $w \in B_{\varepsilon, \mathbb{R}^{m+p}}(0)$  was arbitrary,

$$B_{\varepsilon, \mathbb{R}^{m+p}}(0) \subset \left\{ \begin{pmatrix} g(\bar{x}) + \nabla g(\bar{x})^T s - v \\ \nabla h(\bar{x})^T s \end{pmatrix} : s \in \mathbb{R}^n, v \in (-\infty, 0]^m \right\},$$

which is exactly Robinson's CQ (3.2) for  $\bar{x}$ . □

A particularly interesting case for RCQ is the one where  $\text{int } K \neq \emptyset$ , which allows to obtain a slightly more direct formulation.

**Lemma 3.18.** *If  $\text{int } K \neq \emptyset$ , then Robinson's CQ for  $\bar{x} \in \mathcal{F} = G^{-1}[K]$  is equivalent to the existence of  $h \in X$  with*

$$G(\bar{x}) + G'(\bar{x})h \in \text{int } K, \tag{3.4}$$

*which is also called the Linearized Slater constraint qualification (LSCQ).*

*Proof.* From (3.4) it follows that for  $\delta > 0$  sufficiently small there holds  $G(\bar{x}) + G'(\bar{x})h + B_{Z, \delta}(0) \subset K$ . Hence,

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

and RCQ is satisfied. We next show that if the LSCQ is not satisfied, then the RCQ also fails to hold. So, let (3.4) be violated. This means that the convex sets  $G(\bar{x}) + G'(\bar{x})X$  and  $\text{int } K$  have an empty intersection. (Why is  $\text{int } K$  convex?) By the Hahn-Banach theorem as in Proposition 3.5, the sets can thus be separated by a hyperplane  $[z' = \alpha]$ , i.e., there exist  $z' \in Z^* \setminus \{0\}$  and  $\alpha \in \mathbb{R}$  with

$$\langle z', G(\bar{x}) + G'(\bar{x})h \rangle \geq \alpha \geq \langle z', z \rangle \quad \text{for all } h \in X, z \in K.$$

Now choose a vector  $v \in Z$  with  $\langle z', v \rangle_{Z^*, Z} < 0$ . Then for all  $t > 0$  there holds

$$\langle z', G(\bar{x}) + G'(\bar{x})h - z \rangle \geq 0 > \langle z', tv \rangle \quad \text{for all } h \in X, z \in K,$$

which shows that  $tv \notin G(\bar{x}) + G'(\bar{x})X - K$  for all  $t > 0$ . But this means that Robinson's CQ (3.2) cannot hold.  $\square$

There are also more sufficient conditions and equivalencies to RCQ, in particular for the often-occurring case of multiple “blocks” of constraints with different structure. We collect some useful cases in the following proposition. For the proofs we refer to the exercises.

**Proposition 3.19.** *Let  $\bar{x} \in \mathcal{F} = G^{-1}[K]$  be given.*

- (a) *If  $G'(\bar{x}) \in \mathcal{L}(X; Z)$  is surjective, then RCQ is satisfied for  $\bar{x}$ .*
- (b) *Ljusternik's theorem: Let  $K = \{0_Z\}$ . Then RCQ for  $\bar{x}$  is satisfied if and only if  $G'(\bar{x}): X \rightarrow Z$  is surjective.*

Now let  $G$  be of the form

$$G = \begin{pmatrix} G_1 \\ G_2 \end{pmatrix} : X \rightarrow Z_1 \times Z_2 = Z$$

and let accordingly  $K = K_1 \times K_2$  where  $K_i \subseteq Z_i$  for  $i = 1, 2$ .

- (c) *If  $G'_1(\bar{x}) \in \mathcal{L}(X; Z_1)$  is surjective, then RCQ for  $\bar{x}$  is equivalent to*

$$0 \in \text{int} \left( G_2(\bar{x}) + G'_2(\bar{x})(G'_1(\bar{x})^{-1}[K_1 - G_1(\bar{x})]) - K_2 \right).$$

- (d) *If  $G'_1(\bar{x}) \in \mathcal{L}(X; Z_1)$  is surjective and  $\text{int} K_2 \neq \emptyset$  in  $Z_2$ , then RCQ for  $\bar{x}$  is equivalent to the existence of  $h \in X$  such that*

$$\begin{aligned} G_1(\bar{x}) + G'_1(\bar{x})h &\in K_1, \\ G_2(\bar{x}) + G'_2(\bar{x})h &\in \text{int} K_2. \end{aligned} \tag{3.5}$$

*In particular, if  $K_1 = \{0_{Z_1}\}$ , then the first of the two conditions in (3.5) collapses to  $h \in \ker G'_1(\bar{x})$ .*

We want to show that Robinson's CQ (3.2) implies the ACQ. To do so, we will need an estimate that allows us to bound the distance of a point  $x$  from  $G^{-1}[K]$  by means of the distance of  $G(x)$  from  $K$ .

**Reminder:** The *distance* of a point  $b$  to a set  $A$  in a normed vector space is given by

$$\text{dist}(b, A) := \inf_{a \in A} \|a - b\|.$$

A special case of a celebrated result by Robinson provides such an estimate, cf. [BS00, Thm. 2.87], see also [Ro76, Cor. 1] and [Ro76b]:

**Theorem 3.20.** *Assume that Robinson's CQ is satisfied at  $\bar{x} \in \mathcal{F} = G^{-1}[K]$  and let  $G: X \rightarrow Z$  be continuously  $F$ -differentiable near  $\bar{x}$ . Then there exist constants  $c > 0$  and  $\delta > 0$  such that*

$$\text{dist}(x, G^{-1}[K - z]) \leq c \text{dist}(G(x) + z, K) \quad (3.6)$$

for all  $x \in B_{\delta, X}(\bar{x})$  and  $z \in B_{\delta, Z}(0)$ , where

$$G^{-1}[K - z] = \left\{ x \in X : G(x) + z \in K \right\}.$$

**Remark 3.21.**

- (a) The condition (3.6) is called *metric regularity of  $G$  at  $\bar{x}$  with respect to  $K$* . This is also the origin of calling  $\bar{x}$  *regular* if it satisfies RCQ.
- (b) For the special choice  $z = 0$ , we obtain

$$\text{dist}(x, \mathcal{F}) \leq c \text{dist}(G(x), K)$$

for all  $x \in B_{\delta, X}(\bar{x})$ .

- (c) In the case  $K = \{0\}$ , the metric regularity condition (3.6) is equivalent to

$$\text{dist}(x, \{y \in X : G(y) = z\}) \leq c \|G(x) - z\|_Z$$

for all  $x \in B_{\delta, X}(\bar{x})$  and  $z \in B_{\delta, Z}(0)$ .

Using Theorem 3.20, we finally prove that Robinson's constraint qualification implies the ACQ.

**Theorem 3.22.** *Let  $G$  be continuously  $F$ -differentiable near  $\bar{x} \in \mathcal{F}$  and assume that Robinson's constraint qualification (3.2) holds at  $\bar{x}$ . Then the ACQ is satisfied at  $\bar{x}$ , i.e.,  $T(\mathcal{F}, \bar{x}) = T_{\ell}(G, K, \bar{x})$ .*

*Proof.* We only need to show  $T_{\ell}(G, K, \bar{x}) \subseteq T(\mathcal{F}; \bar{x})$ , so consider an arbitrary direction  $h \in T_{\ell}(G, K, \bar{x})$ . Then, by definition,  $G'(\bar{x})h = v \in T(K, G(\bar{x}))$  and there exist sequences  $(z_k) \subseteq K$ ,  $(\eta_k) > 0$  such that  $z_k \rightarrow G(\bar{x})$  and  $v_k := \eta_k(z_k - G(\bar{x})) \rightarrow v$  as  $k \rightarrow \infty$ . We can always choose  $(z_k)$  and  $(\eta_k)$  such that  $\eta_k \rightarrow \infty$  (why?) and we suppose that this is the case from now on.

By Taylor expansion, or the definition of F-differentiability, we obtain

$$\begin{aligned} G(\bar{x} + \eta_k^{-1}h) &= G(\bar{x}) + G'(\bar{x})\eta_k^{-1}h + r_k(h) = G(\bar{x}) + \eta_k^{-1}v + r_k(h) \\ &= z_k + \eta_k^{-1}(v - v_k) + r_k(h), \end{aligned}$$

where  $r_k(h) \in Z$  with  $\|r_k(h)\|_Z = o(\eta_k^{-1}\|h\|_X)$  as  $k \rightarrow \infty$ . Hence, using [Theorem 3.20](#) with  $z = 0$  (cf. [Remark 3.21](#)) and  $\eta_k^{-1} \rightarrow 0$ , there exists  $c > 0$  and  $\ell > 0$  such that for all  $k \geq \ell$  we have

$$\begin{aligned} \text{dist}(\bar{x} + \eta_k^{-1}h, \mathcal{F}) &\leq c \text{dist}(G(\bar{x} + \eta_k^{-1}h), K) \\ &\leq c \|G(\bar{x} + \eta_k^{-1}h) - z_k\|_Z = c \|\eta_k^{-1}(v - v_k) + r_k(h)\|_Z. \end{aligned}$$

Now, for each  $k \geq \ell$ , there exists an infimal sequence  $(x_m^k)_m \subseteq \mathcal{F}$  for  $\text{dist}(\bar{x} + \eta_k^{-1}h, \mathcal{F})$ . Thus there is  $m_0(k)$  such that

$$\|\bar{x} + \eta_k^{-1}h - x_m^k\|_X \leq \text{dist}(\bar{x} + \eta_k^{-1}h, \mathcal{F}) + \frac{1}{k\eta_k} \quad \text{for } m \geq m_0(k).$$

Accordingly, the sequence  $(x_k) := (x_{m_0(k)}^k)_k \subseteq \mathcal{F}$  satisfies

$$\|\bar{x} + \eta_k^{-1}h - x_k\|_X \leq c \|\eta_k^{-1}(v - v_k) + r_k(h)\|_Z + \frac{1}{k\eta_k} \quad \text{for } k \geq \ell.$$

But this implies

$$\|\eta_k(x_k - \bar{x}) - h\|_X \leq c \|v - v_k\|_Z + c\eta_k o(\eta_k^{-1}\|h\|_X) + \frac{1}{k} \longrightarrow 0 \quad \text{as } k \rightarrow \infty,$$

so  $\eta_k(x_k - \bar{x}) \rightarrow h$  and also  $x_k \rightarrow \bar{x}$  (recall that  $\eta_k \rightarrow \infty$ ). This proves  $h \in T(\mathcal{F}, \bar{x})$ .  $\square$

## 3.2 The Karush-Kuhn-Tucker conditions

We will take advantage of the easier structure of  $T_\ell(G, K, \bar{x})$  to derive optimality conditions of Karush-Kuhn-Tucker type under a constraint qualification. This will be achieved by applying a separation theorem, cf. [Proposition 3.5](#). To prove the required interior point condition, Robinson's CQ ([3.2](#)) is used once again in a form that will be derived now.

We give the first result in an abstract formulation to make the idea more transparent. It builds upon the following extension of the open mapping theorem for multi-valued functions, again due to Robinson ([\[Ro72\]](#), [\[Ro76b, Thm. 1\]](#)).

**Reminder:** The fundamental *open mapping theorem* says that if  $A \in \mathcal{L}(X; Y)$  is a continuous linear *surjective* operator between Banach spaces  $X$  and  $Y$ , then it is an open mapping, i.e.,  $0 \in \text{int}(AB_{r,X}(0))$  for all  $r > 0$ . The Banach space property for  $X$  and  $Y$  is crucial here.

Note that the graph of a continuous linear operator  $A \in \mathcal{L}(X; Y)$ , so the set  $\{(x, y) \in X \times Y : Ax = y\}$ , is closed and convex, and if  $A$  is surjective, then  $0 \in \text{int}(AB_{r,X}(0))$  for some (and thus all)  $r > 0$ . Moreover,  $A$  is surjective if and only if  $0 \in \text{int}(AX)$ . In this sense, the following result is a true generalization of the classical open mapping theorem.

**Proposition 3.23** (Generalized open mapping theorem). *Let  $\Psi: X \rightrightarrows Z$  be a set-valued function (i.e.,  $\Psi(x) \subseteq Z$  for all  $x \in X$ ) between Banach spaces  $X$  and  $Z$  whose graph*

$$\text{graph } \Psi = \left\{ (x, z) \in X \times Z : z \in \Psi(x) \right\}.$$

*is a closed convex set, and let  $z \in \text{int } \Psi(X)$ . Then  $z \in \text{int } \Psi(B_{r,X}(x))$  for all  $r > 0$  and all  $x \in \Psi^{-1}[\{z\}]$ .*

**Lemma 3.24.** *Let  $A \in \mathcal{L}(X; Z)$  with Banach spaces  $X$  and  $Z$  and let  $C \subseteq Z$  be a closed convex set with  $0 \in C$ . Then the following assertions are equivalent:*

- (i)  $AX + \text{cone}(C) = Z$ , and
- (ii)  $0 \in \text{int}\left(\overline{AB_X(0)} + (C \cap \overline{B_Z(0)})\right)$ .

*Proof.* We start with (ii)  $\implies$  (i). Let  $z \in Z$ . Choosing  $\varepsilon > 0$  sufficiently small, (ii) means that

$$\varepsilon z \in \overline{AB_X(0)} + (C \cap \overline{B_Z(0)}) \implies z \in \overline{AB_{\varepsilon^{-1},X}(0)} + \varepsilon^{-1}(C \cap \overline{B_Z(0)})$$

and thus  $z \in AX + \text{cone}(C)$ . Since  $z$  was arbitrary, this implies (i).

Now we turn to (i)  $\implies$  (ii). We employ the generalized open mapping theorem, setting

$$\Psi: X \times \mathbb{R} \rightrightarrows Z, \quad \Psi(x, t) := \begin{cases} Ax + t(C \cap \overline{B_Z(0)}) & \text{if } t \geq 0, \\ \emptyset & \text{otherwise.} \end{cases}$$

Then by (i),  $\Psi(X, \mathbb{R}) = AX + \text{cone}(C) = Z$  and thus clearly  $\Psi(0, 0) = 0 \in \text{int } \Psi(X, \mathbb{R})$ . Moreover, the graph of  $\Psi$  is closed and convex and, due to [Proposition 3.23](#), we obtain

$$\begin{aligned} 0 \in \text{int } \Psi(B_X(0), B(0)) &= \text{int}\left(\overline{AB_X(0)} + [0, 1](C \cap \overline{B_Z(0)})\right) \\ &\subset \text{int}\left(\overline{AB_X(0)} + (C \cap \overline{B_Z(0)})\right), \end{aligned}$$

where we have used that  $0 \in C$ . □

With [Lemma 3.24](#), we are now able to give an alternative form of the RCQ: the *Zowe-Kurcyusz constraint qualification* (ZKQC) [[ZK79](#)] in a feasible point  $\bar{x} \in \mathcal{F} = G^{-1}[K]$  given by

$$Z = G'(\bar{x})X - \text{cone}(K, G(\bar{x})). \quad (3.7)$$

It is indeed equivalent to (RCQ), and even to the restricted version

$$0 \in \text{int}\left(G'(\bar{x})\overline{B_X(0)} - ((K - G(\bar{x})) \cap \overline{B_Z(0)})\right), \quad (3.8)$$

as the following lemma shows:

**Lemma 3.25.** *Robinson's constraint qualification (3.2) in  $\bar{x} \in \mathcal{F}$  is equivalent to both forms of the Zowe-Kurcyusz constraint qualification, that is*

$$(3.2) \iff (3.7) \iff (3.8).$$

*Proof.* The equivalence of (3.7) and (3.8) is exactly the statement of [Lemma 3.24](#) for the choices  $A = G'(\bar{x})$  and  $C = G(\bar{x}) - K$ , the latter being a closed convex set satisfying  $0 \in C$  thanks to  $G(\bar{x}) \in K$ .

Robinson's CQ (3.2) follows from (3.8) immediately due to the existence of  $\varepsilon > 0$  such that

$$B_{\varepsilon, Z}(0) \subset G'(\bar{x})\overline{B_X(0)} - ((K - G(\bar{x})) \cap \overline{B_Z(0)}) \subset G(\bar{x}) + G'(\bar{x})X - K.$$

Lastly, from Robinson's CQ (3.2) we infer the (ZKQC) (3.7) by considering  $z \in Z$  and observing that for  $\varepsilon > 0$  sufficiently small we find  $\varepsilon z \in G(\bar{x}) + G'(\bar{x})X - K$  and thus

$$z \in G'(\bar{x})\varepsilon^{-1}X - \varepsilon^{-1}(K - G(\bar{x})) \subset G'(\bar{x})X - \text{cone}(K, G(\bar{x})).$$

Since  $z \in Z$  was arbitrary, this implies (3.7). □

For stating the main result of this section, we need the notion of the polar cone.

**Definition 3.26** (Polar cone). Let  $\emptyset \neq C \subseteq Z$ . Then the set

$$C^\circ := \left\{ z' \in Z^* : \langle z', z \rangle_{Z^*, Z} \leq 0 \text{ for all } z \in C \right\} \subseteq Z^*$$

denotes the *polar cone* of  $C$ , which is a closed convex cone.

**Theorem 3.27** (First-order necessary optimality conditions). *Let  $X$  and  $Z$  be Banach spaces and  $K \subseteq Z$  be closed and convex. Further, let  $\bar{x}$  be a local solution of (P) at which  $f: X \rightarrow \mathbb{R}$  and  $G: X \rightarrow Z$  are continuously  $F$ -differentiable. Assume that Robinson's constraint qualification (3.2) is satisfied at  $\bar{x}$ .*

*Then there exists a Lagrange multiplier  $\bar{\lambda} \in Z^*$  such that the Karush-Kuhn-Tucker*

conditions for (P)

$$f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} = 0, \quad (3.9)$$

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in T(K, G(\bar{x}))^\circ, \quad (3.10)$$

are satisfied.

December 7, 2022

8 →

*Proof.* To show the existence of a Lagrange multiplier, we define the set  $M \subseteq \mathbb{R} \times Z$  as follows:

$$M = \left\{ \left( \langle f'(\bar{x}), h \rangle_{X^*, X} + \sigma, G'(\bar{x})h - v \right) : h \in X, \sigma \geq 0, v + G(\bar{x}) \in K \right\}$$

The idea is to separate  $\text{int } M$  from the origin  $(0, 0)$  in  $\mathbb{R} \times Z$  and to derive a Lagrange multiplier from the separating hyperplane. Due to the linear appearances of  $h, \sigma$  and  $v$  and the convexity of  $K$ , the set  $M$  is quite obviously convex. Accordingly,  $\text{int } M$  is also convex (why?) and clearly open. In order to apply the Hahn Banach separation theorem as stated in [Proposition 3.5](#), we need to show that  $(0, 0) \notin \text{int } M$  and that  $\text{int } M$  is in fact nonempty.

We first claim that  $(0, 0)$  is a boundary point of  $M$ , so  $(0, 0) \notin \text{int } M$ . It is evident that  $(0, 0) \in M$ , but every open neighborhood of  $(0, 0)$  in  $\mathbb{R} \times Z$  must contain elements which are not in  $M$ , namely at least those of the form  $(-\tau, 0)$  for  $\tau > 0$ . In fact, assume that  $(-\tau, 0) \in M$  for some  $\tau > 0$ . We show that this contradicts the local optimality of  $\bar{x}$  expressed by [Theorem 3.9](#), which was

$$\langle f'(\bar{x}), d \rangle_{X^*, X} \geq 0 \quad \text{for all } d \in T(\mathcal{F}, \bar{x}) = T_\ell(G, K, \bar{x}),$$

where we have already used that Robinson's CQ implies the ACQ ([Theorem 3.22](#)). Indeed, if  $(-\tau, 0) \in M$  for some  $\tau > 0$ , then there exist  $h \in X$  and  $\sigma \geq 0$  with  $\langle f'(\bar{x}), h \rangle_{X^*, X} + \sigma = -\tau < 0$  and  $G'(\bar{x})h = v \in K - G(\bar{x})$ . But then  $h \in T_\ell(G, K, \bar{x})$  and  $\langle f'(\bar{x}), h \rangle_{X^*, X} < 0$ . This is the contradiction.

Next, we show that  $M$  has a nonempty interior. To this end, we use that [\(3.2\)](#) implies [\(3.8\)](#) as in [Lemma 3.25](#). By [\(3.8\)](#), there exists  $\delta > 0$  such that for any  $z \in B_{\delta, Z}(0)$  there exist  $h \in \overline{B_X(0)}$  and  $v \in \overline{K - G(\bar{x})}$  with  $G'(\bar{x})h - v = z$ . Moreover,  $\langle f'(\bar{x}), h \rangle_{X^*, X} \leq \|f'(\bar{x})\|_{X^*}$  due to  $h \in \overline{B_X(0)}$ . This shows that  $[\|f'(\bar{x})\|_{X^*}, \infty) \times B_{\delta, Z}(0) \subseteq M$ , hence  $M$  has nonempty interior.

Now [Proposition 3.5](#) shows that  $(0, 0)$  and  $\text{int } M$  can be separated by a hyperplane  $[(\alpha, z') = \beta]$ , i.e., there exist  $\alpha \in \mathbb{R}$  and  $z' \in Z^*$  with  $(\alpha, z') \neq (0, 0)$  such that

$$\left\langle \begin{pmatrix} \alpha \\ z' \end{pmatrix}, \begin{pmatrix} t \\ z \end{pmatrix} \right\rangle = \alpha t + \langle z', z \rangle_{Z^*, Z} \geq \beta \geq \left\langle \begin{pmatrix} \alpha \\ z' \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right\rangle = 0 \quad \text{for all } (t, z) \in M.$$



Inserting the definition of  $M$ , the foregoing inequality means that

$$\alpha(\langle f'(\bar{x}), h \rangle_{X^*, X} + \sigma) + \langle z', G'(\bar{x})h - v \rangle_{Z^*, Z} \geq 0$$

for all  $h \in X$ ,  $\sigma \geq 0$ ,  $v \in K - G(\bar{x})$ . (3.11)

We will derive the KKT conditions (3.9) and (3.10) from this inequality. For  $h = 0$ ,  $v = 0$  and  $\sigma > 0$  we obtain  $\alpha \geq 0$ . Now assume that  $\alpha = 0$ . Then

$$\langle z', G'(\bar{x})\lambda h - \lambda v \rangle_{Z^*, Z} \geq 0 \quad \text{for all } h \in X, v \in K - G(\bar{x}), \lambda \geq 0,$$

from which we find that for  $\alpha = 0$ , (3.11) implies that

$$\langle z', z \rangle_{Z^*, Z} \geq 0 \quad \text{for all } z \in G'(\bar{x})X - \text{cone}(K, G(\bar{x})).$$

From the (ZKCQ) (3.7) we obtain the contradiction  $z' = 0$ , so we have  $\alpha > 0$ . This allows to multiply (3.11) by  $\alpha^{-1}$  and to obtain, setting  $\bar{\lambda} = \alpha^{-1}z'$ :

$$(\langle f'(\bar{x}), h \rangle_{X^*, X} + \sigma) + \langle \bar{\lambda}, G'(\bar{x})h - v \rangle_{Z^*, Z} \geq 0 \quad \text{for all } h \in X, \sigma \geq 0, v \in K - G(\bar{x}).$$

The choice  $h = 0$  and  $\sigma = 0$  shows that  $\bar{\lambda} \in \text{cone}(K, G(\bar{x}))^\circ = T(K, G(\bar{x}))^\circ$ .

Further, choosing  $\sigma = 0$  and  $v = 0$  shows that

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

which is the same as

$$\langle f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda}, h \rangle_{X^*, X} \geq 0 \quad \text{for all } h \in X.$$

This implies

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

Thus, the existence of a Lagrange multiplier is proved. □

So far, we have derived the KKT conditions under Robinson's CQ and used the equivalences to other conditions as proven in Lemma 3.25. One might wonder how restrictive Robinson's CQ or the equivalent (ZKCQ) (3.7) actually is. The following lemma proves necessary properties for the set of Lagrange multipliers associated to (P) and shows that these properties are also *nearly* sufficient for the (ZKCQ).

**Lemma 3.28.** *Let  $\bar{x} \in \mathcal{F}$  be given and assume that the set of Lagrange multipliers associated to (P) given by*

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

*is nonempty. It is a closed and convex set characterized by the following assertions:*

(1) *If the (ZKCQ)*

$$Z = G'(\bar{x})X - \text{cone}(K, G(\bar{x})) \tag{3.7}$$

or equivalently (3.8) or Robinson's CQ (3.2) are satisfied in  $\bar{x}$ , then  $\Lambda(\bar{x})$  is bounded.

(2) Let conversely  $\Lambda(\bar{x})$  be bounded. Then

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

i.e.,  $G'(\bar{x})X - \text{cone}(K, G(\bar{x}))$  is dense in  $Z$ .

December 14, 2022

9 →

*Proof.* Closedness and convexity of  $\Lambda(\bar{x})$  are evident since the polar cone is closed and convex (as an infinite intersection of closed half spaces) and the defining equation, so  $f'(\bar{x}) + G'(\bar{x})^* \lambda = 0$ , is linear and continuous w.r.t.  $\lambda \in Z^*$ .

We start with (1), so let one of the equivalent constraint qualifications (3.2), (3.7), or (3.8) as well as the KKT conditions (3.9) and (3.10) be satisfied, i.e.,

$$f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} = 0, \quad (3.9)$$

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in T(K, G(\bar{x}))^\circ. \quad (3.10)$$

Due to  $\bar{\lambda} \in T(K, G(\bar{x}))^\circ = \text{cone}(K, G(\bar{x}))^\circ$ , we know that

$$\langle \bar{\lambda}, v - G(\bar{x}) \rangle_{Z^*, Z} \leq 0 \quad \text{for all } v \in K.$$

Moreover, from (3.8) there exists  $\delta > 0$  such that for all  $z \in B_{\delta, Z}(0)$  there are  $h \in \overline{B_X(0)}$  and  $v \in K$  such that  $-z = G(\bar{x}) + G'(\bar{x})h - v$ . Applying  $\bar{\lambda}$  to  $z$  yields

$$\begin{aligned} \langle \bar{\lambda}, z \rangle_{Z^*, Z} &= \langle \bar{\lambda}, v - G(\bar{x}) - G'(\bar{x})h \rangle_{Z^*, Z} \\ &= \langle \bar{\lambda}, v - G(\bar{x}) \rangle_{Z^*, Z} - \langle f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda}, h \rangle_{X^*, X} + \langle f'(\bar{x}), h \rangle_{X^*, X} \\ &\leq \|f'(\bar{x})\|_{X^*} \|h\|_X \leq \|f'(\bar{x})\|_{X^*}. \end{aligned}$$

But this shows that

$$\langle \bar{\lambda}, \bar{z} \rangle_{Z^*, Z} \leq \delta^{-1} \|f'(\bar{x})\|_{X^*} \quad \text{for all } \bar{z} \in B_Z(0)$$

since  $z \in B_{\delta, Z}(0)$  was arbitrary, and thus

$$\|\bar{\lambda}\|_{Z^*} \leq \delta^{-1} \|f'(\bar{x})\|_{X^*}.$$

Since this estimate is uniform in  $\bar{\lambda}$ , the set  $\Lambda(\bar{x})$  is bounded.

Now assume that  $\Lambda(\bar{x})$  is nonempty and bounded. We argue via contradiction, so assume that there exists

$$\bar{z} \in Z \setminus M \quad \text{with} \quad M := \overline{G'(\bar{x})X - \text{cone}(K, G(\bar{x}))}.$$

The set  $M$  is clearly closed and convex and contains 0, so it is nonempty. We use [Proposition 3.5](#) to separate  $\bar{z}$  and  $M$ : There exists an hyperplane  $[z' = \alpha]$  such that

$$\langle z', z \rangle_{Z^*, Z} \geq \alpha \geq \langle z', \bar{z} \rangle_{Z^*, Z} \quad \text{for all } z \in M.$$

(In fact, [Proposition 3.5](#) states that there is a hyperplane  $H$  which even *strictly* separates  $\bar{z}$  and  $M$ . We will however not need the strict separation.) Since  $G'(\bar{x})X - \text{cone}(K, G(\bar{x}))$  is a cone, so is its closure  $M$ , hence the foregoing inequality implies  $\langle z', z \rangle_{Z^*, Z} \geq 0$  for all  $z \in M$ : The right-hand side is a fixed nonpositive number—recall  $0 \in M$ —and we are allowed to scale the left-hand side by an arbitrary number  $\lambda > 0$  by inserting  $\lambda z$  for  $z \in M$ . The inequality can then only be true if  $\langle z', z \rangle_{Z^*, Z} \geq 0$  for all  $z \in M$ , which in turn implies

$$\langle z', G'(\bar{x})h + G(\bar{x}) - v \rangle_{Z^*, Z} \geq 0 \quad \text{for all } h \in X, v \in K. \quad (3.12)$$

We will derive that  $\bar{\lambda} + \beta z' \in \Lambda(\bar{x})$  for every  $\bar{\lambda} \in \Lambda(\bar{x})$  and  $\beta \geq 0$  from this inequality. Choosing  $v = G(\bar{x}) \in K$  in [\(3.12\)](#) shows that

$$\langle z', G'(\bar{x})h \rangle_{Z^*, Z} = \langle G'(\bar{x})^* z', h \rangle_{X^*, X} \geq 0 \quad \text{for all } h \in X$$

and thus (insert  $h \in X$  and  $-h \in X$ )

$$G'(\bar{x})^* z' = 0 \quad \text{in } X^*. \quad (3.13)$$

Conversely, inserting  $h = 0$  in [\(3.12\)](#) implies

$$\langle z', v - G(\bar{x}) \rangle_{Z^*, Z} \leq 0 \quad \text{for all } v \in K,$$

so

$$z' \in \text{cone}(K, G(\bar{x}))^\circ = T(K, G(\bar{x}))^\circ. \quad (3.14)$$

Now finally consider  $\bar{\lambda} \in \Lambda(\bar{x})$  satisfying

$$f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} = 0, \quad (3.9)$$

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in T(K, G(\bar{x}))^\circ. \quad (3.10)$$

From [\(3.13\)](#) and [\(3.14\)](#) we observe that

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

so  $\bar{\lambda} + \beta z' \in \Lambda(\bar{x})$  for every  $\beta \geq 0$ . Letting  $\beta \rightarrow \infty$  gives a contradiction to the boundedness of  $\Lambda(\bar{x})$ .  $\square$

**Remark 3.29.** A couple  $(0, 0) \neq (\alpha, \bar{\lambda}) \in \mathbb{R}^+ \times Z^*$  satisfying the generalized KKT conditions

$$\alpha f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} = 0, \quad (3.15)$$

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in T(K, G(\bar{x}))^\circ. \quad (3.16)$$

is called *generalized Lagrange multiplier*. The functional  $z'$  constructed in the second point of the foregoing proof is a particular instance of a so-called *singular Lagrange multiplier*  $(\alpha, \bar{\lambda}) = (0, z') \in \mathbb{R} \times Z^*$  satisfying the generalized KKT conditions (3.15) and (3.16) for  $\alpha = 0$  and  $\bar{\lambda} \neq 0$ . If such a singular Lagrange multiplier exists, the set  $\Lambda(\bar{x})$  can never be bounded, as seen in the foregoing proof.

Even more, the existence of a singular Lagrange multiplier implies that

$$Z \neq \overline{G'(\bar{x})X - \text{cone}(K, G(\bar{x}))}.$$

This is seen as follows: Let  $(0, \bar{\lambda}) \neq (0, 0)$  be a singular Lagrange multiplier. Then we have by (3.15) and (3.16)

$$\langle \bar{\lambda}, G'(\bar{x})h - v \rangle_{Z^*, Z} \geq 0 \quad \text{for all } h \in X, v \in \text{cone}(K, G(\bar{x})),$$

and thus  $-\bar{\lambda} \in (G'(\bar{x})X - \text{cone}(K, G(\bar{x})))^\circ$  from which there would follow  $\bar{\lambda} = 0$  if  $G'(\bar{x})X - \text{cone}(K, G(\bar{x}))$  was dense in  $Z$  (why?).

**Lemma 3.30.** *If  $\bar{x}$  is a KKT-point, so  $\Lambda(\bar{x}) \neq \emptyset$ , then*

$$\langle f'(\bar{x}), d \rangle_{X^*, X} \geq 0 \quad \text{for all } d \in T_\ell(G, K, \bar{x}).$$

*Proof.* For  $d \in T_\ell(G, K, \bar{x})$  there holds  $G'(\bar{x})d \in T(K, G(\bar{x}))$ . Now with  $\bar{\lambda} \in \Lambda(\bar{x})$ , we have

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

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

The KKT conditions can be written very concisely by means of the Lagrange function which we have already encountered in Nonlinear Optimization:

**Definition 3.31** (Lagrangian). The *Lagrange function* or *Lagrangian*  $L: X \times Z^* \rightarrow \mathbb{R}$  for (P) is given by

$$L(x, \lambda) = f(x) + \langle \lambda, G(x) \rangle_{Z^*, Z}.$$

Using the Lagrangian, the first KKT expression can be expressed quite comfortably by

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

We next turn to particular cases of (P) and their associated KKT conditions.

### 3.2.1 The case of cone constraints

We derive an easier representation of  $T(K, \bar{z})^\circ$  for  $\bar{z} \in K$  when  $K$  is a closed convex cone.

**Reminder:** Recall that a cone  $K$  is convex *if and only if* from  $x, y \in K$  it follows that  $x + y \in K$ . Indeed, if  $K$  is convex, then  $\frac{1}{2}(x + y) \in K$ , and since  $K$  is a cone,  $2 \cdot \frac{1}{2}(x + y) = x + y \in K$ . Conversely, if  $x, y \in K$ , then also  $(1 - t)x \in K$  and  $ty \in K$  if  $t \in [0, 1]$ , and then by assumption also  $(1 - t)x + ty \in K$  and  $K$  is convex.

This is done using the *annihilator*.

**Reminder:** The *annihilator*  $A^\perp$  of a set  $A \subseteq X$  is given by

$$A^\perp := \left\{ x' \in X^* : \langle x', x \rangle = 0 \text{ for all } x \in A \right\},$$

so the collection of all functionals  $x' \in X^*$  for which  $A \subseteq \ker x'$ .

**Lemma 3.32.** *If  $K$  is a closed convex cone and  $\bar{z} \in K$ , then*

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

*Proof.* The first equality follows from convexity of  $K$  (Lemma 3.8) and  $A^\circ = \overline{A}^\circ$  for every set  $A$ ; see the exercises.

For the second equality: For every  $z' \in K^\circ \cap \{\bar{z}\}^\perp$ , there holds for all  $t > 0$  and all  $z \in K$ :

$$\langle z', t(z - \bar{z}) \rangle_{Z^*, Z} = t \langle z', z \rangle_{Z^*, Z} - t \langle z', \bar{z} \rangle_{Z^*, Z} = t \langle z', z \rangle_{Z^*, Z} \leq 0.$$

Hence,  $\text{cone}(K, \bar{z})^\circ \supseteq K^\circ \cap \{\bar{z}\}^\perp$ .

Conversely, let  $z' \in \text{cone}(K, \bar{z})^\circ$ . Due to  $K$  being closed, we have  $0 \in K$  (why?) and thus  $-\bar{z} = 1 \cdot (0 - \bar{z}) \in \text{cone}(K, \bar{z})$ . This implies  $\langle z', -\bar{z} \rangle_{Z^*, Z} \leq 0$ . On the other hand, as above,

$$\langle z', t(z - \bar{z}) \rangle_{Z^*, Z} \leq 0 \quad \text{for all } z \in K, t > 0. \quad (3.17)$$

Since  $K$  is a cone,  $z = 2\bar{z} \in K$ , so  $\langle z', \bar{z} \rangle_{Z^*, Z} \leq 0$ . But this means that both  $\langle z', \bar{z} \rangle_{Z^*, Z} \leq 0$  and  $\langle z', -\bar{z} \rangle_{Z^*, Z} \leq 0$ , hence  $z' \in \{\bar{z}\}^\perp$ . Then, in (3.17),  $\langle z', z \rangle_{Z^*, Z} \leq 0$  for all  $z \in K$ , which is exactly the definition of  $z' \in K^\circ$ .  $\square$

Hence, in the case of a closed convex cone  $K$ , the KKT condition

$$G(\bar{x}) \in K, \quad \bar{\lambda} \in T(K, G(\bar{x}))^\circ \quad (3.10)$$

can be written equivalently as a cone *complementarity condition*

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

As mentioned earlier, when  $K$  is a closed convex cone, then the condition  $G(x) \in K$  can be viewed as an abstract inequality constraint. Indeed, we define a relation  $\leq_K$  induced by  $-K$  by

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

If  $K$  is *pointed* (**spitz**), i.e., if  $z, -z \in K$  implies that  $z = 0$ , then this is indeed a partial ordering in which we can cancel positive factors on both sides and have the relation preserved under adding any element to both sides; in fact, it is enough to suppose that  $0 \in K$  instead of  $K$  closed. (A non-pointed cone  $K$  is also sometimes called *flat* or *blunt*; there is an unfortunate amount of different notions for properties of cones in general.) A partial ordering with the foregoing properties also defines a cone. See the exercises.

**Example 3.33.** In a function space  $X$  consisting of functions defined on some set  $\Omega \subset \mathbb{R}^d$ , the cone of nonpositive functions

$$K_- := \left\{ f \in X : f(x) \leq 0 \text{ for almost all } x \in \Omega \right\}$$

is pointed and induces the usual pointwise ordering of functions, so  $f \leq_{K_-} g$  if and only if  $f(x) \leq g(x)$  for almost all  $x \in \Omega$ .

**Remark 3.34.** The seemingly unnecessarily complicated definition of the ordering  $\leq_K$  is tailored to the classical notion of nonlinear programs:  $z \leq_K 0$  means exactly that  $z \in K$ ; so in case of cone constraints the standard constraint  $G(x) \in K$  can be written as  $G(x) \leq_K 0$ .

Note however that the notation can be misleading, as the partial ordering induced by the cone of *nonnegative* functions  $K_+$  (cf. [Example 3.33](#)) is given by  $f \leq_{K_+} g$  if

and only if  $f(x) \geq g(x)$  for almost all  $x \in \Omega \dots$

Using the cone ordering notation from above, we can rewrite the KKT conditions (3.10) yet again to

$$G(\bar{x}) \leq_K 0, \quad \bar{\lambda} \geq_{K^+} 0, \quad \langle \bar{\lambda}, G(\bar{x}) \rangle_{Z^*, Z} = 0,$$

where  $K^+ := -K^\circ$  is the *dual cone* to  $K$ . The last representation quite exactly resembles the classical KKT conditions from Nonlinear Optimization.

### 3.2.2 The Slater condition

We now consider the case of a convex problem, for which we first give a notion of convexity in function spaces.

**Definition 3.35** (Generalized convexity). Let  $K \subseteq Z$  be a closed convex cone. We say that  $G: X \rightarrow Z$  is *convex with respect to*  $-K$  (or  $\leq_K$ ), if

$$G((1-t)x + ty) \leq_K (1-t)G(x) + tG(y) \quad \text{for all } x, y \in X, t \in [0, 1],$$

or equivalently

$$(1-t)G(x) + tG(y) - G((1-t)x + ty) \in -K \quad \text{for all } x, y \in X, t \in [0, 1].$$

In case of  $F$ -differentiable and convex  $G$ , the *Slater constraint qualification* (or *Slater condition*), well-known from Nonlinear Optimization,

$$\text{there exists a } \mathbf{x} \in \mathcal{F}: \quad G(\mathbf{x}) \in \text{int } K \tag{3.18}$$

implies the Robinson CQ:

**Lemma 3.36.** *Let  $K \subseteq Z$  be a closed convex cone and let  $G$  be  $F$ -differentiable and convex with respect to  $\leq_K$ . Assume that the Slater condition (3.18) is satisfied. Then RCQ (3.2) in the form of the Linearized Slater CQ (3.4) is satisfied in every feasible point  $\bar{x} \in \mathcal{F} = G^{-1}[K]$ .*

*Proof.* Using convexity of  $G$  and the fact that  $K$  is a cone, we observe that for every  $t \in (0, 1]$  and  $\bar{x} \in \mathcal{F}$  we have

$$\frac{G((1-t)\bar{x} + t\mathbf{x}) - (1-t)G(\bar{x}) - tG(\mathbf{x})}{t} \in K.$$

Taking the limit  $t \searrow 0$  gives (note that  $K$  is closed)

$$G(\bar{x}) + G'(\bar{x})(\mathbf{x} - \bar{x}) - G(\mathbf{x}) \in K \quad \iff \quad G(\bar{x}) + G'(\bar{x})(\mathbf{x} - \bar{x}) \in K + G(\mathbf{x}).$$

By the Slater condition, there exists  $\varepsilon > 0$  with  $G(\mathbf{x}) + B_{\varepsilon, Z}(0) \subseteq K$ . Therefore,

$$G(\bar{x}) + G'(\bar{x})(\mathbf{x} - \bar{x}) + B_{\varepsilon, Z}(0) \subseteq K + G(\mathbf{x}) + B_{\varepsilon, Z}(0) \subseteq K + K = K.$$

Thus,

$$G(\bar{x}) + G'(\bar{x})(\mathbf{x} - \bar{x}) \in \text{int } K,$$

which is exactly the LSCQ (3.4) and thus equivalent to the Robinson CQ by Lemma 3.18.  $\square$

**Remark 3.37.** In the proof of Lemma 3.36, we have shown along the way that the differentiable convex function  $G$  satisfies

$$G(y) - G(x) - G'(x)(y - x) \in -K \quad \text{or} \quad G'(x)(y - x) \leq_K G(y) - G(x) \quad (3.19)$$

for all  $x, y \in \mathcal{F} = G^{-1}[K]$ . This is exactly the analogue of the well-known characterization of classical convex functions  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , given by

$$\nabla f(x)^T(y - x) \leq f(y) - f(x).$$

The proof that (3.19) also implies convexity works again analogously to the classical case.

### 3.2.3 Applications

We proceed by giving two examples for the KKT conditions in an optimal control setting. The first one is still of rather abstract nature and incorporates control constraints, while the second one is slightly more specialized and has state constraints.

**Reminder:** Let  $H$  be a Hilbert space. The *Fréchet-Riesz representation theorem* says that there exists an isometric isomorphism  $T \in \mathcal{L}(H^*; H)$ —the *Riesz isomorphism*—such that

$$\langle g, v \rangle_{H^*, H} = (Tg, v)_H \quad \text{for all } g \in H^*, v \in H,$$

where  $(\cdot, \cdot)_H$  is the inner product on the Hilbert space  $H$ . In particular,  $\|g\|_{H^*} = \|Tg\|_H$ . In this sense, we can always identify a Hilbert space  $H$  with its dual  $H^*$  up to the application of the Riesz isomorphism.

**Example 3.38** (Optimal control problem with control constraints). We consider a control-constrained optimal control problem

$$\min_{(y, u) \in Y \times U} J(y, u) \quad \text{s.t.} \quad E(y, u) = 0, \quad u \in U_{\text{ad}} \quad (3.20)$$



governed by the state equation

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

where  $E: Y \times U \rightarrow W$  is continuously differentiable,  $Y$  and  $W$  are Banach spaces, and  $U = L^2(\Upsilon)$  for a measure space  $\Upsilon$ . The objective function  $J: Y \times U \rightarrow \mathbb{R}$  is assumed to be F-differentiable and the control constraints are given by

$$U_{\text{ad}} = \left\{ u \in U : a \leq u \leq b \text{ a.e. in } \Upsilon \right\}$$

with  $a, b \in L^2(\Upsilon)$  and  $a \leq b$  almost everywhere on  $\Upsilon$ , so  $U_{\text{ad}} \neq \emptyset$ . (Usually we will have  $\Upsilon = \Omega \subseteq \mathbb{R}^n$  for a bounded domain  $\Omega \subseteq \mathbb{R}^n$  with the Lebesgue measure or  $\Upsilon = \partial\Omega$  with the boundary measure.)

Let  $(\bar{y}, \bar{u}) \in Y \times U$  be a local solution to (3.20). We need a constraint qualification to be satisfied in order to derive a KKT characterization for  $(\bar{y}, \bar{u})$ . To this end, we first transfer the problem to standard form. This is obtained by identifying

$$X = Y \times U, \quad Z = W \times U \quad \text{and} \quad K = \{0_W\} \times U_{\text{ad}} \subseteq Z,$$

as well as  $\bar{x} = (\bar{y}, \bar{u})$  and  $G: X \rightarrow Z$  given by  $G(x) = \begin{pmatrix} E(y,u) \\ u \end{pmatrix}$ .

We have seen in Proposition 3.19 that  $G'(\bar{x})$  being surjective is a constraint qualification by implying RCQ. Note that  $\text{int} U_{\text{ad}} = \emptyset$  in  $U = L^2(\Upsilon)$  as seen in the exercises, so there is no other (practical) characterization of RCQ than surjectivity of  $G'(\bar{x})$  available. To show that  $G'(\bar{x})$  is indeed surjective, we need to show that for any  $(w, u) \in Z = W \times U$  there exists  $(h_y, h_u) \in X = Y \times U$  such that

$$G'(\bar{x})h = \begin{pmatrix} E'_y(\bar{y}, \bar{u}) & E'_u(\bar{y}, \bar{u}) \\ 0 & \text{id}_U \end{pmatrix} \begin{pmatrix} h_y \\ h_u \end{pmatrix} = \begin{pmatrix} w \\ u \end{pmatrix},$$

where we have identified  $G'(\bar{x})$  with its Jacobian matrix type representation. (See the exercises.) Due to the upper triangular form of this Jacobian of  $G'(\bar{x})$ , it will turn out that it is both sufficient and necessary to assume that  $E'_y(\bar{y}, \bar{u})$  is surjective in order to have  $G'(\bar{x})$  surjective. In fact, looking at the second row in the foregoing equality, we immediately see that necessarily  $h_u = u$ . Thus, surjectivity of  $G'(\bar{x})$  is equivalent to, for every  $(w, u) \in W \times U$ , being able to find  $h_y \in Y$  such that

$$E'_y(\bar{y}, \bar{u})h_y = w - E'_u(\bar{y}, \bar{u})u.$$

But this is exactly the question of surjectivity of  $E'_y(\bar{y}, \bar{u})$ . (Choose any  $u \in U$  and consider  $w := E'_u(\bar{y}, \bar{u})u + v$  for arbitrary  $v \in W$ .)

So, if we *assume* that

$$E'_y(\bar{y}, \bar{u}) \text{ is surjective,} \tag{3.21}$$

then  $G'(\bar{x})$  is surjective and RCQ for (3.20) is satisfied. Hence, there exist a Lagrange multiplier (pair)

$$\bar{\lambda} = (\bar{p}, \bar{\mu}) \in Z^* = (W \times U)^* = W^* \times U^* = W^* \times L^2(\Upsilon)$$

such that (3.9) and (3.10) are satisfied, i.e.,

$$\begin{aligned} J'_y(\bar{y}, \bar{u}) + E'_y(\bar{y}, \bar{u})^* \bar{p} &= 0 && \text{in } Y^*, \\ J'_u(\bar{y}, \bar{u}) + E'_u(\bar{y}, \bar{u})^* \bar{p} + \bar{\mu} &= 0 && \text{in } L^2(\Upsilon), \\ E(\bar{y}, \bar{u}) &= 0 && \text{in } W, \\ \bar{u} \in U_{\text{ad}}, \quad \bar{\mu} \in T(U_{\text{ad}}, \bar{u})^\circ. &&& \end{aligned}$$

Note that  $\{0_W\}^\circ = W^*$ , thus the condition  $\bar{p} \in \{0_W\}^\circ$  is void. Here we have used that, with the Jacobian matrix type representations for  $G'(\bar{x})$  and  $f'(\bar{x})$ ,

$$f'(\bar{x}) = \begin{pmatrix} J_y(\bar{y}, \bar{u}) \\ J_u(\bar{y}, \bar{u}) \end{pmatrix} \quad \text{and} \quad G'(\bar{x})^* = \begin{pmatrix} E_y(\bar{y}, \bar{u})^* & 0 \\ E_u(\bar{y}, \bar{u})^* & \text{id}_U \end{pmatrix}.$$

Both  $E_y(\bar{y}, \bar{u})$  and its adjoint correspond again to (partial) differential operators if  $E$  did so; more precisely, they correspond to the linearization of the operator represented by  $E$  and the adjoint of its linearization. See the coming [Example 3.39](#).

We next consider  $T(U_{\text{ad}}, \bar{u})^\circ$ . Since  $U_{\text{ad}}$  is convex, we have  $T(U_{\text{ad}}, \bar{u}) = \overline{\text{cone}(U_{\text{ad}}, \bar{u})}$  ([Lemma 3.8](#)). Using this, it is easy to see (and an exercise) that

$$T(U_{\text{ad}}, \bar{u}) = \left\{ h \in L^2(\Upsilon) : h|_{[\bar{u}=a]} \geq 0, h|_{[\bar{u}=b]} \leq 0 \right\},$$

where  $[\bar{u} = a] = \{x \in \Upsilon : \bar{u}(x) = a(x)\}$  and analogously for  $[\bar{u} = b]$ . Thus,  $T(U_{\text{ad}}, \bar{u})^\circ$  consists exactly of all functions  $s \in L^2(\Upsilon)$  such that

$$\int_{\Upsilon} h(x)s(x) dx \leq 0 \quad \text{for all } h \in L^2(\Upsilon) \text{ with } h|_{[\bar{u}=a]} \geq 0, h|_{[\bar{u}=b]} \leq 0.$$

We are going to derive pointwise properties for the functions  $s$  from this integral variational inequality. In fact, we are going to show that

$$T(U_{\text{ad}}, \bar{u})^\circ = \left\{ s \in L^2(\Upsilon) : s|_{[\bar{u}=a]} \leq 0, s|_{[\bar{u}=b]} \geq 0, s|_{[a < \bar{u} < b]} = 0 \right\}.$$

It is easily verified that an element  $s$  of the set on the right will satisfy the above variational inequality, so “ $\supseteq$ ” is certainly true. We show the other inclusion. So let  $s \in T(U_{\text{ad}}, \bar{u})^\circ$ . Consider  $M \subseteq [\bar{u} = a]$ . Then  $\chi_M \in T(U_{\text{ad}}, \bar{u})$ , where

$$\chi_M(x) = \begin{cases} 1 & \text{if } x \in M, \\ 0 & \text{otherwise.} \end{cases}$$

Suppose that  $s|_M > 0$  if  $M$  has nonzero measure. We obtain

$$\int_{\Upsilon} \chi_M(x) s(x) \, dx = \int_M s(x) \, dx > 0,$$

which is a contradiction to  $s \in T(U_{\text{ad}}, \bar{u})^\circ$ . Thus  $s|_{[\bar{u}=a]} \leq 0$  almost everywhere. By the analogous argument with  $-\chi_M \in T(U_{\text{ad}}, \bar{u})$  for  $M \subseteq [\bar{u} = b]$ , we find that then  $s|_{[\bar{u}=b]} \geq 0$  almost everywhere. Finally, consider  $M \subseteq [a < \bar{u} < b]$ . Then  $\pm \chi_M \in T(U_{\text{ad}}, \bar{u})$ , and it follows that  $s|_{[a < \bar{u} < b]} = 0$ .

Hence, the multiplier  $\bar{\mu}$  for the control constraints is an  $L^2(\Upsilon)$  function satisfying

$$\bar{\mu}|_{[\bar{u}=a]} \leq 0, \quad \bar{\mu}|_{[\bar{u}=b]} \geq 0, \quad \bar{\mu}|_{[a < \bar{u} < b]} = 0. \quad (3.22)$$

A potential and common interpretation here is that  $\bar{\mu}$  acts like an indicator for how the restriction upon  $\bar{u}$  induced by the respective constraints  $a \leq \bar{u}$  and  $\bar{u} \leq b$  acts. Of course, on the set  $[a < \bar{u} < b]$ , the constraints do not restrict  $\bar{u}$  at all, which is indicated by  $\bar{\mu} = 0$  on this set.

**Example 3.39.** We consider the general KKT-conditions derived in the foregoing example for the semilinear elliptic optimal control problem as in [Example 2.9](#). Then  $Y = H^1(\Omega)$  and  $\Upsilon = \partial\Omega$  with the boundary measure; also  $U = L^2(\partial\Omega)$  and  $W = H^1(\Omega)^* \times L^2(\partial\Omega)$ . In [Example 3.4](#), we have seen that for this problem, the linearized operator  $E'_y(\bar{y}, \bar{u}) \in \mathcal{L}(H^1(\Omega); H^1(\Omega)^*)$  and the equation  $E'_y(\bar{y}, \bar{u})z = w$  give rise to the weak formulation of the *linear* PDE

$$\begin{aligned} -\Delta z + 3\bar{y}^2 z &= f && \text{on } \Omega, \\ \frac{\partial z}{\partial \nu} + z &= g && \text{on } \partial\Omega \end{aligned}$$

for the functional  $w \in H^1(\Omega)^*$  defined by

$$\langle w, v \rangle := \int_{\Omega} f(x)v(x) \, dx + \int_{\partial\Omega} g(x)v(x) \, d\omega(x).$$

Hence, the assumption that  $E'_y(\bar{y}, \bar{u})$  be surjective in truth corresponds to the question whether we can solve (in a weak sense) the above PDE problem for all suitable right-hand sides  $f$  and  $g$ . With the Lax-Milgram lemma, we can easily show that this is the case for every  $(\bar{y}, \bar{u})$ . The associated bilinear form including the Robin boundary condition is uniformly coercive since the zero-order term is induced by  $3\bar{y}^2$ , which is nonnegative, and there is a strictly positive factor  $1 \cdot z$  in the Robin boundary condition term. (Use the generalized Friedrichs' inequality [\(2.2\)](#).)

For this particular example, one moreover easily checks from the formula for the

derivative  $E'_y(\bar{y}, \bar{u}) \in \mathcal{L}(H^1(\Omega); H^1(\Omega)^*)$  in  $(\bar{y}, \bar{u})$ , which was given by

$$E'_y(\bar{y}, \bar{u})z: \left[ v \mapsto (\nabla z, \nabla v)_{L^2(\Omega)^n} + (z, v)_{L^2(\partial\Omega)} + 3(\bar{y}^2 z, v)_{L^2(\Omega)} \right] \in H^1(\Omega)^*,$$

that the adjoint operator  $E'_y(\bar{y}, \bar{u})^* \in \mathcal{L}(H^1(\Omega); H^1(\Omega)^*)$  in fact induces the *same* weak formulation. This is however not generic; here it happens because we consider the particular case of the Laplacian, and because there is only an additional zero-order term  $3\bar{y}^2 z$  as opposed to possible first-order terms. (See also the next example.)

We further recall that  $J(y, u) = \frac{1}{2}\|y - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha}{2}\|u\|_{L^2(\partial\Omega)}^2$ . Thus, by [Example 3.2](#) (bilinear form!),  $J$  is continuously F-differentiable  $L^2(\Omega) \times L^2(\partial\Omega) \rightarrow \mathbb{R}$ , and we have

$$J'_y(\bar{y}, \bar{u})z = (\bar{y} - y_d, z)_{L^2(\Omega)}, \quad \text{so} \quad J'_y(\bar{y}, \bar{u}) \cong \bar{y} - y_d$$

and

$$J'_u(\bar{y}, \bar{u})h = \alpha(\bar{u}, h)_{L^2(\partial\Omega)}, \quad \text{so} \quad J'_u(\bar{y}, \bar{u}) \cong \alpha\bar{u}$$

where we have used the Fréchet-Riesz isomorphism for each case to identify the derivatives with the respective Hilbert space element. (In this sense, the  $\cong$  identification yields the respective *gradient*.)

Recall from [Example 3.4](#) that  $E'_u(\bar{y}, \bar{u})h$  was given by  $-(h, \cdot)_{L^2(\partial\Omega)} \in L^2(\Omega)^*$ , and, thus, so is  $E'_u(\bar{y}, \bar{u})^*p = -(p, \cdot)_{L^2(\partial\Omega)}$ . It follows that the gradient equations

$$\begin{aligned} J'_y(\bar{y}, \bar{u}) + E'_y(\bar{y}, \bar{u})^*\bar{p} &= 0 && \text{in } H^1(\Omega)^*, \\ J'_u(\bar{y}, \bar{u}) + E'_u(\bar{y}, \bar{u})^*\bar{p} + \bar{\mu} &= 0 && \text{in } L^2(\partial\Omega) \end{aligned}$$

in the KKT conditions in [Example 3.38](#) correspond to the weak formulation of

$$\begin{aligned} -\Delta\bar{p} + 3\bar{y}^2\bar{p} &= y_d - \bar{y} && \text{on } \Omega, \\ \frac{\partial\bar{p}}{\partial\nu} + \bar{p} &= 0 && \text{on } \partial\Omega \end{aligned}$$

and, with deliberate use of the Fréchet-Riesz isomorphism for  $L^2(\partial\Omega)$ ,

$$\alpha\bar{u} - \bar{p} + \bar{\mu} = 0 \quad \text{on } \partial\Omega.$$

Note that the latter gives an explicit formula for the desired object  $\bar{u}$  in terms of the so-called *adjoint state*  $\bar{p} \in H^1(\Omega)$  which is again a weak solution of a PDE. In the most simple case when  $\bar{\mu}$ —this was the Lagrange multiplier associated to  $U_{\text{ad}}$ —is zero, this shows that the optimal control  $\bar{u}$  is in fact as regular as the boundary trace of an  $H^1(\Omega)$  function, so in  $H^{1/2}(\partial\Omega)$ .

But even with nonzero  $\bar{\mu}$  associated to box constraints  $U_{\text{ad}}$  as in [\(3.22\)](#) derived in [Example 3.4](#), we obtain a very useful formula for  $\bar{u}$ , which is that for  $\omega$ -almost all

$x \in \partial\Omega$ , we have

$$\bar{u}(x) = \text{proj}_{[a(x), b(x)]} \left( \frac{1}{\alpha} \bar{p}(x) \right) = \begin{cases} a(x) & \text{if } \frac{1}{\alpha} \bar{p}(x) \leq a(x), \\ \frac{1}{\alpha} \bar{p}(x) & \text{if } a(x) < \frac{1}{\alpha} \bar{p}(x) < b(x), \\ b(x) & \text{if } \frac{1}{\alpha} \bar{p}(x) \geq b(x). \end{cases}$$

This is a very interesting formula with far-reaching implications. (Check for instance what happens for  $\alpha \searrow 0$ .) Also, we see that a certain regularity transfer from  $\bar{p}$  to  $\bar{u}$  still happens also in the case when  $\bar{\mu} \neq 0$ ; in particular, this is the case if the functions  $a$  and  $b$  defining  $U_{\text{ad}}$  are sufficiently regular, for example, also in  $H^{1/2}(\partial\Omega)$ , so boundary traces of  $H^1(\Omega)$  functions. See the exercises.

January 18, 2023

← 12

**Example 3.40** (Elliptic optimal control problem with state constraints). We consider an optimal control problem with an abstract elliptic state equation, control on the right-hand side and pointwise state constraints, that is, the constraints are of the form

$$Ay = Bu + b \quad \text{and} \quad y \leq \psi.$$

The state equation is to be seen as an abstract form of an elliptic partial differential equation in weak formulation by the following assumptions:

- $\Omega \subset \mathbb{R}^n$ , where  $1 \leq n \leq 3$ , is a bounded Lipschitz domain,
- $B \in \mathcal{L}(U; L^2(\Omega))$ , where the control space  $U$  is a Banach space, and  $b \in L^2(\Omega)$ ,
- $A \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega))$  with  $A^{-1} \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$ , and additionally, the mapping  $y \mapsto Ay$  defines a bounded, injective and surjective operator from the state space  $Y = H_0^1(\Omega) \cap H^2(\Omega)$  to  $L^2(\Omega)$ . As a consequence,  $v \mapsto A^{-1}v$  defines the *solution operator*  $S \in \mathcal{L}(L^2(\Omega); Y)$ . See also [Remark 3.41](#) below.

For the upper bound  $\psi$  in the state constraint we suppose  $\psi \in C(\bar{\Omega})$ , and we describe the state constraint by

$$\mathcal{E}y - \psi \in K_-$$

with the well-known cone of nonnegative functions in  $C(\bar{\Omega})$

$$K_- = \left\{ q \in C(\bar{\Omega}) : q(x) \leq 0 \text{ for all } x \in \bar{\Omega} \right\}$$

and the embedding  $\mathcal{E} \in \mathcal{L}(Y; C(\bar{\Omega}))$  granted by the Sobolev embedding theorem. It will be crucial that  $\text{int} K_- \neq \emptyset$  in  $C(\bar{\Omega})$ . Together with the state equation  $E: Y \times U \rightarrow L^2(\Omega)$  given by  $E(y, u) = Ay - Bu - b$ , we collect all constraints in the function  $G: X \rightarrow Z$  by setting  $X = Y \times U$  and  $Z = L^2(\Omega) \times C(\bar{\Omega})$  and  $G(x) = \begin{pmatrix} E(y, u) \\ \mathcal{E}y - \psi \end{pmatrix} \in K$  with  $K = \{0_{L^2(\Omega)}\} \times K_-$  and  $x = (y, u) \in X$ . Then  $G_1(x) = E(y, u)$  and  $G_2(x) = \mathcal{E}y - \psi$ .

Since  $\text{int } K_- \neq \emptyset$  in  $C(\overline{\Omega})$ , by (3.5) in Proposition 3.19 the following linearized Slater-type assumption will be sufficient for RCQ at *any* feasible point  $(\bar{y}, \bar{u})$  for this problem: There exist  $\hat{u} \in U$  and  $\hat{y} \in Y$  with  $A\hat{y} = B\hat{u} + b$  and  $\mathcal{E}\hat{y} - \psi < 0$  on  $\overline{\Omega}$ , so  $(\hat{y}, \hat{u})$  is feasible for the equality constraint and satisfies the inequality constraint strictly.

Indeed, since  $E'_y(\bar{y}, \bar{u}) = A: Y \rightarrow L^2(\Omega)$  is continuously invertible, it is in particular surjective, and so is  $E'(\bar{y}, \bar{u}) = G'_1(\bar{x})$  for *all*  $\bar{x} = (\bar{y}, \bar{u}) \in Y \times U = X$ . Let  $(\bar{y}, \bar{u}) \in Y \times U$  be an arbitrary feasible pair. Then  $(z, h) := (\hat{y} - \bar{y}, \hat{u} - \bar{u})$  satisfies

$$E'(\bar{y}, \bar{u})(z, h) = Az - Bh = A\hat{y} - B\hat{u} - (A\bar{y} - B\bar{u}) = b - b = 0.$$

Hence, we have  $(z, h) \in \ker E'(\bar{y}, \bar{u}) = \ker G'_1(\bar{x})$ . In particular,

$$G_1(\bar{x}) + G'_1(\bar{x})(z, h) = 0 + 0 = 0.$$

Moreover, from  $\mathcal{E}\hat{y} - \psi \in \text{int } K_-$  we infer

$$G_2(\bar{x}) + G'_2(\bar{x})(z, h) = \mathcal{E}\bar{y} - \psi + \mathcal{E}z = \mathcal{E}\hat{y} - \psi \in \text{int } K_-.$$

In (3.5) in Proposition 3.19 we have seen that surjectivity of  $E'(\bar{y}, \bar{u}) = G_1(\bar{x})$  and the existence of  $(z, h) \in Y \times U$  with  $G_1(\bar{x}) + G'_1(\bar{x})(z, h) = 0$  and  $G_2(\bar{x}) + G'_2(\bar{x})(z, h) \in \text{int } K_-$  is sufficient for (in fact, even equivalent to) RCQ.

Hence, under the assumption that  $(\hat{y}, \hat{u})$  as above exists and letting  $J: Y \times U \rightarrow \mathbb{R}$  be F-differentiable, we obtain a KKT characterization of *any* locally optimal control  $(\bar{y}, \bar{u})$  for the problem

$$\min_{(y,u) \in Y \times U} J(y, u) \quad \text{s.t.} \quad Ay = Bu + b, \quad y \leq \psi$$

as follows:

There exist  $\bar{p} \in L^2(\Omega)^* = L^2(\Omega)$  and  $\bar{\mu} \in C(\overline{\Omega})^*$  such that

$$\begin{aligned} J'_y(\bar{y}, \bar{u}) + A^*\bar{p} + \mathcal{E}^*\bar{\mu} &= 0 && \text{in } Y^*, \\ J'_u(\bar{y}, \bar{u}) - B^*\bar{p} &= 0 && \text{in } U^*, \\ A\bar{y} &= B\bar{u} + b && \text{in } L^2(\Omega), \\ \mathcal{E}\bar{y} - \psi &\leq 0 && \text{in } C(\overline{\Omega}), \\ \bar{\mu} \in T(K_-, \mathcal{E}\bar{y} - \psi)^\circ &= K_-^\circ \cap (\mathcal{E}\bar{y} - \psi)^\perp && \text{in } C(\overline{\Omega})^*, \end{aligned}$$

where the last equality follows from Lemma 3.32, since  $K_-$  is a closed convex cone.

The last conditions can be rewritten as a complementarity condition

$$\mathcal{E}\bar{y} - \psi \leq 0, \quad \bar{\mu} \in K_-^\circ, \quad \langle \bar{\mu}, \mathcal{E}\bar{y} - \psi \rangle_{C(\overline{\Omega})^*, C(\overline{\Omega})} = 0.$$

Here, one can show [La93, Ch. IX, Thm. 4.2] that there indeed holds  $C(\bar{\Omega})^* = M(\bar{\Omega})$ , where  $M(\bar{\Omega})$  is the space of (real, signed) regular Borel measures on  $\bar{\Omega}$  with the *total variation* norm

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

Thereby, the dual pairing of  $\mu \in M(\bar{\Omega}) = C(\bar{\Omega})^*$  with  $f \in C(\bar{\Omega})$  is given by

$$\langle \mu, f \rangle_{M(\bar{\Omega}), C(\bar{\Omega})} = \int_{\bar{\Omega}} f(x) d\mu(x).$$

Further,  $\bar{\mu} \in K_-^\circ$  means that  $\bar{\mu} \geq 0$  in the sense of  $\langle \bar{\mu}, q \rangle \leq 0$  for all functions  $q \in K_-$ .

From these properties it follows that the above complementarity condition can be written as

$$\mathcal{E}\bar{y} - \psi \leq 0, \quad \bar{\mu} \geq 0, \quad \bar{\mu}([\mathcal{E}\bar{y} - \psi < 0]) = 0,$$

so that the support of  $\bar{\mu}$  is concentrated on  $[\mathcal{E}\bar{y} = \psi]$ . In particular,

$$T(K_-, \mathcal{E}\bar{y} - \psi)^\circ = K_-^\circ \cap (\mathcal{E}\bar{y} - \psi)^\perp = \left\{ \mu \in M(\bar{\Omega}) : \mu \geq 0, \mu([\mathcal{E}\bar{y} - \psi < 0]) = 0 \right\}.$$

**Remark 3.41.** The regularity assumptions on  $A$  in Example 3.40 are to be understood as follows: The basic assumption says that for every  $f \in H^{-1}(\Omega)$ , there exist a unique  $y \in H_0^1(\Omega)$  such that  $Ay = f$  and a constant  $C > 0$  independent of  $f$  such that  $\|y\|_{H_0^1(\Omega)} \leq C\|f\|_{H^{-1}(\Omega)}$ .

This is the classical result obtained by the Lax-Milgram lemma for large classes of elliptic partial differential operators [Br10, Ch. 5.3], such as exemplarily the divergence-gradient operators  $y \mapsto -\operatorname{div}(\mu \nabla y)$ , complemented with homogeneous Dirichlet boundary conditions, in their weak form given by

$$\langle Ay, \varphi \rangle = \int_{\Omega} (\mu \nabla y) \cdot \nabla \varphi dx \quad \text{for } \varphi \in H_0^1(\Omega),$$

where  $\mu \in L^\infty(\Omega; \mathbb{R}^{n \times n})$  takes its values in the space of  $n \times n$ -matrices and satisfies the *coercivity- or ellipticity condition*: There exists  $\alpha > 0$  such that

$$v^T \mu(x) v \geq \alpha \|v\|_2^2 \quad \text{for all } v \in \mathbb{R}^n \quad \text{for almost all } x \in \Omega.$$

(Compare also with the weak form of the negative Laplacian  $-\Delta$ , so the divergence-gradient operator with  $\mu$  being the  $n \times n$ -identity matrix, in Example 2.9.) See the exercises.

Since the underlying partial differential operators are of order two, the assumption  $A^{-1} \in \mathcal{L}(H^{-1}(\Omega); H_0^1(\Omega))$  can be seen as a *maximal Sobolev regularity* result in the sense that  $A$  and  $A^{-1}$  operate exactly between Sobolev spaces with a gap of

differentiability of order two, namely  $H_0^1(\Omega)$  and  $H^{-1}(\Omega)$ . The particular power of the Lax-Milgram lemma manifests in the fact that a second-order elliptic partial differential operator will always have this nice property with respect to the Hilbert spaces  $H_0^1(\Omega)$  and  $H^{-1}(\Omega)$ .

The additional assumption on  $A$  then requires the following: If the right-hand side  $f$  in the equation  $Ay = f$  is actually from the better space  $L^2(\Omega) \hookrightarrow H^{-1}(\Omega)$ , then this additional  $L^2$ -regularity for the data implies additional  $H^2$ -regularity for the state  $y$ , so  $y \in H_0^1(\Omega) \cap H^2(\Omega)$  with a continuous dependence of  $y$  on  $f$ , in the sense of the existence of a constant  $C > 0$  independent of  $f$  such that

$$\|y\|_{H_0^1(\Omega)} + \|y\|_{H^2(\Omega)} \leq C\|f\|_{L^2(\Omega)}.$$

Of course, this additional assumption is again a maximal Sobolev regularity property, but unfortunately, proving such a property is quite difficult, even for the Laplacian. It can be shown, e.g., if  $\Omega$  is of class  $C^{1,1}$  and the coefficients  $\mu_{ij}$  of the operator  $A$  are uniformly continuous on  $\bar{\Omega}$ , cf. [GT01, Thm. 9.15].

**Remark 3.42.** In the foregoing [Example 3.40](#), consider again the KKT condition or *adjoint equation*

$$A^*\bar{p} + \mathcal{E}^*\bar{\mu} = -J'_y(\bar{y}, \bar{u}).$$

The operator  $A^* \in \mathcal{L}(H_0^1(\Omega); H^{-1}(\Omega)) \cap \mathcal{L}(L^2(\Omega); Y^*)$  can also be considered as a linear elliptic second-order partial differential operator. (Also, since  $A$  is continuously invertible, so is  $A^*$ .) Thus, the foregoing adjoint equation is also a linear elliptic second-order elliptic PDE for the *adjoint state*  $\bar{p}$  which however involves the measure  $\bar{\mu}$ . One must thus deal with differential equations with measure data in the context of state constraints in optimal control problems. On the other hand, from [Example 3.40](#) we know that  $\bar{\mu}$  is not a totally generic measure; for example, we had seen that  $\bar{\mu} \geq 0$ . Often it is possible to can leverage such additional structural information in PDE analysis.

The presence of the state constraint  $y \leq \psi$  is also the reason why we have introduced the space  $Y = H^2(\Omega) \cap H_0^1(\Omega)$  and thus also why we end up with the above adjoint equation in  $Y^*$ . (This amounts to a *very weak* formulation for the problem at hand). In particular, for the right hand side in the adjoint equation, we only have  $J'_y(\bar{y}, \bar{u}) \in Y^*$  in general. But if  $J$  in fact is already F-differentiable from  $Y_0 \times U$  to  $\mathbb{R}$ , where  $Y$  is densely embedded into  $Y_0$ —for instance, if  $Y_0 = H_0^1(\Omega)$  or  $Y_0 = L^2(\Omega)$ —, then  $J'_y(\bar{y}, \bar{u}) \in Y_0^* \hookrightarrow Y^*$  and additional regularity for the adjoint state  $\bar{p}$  may be derived from the adjoint equation. We have seen this also in [Example 3.39](#) where the adjoint state  $\bar{p}$  was the weak solution to a linear elliptic equation with  $J'_y(\bar{y}, \bar{u}) = \bar{y} - y_d \in L^2(\Omega)$  on the right-hand side. In the present case however, the regularity of  $\bar{p}$  is limited by the presence of the measure  $\bar{\mu}$ , so we will not be able to obtain full regularity in general.



Note that a similar effect can also occur for the optimal control  $\bar{u}$  from the KKT condition  $J'_u(\bar{y}, \bar{u}) = B^*\bar{p}$ , since often  $J'_u(\bar{y}, \bar{u}) = \beta\bar{u}$  for some  $\beta > 0$ . Then

$$J'_u(\bar{y}, \bar{u}) = B^*\bar{p} \iff \bar{u} = \frac{1}{\beta}B^*\bar{p}.$$

See [Example 3.39](#) for the particular example, and also the exercises.

### 3.3 Sufficient optimality conditions

The KKT conditions give a satisfying characterization of (first-order) *necessary* optimality conditions. Of course, we are also interested in *sufficient* conditions. Such conditions often include that the optimal point in question is an isolated optimum, which is an extremely useful property in the numerical analysis of optimal control problems.

We will however see that the situation is quite more delicate than in the finite-dimensional case and that the standard second-order sufficient condition from Nonlinear Optimization will in general not be sufficient in infinite-dimensional settings any more. We begin with the particular and important case of a *convex* problem, where everything still works out just fine.

#### 3.3.1 The convex case

In the convex case, the KKT conditions alone already turn out to be also sufficient for global optimality:

**Theorem 3.43.** *Let  $X$  and  $Z$  be Banach spaces. Suppose that  $f: X \rightarrow \mathbb{R}$  is convex and  $F$ -differentiable,  $G: X \rightarrow Z$  is  $F$ -differentiable and convex w.r.t.  $-K$ , where  $K$  is a closed convex cone. Further, let  $\bar{x} \in \mathcal{F} = G^{-1}[K]$  be a KKT-point of  $(\mathbf{P})$ , so  $\Lambda(\bar{x}) \neq \emptyset$ . Then  $\bar{x}$  is a global solution of  $(\mathbf{P})$ .*

*Proof.* Recall from (3.19) that since  $f: X \rightarrow \mathbb{R}$  and  $G: X \rightarrow Z$  are convex—the latter w.r.t.  $\leq_K$ —, we have the inequality

$$f(x) - f(\bar{x}) \geq \langle f'(\bar{x}), x - \bar{x} \rangle_{X^*, X} \quad \text{for all } x \in X$$

and the inclusion

$$G(\bar{x}) - G(x) + G'(\bar{x})(x - \bar{x}) \in K \quad \text{for all } x \in \mathcal{F}$$

at hand. Now, for  $x \in \mathcal{F}$ , we have  $K + G(x) \subset K$  since  $K$  is a *convex* cone, and hence

$$G'(\bar{x})(x - \bar{x}) \in (K + G(x)) - G(\bar{x}) \subset K - G(\bar{x}) \subset \text{cone}(K, G(\bar{x})) \subset T(K, G(\bar{x}))$$

for all  $x \in \mathcal{F}$ . (Recall [Lemma 3.8](#).) But then optimality of  $\bar{x}$  is immediate from the KKT conditions [\(3.15\)](#) and [\(3.16\)](#) for all  $x \in \mathcal{F}$  as follows:

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

since  $G'(\bar{x})(x - \bar{x}) \in T(K, G(\bar{x}))$  and  $\bar{\lambda} \in T(K, G(\bar{x}))^\circ$  as seen above.  $\square$

**Remark 3.44.** Note that we have only supposed  $\Lambda(\bar{x}) \neq \emptyset$  in [Theorem 3.43](#) instead of  $\bar{x}$  regular, i.e., we have *not* needed a constraint qualification.

### 3.3.2 Second-order sufficient optimality conditions

We next use second-order derivative—so: *curvature*—information to formulate *sufficient optimality conditions*, also in the case of a *nonconvex* problem. Already in high school we learn that if  $f: \mathbb{R} \rightarrow \mathbb{R}$  and we have a point  $\bar{x} \in \mathbb{R}$  with  $f'(\bar{x}) = 0$  and  $f''(\bar{x}) > 0$ , then  $\bar{x}$  is a minimum of  $f$ —this is a second-order sufficient optimality condition! Note that such a curvature condition for  $\bar{x}$  will also imply that  $\bar{x}$  is a *strict* (and thus isolated) local minimum and that the function grows at least quadratically around  $\bar{x}$ , that is, for  $\varepsilon > 0$  sufficiently small, there is a  $\gamma > 0$  such that

$$f(x) \geq f(\bar{x}) + \gamma \|\bar{x} - x\|^2 \quad \text{for all } x \in B_\varepsilon(\bar{x}).$$

Of course, when aiming for sufficient conditions, we can freely assume that the designated optimal solution  $\bar{x} \in \mathcal{F}$  satisfies *necessary* optimality conditions. We will thus consider only KKT-points  $\bar{x}$  in the following.

**Reminder:** The second derivative  $h''(\bar{x})$  of a twice F-differentiable function  $h: X \rightarrow Z$  is given by the F-derivative of the mapping  $h': X \rightarrow \mathcal{L}(X; Z)$ , and is thus a mapping

$$h'': X \rightarrow \mathcal{L}(X; \mathcal{L}(X; Z)) \cong \mathcal{L}^2(X \times X; Z),$$

where  $\mathcal{L}^2(X \times X; Z)$  denotes the space of continuous bilinear forms on  $X \times X$  mapping into  $Z$ . This is compatible with the notion of the *Hessian matrix* for  $X = \mathbb{R}^n$  and  $Z = \mathbb{R}$ , since there is a one-to-one correspondence between bilinear forms on  $\mathbb{R}^n \times \mathbb{R}^n$  and matrices  $\mathbb{R}^{n \times n}$ .

We want to pose a positive definiteness condition on second-order derivatives in  $\bar{x}$ . It is clear that it will not be necessary to require such a condition on the whole space  $X$  but only in certain directions from a tangent cone at the designated optimal point  $\bar{x}$ . (A

function whose second derivatives at every point are positive definite on the whole  $X$  would be convex.)

The intuition here could be as follows: Given a KKT-point  $\bar{x}$ , we have (Lemma 3.30)

$$\langle f'(\bar{x}), d \rangle_{X^*, X} \geq 0 \quad \text{for all } d \in T_\ell(G, K, \bar{x}).$$

If we want to pose additional assumptions in certain directions  $d \in T_\ell(G, K, \bar{x})$ , then the directions for which  $\langle f'(\bar{x}), d \rangle_{X^*, X} > 0$  should not be a problem for optimality since  $f$  “strictly” increases in these directions. It should thus be enough to consider only directions in  $T_\ell(G, K, \bar{x})$  for which the directional derivative of  $f$  is nonpositive.

This leads to the following definition of the *critical cone* as a subset of the linearizing cone:

**Definition 3.45** (Critical cone). The *critical cone* at  $\bar{x} \in \mathcal{F}$  is defined by

$$C(\bar{x}) = \left\{ d \in T_\ell(G, K, \bar{x}) : \langle f'(\bar{x}), d \rangle_{X^*, X} \leq 0 \right\}.$$

In fact, for a KKT-point  $\bar{x}$ , we have (Lemma 3.30)

$$\langle f'(\bar{x}), d \rangle_{X^*, X} \geq 0 \quad \text{for all } d \in T_\ell(G, K, \bar{x}).$$

Accordingly, in this case, the critical cone becomes

$$C(\bar{x}) = \left\{ d \in T_\ell(G, K, \bar{x}) : \langle f'(\bar{x}), d \rangle_{X^*, X} = 0 \right\}. \quad (3.23)$$

Since we will generally suppose that  $\bar{x}$  is a KKT point in the following, we could have defined the critical cone also in the latter form. However, the original definition as in Definition 3.45 will generalize more nicely later when we will have to readjust it a little bit.

Now, regarding which condition we actually pose along  $C(\bar{x})$ , recall that the first part (3.9) of the KKT conditions could be rewritten to  $L'_x(\bar{x}, \bar{\lambda}) = 0$  in  $X^*$  with the Lagrangian function  $L$  as in Definition 3.31. Thus, in order to formulate a sufficient condition for a given KKT point  $\bar{x}$  to be a local solution, so a *minimizer*, it is natural to require that there are no directions of nonpositive curvature w.r.t.  $x$  of the Lagrangian function in  $C(\bar{x})$ ; that is,

$$L''_{xx}(\bar{x}, \bar{\lambda})(d, d) > 0 \quad \text{for all } d \in C(\bar{x}) \setminus \{0\}. \quad (3.24)$$

This is the classical second-order sufficient condition (SOSC) from Nonlinear Optimization. It is not sensible to pose a condition on the curvature of the objective function only; this is already true for finite-dimensional problems, see Figure 2. (The example is taken from [UU12, P. 103].) It is also imperative to compare the classical SOSC (3.24) with the classical *necessary* second-order condition which says that a locally optimal KKT-point  $\bar{x}$  satisfies

$$L''_{xx}(\bar{x}, \bar{\lambda})(d, d) \geq 0 \quad \text{for all } d \in C(\bar{x}).$$

(Unfortunately, the proof of such a necessary second-order condition in infinite dimensions is very hard and only known for particular cases; even the proof for finite-dimensional problems is not very pleasant.) Clearly, the classical SOS (3.24) admits the minimal possible gap between second-order conditions of necessary and sufficient type.

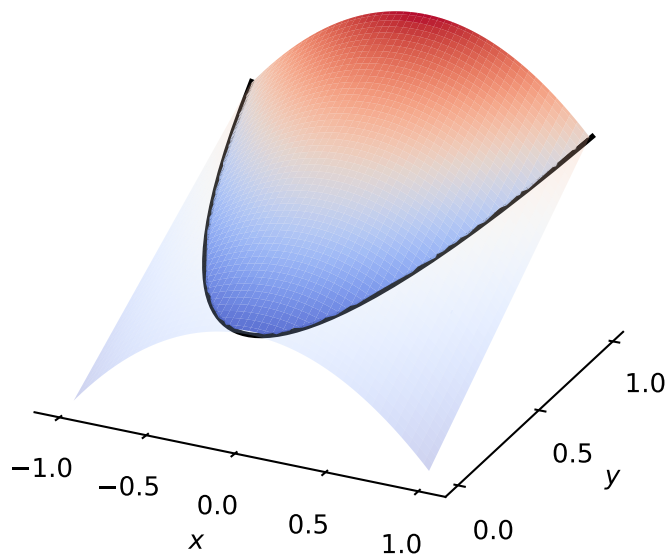


Figure 2: The graph of  $f(x, y) = -x^2 + 2y$ . The black line indicates the boundary of the feasible set  $\mathcal{F}$  defined by  $g(x, y) = x^2 - y \leq 0$ . The minimum of  $f$  over  $\mathcal{F}$  is  $\bar{x} = 0$  which is a KKT-point with multiplier  $\bar{\lambda} = 2$ . The critical cone in  $\bar{x}$  is  $C(\bar{x}) = \mathbb{R} \times \{0\}$ . Clearly, the curvature of  $f$  along  $C(\bar{x})$  is nonpositive, but  $L''_{xx}(\bar{x}, \bar{\lambda})$  is positive definite along  $C(\bar{x}) \setminus \{0\}$  and (3.24) is satisfied.

However, unfortunately, the next example shows that the classical condition (3.24) is in general *not* strong enough for general infinite-dimensional optimization problems to indeed obtain a sufficient optimality condition.

**Example 3.46.** Let  $X = Z = \ell^2$ , where  $\ell^2 = \ell^2(\mathbb{R})$  is the Hilbert space of  $\mathbb{R}$ -valued square summable sequences whose inner product and norm are given by

$$(x, y)_{\ell^2} := \sum_{i=1}^{\infty} x_i y_i, \quad \text{so} \quad \|x\|_{\ell^2} = \left( \sum_{i=1}^{\infty} x_i^2 \right)^{1/2}.$$

Consider the infinite-dimensional optimization problem

$$\min_{x \in \ell^2} (c, x)_{\ell^2} - (x, x)_{\ell^2} \quad \text{s.t.} \quad x_i \geq 0 \text{ for all } i \in \mathbb{N}$$

with the objective function  $f(x) = (c, x)_{\ell^2} - (x, x)_{\ell^2}$ , where  $c \in \ell^2$  satisfies  $c_i > 0$  for all  $i \in \mathbb{N}$ . The constraints are given by  $G(x) \in K$  with  $K = \{x \in \ell^2 : x_i \geq 0 \text{ for all } i \in \mathbb{N}\}$  and  $G(x) = x$ , the identity mapping. (Hence  $\mathcal{F} = K$ .)

Since  $G'(x) = \text{id}_{\ell^2}$  is clearly surjective for all  $x \in \ell^2$ , RCQ holds in every feasible  $x \in K$ , see [Proposition 3.19](#).

We consider  $\bar{x} = 0$  and claim that it is a KKT point with  $-c$  being the unique Lagrange multiplier, so  $\{-c\} = \Lambda(\bar{x}) \neq \emptyset$ . First, concerning the multiplier rule (3.9) in the KKT conditions, observe that  $f'(\bar{x}) \in (\ell^2)^*$  can be identified with  $c - 2\bar{x} \in \ell^2$ . (This is the Fréchet-Riesz representation theorem; in fact, we can regard  $c - 2\bar{x}$  as the gradient  $\nabla f(\bar{x})$ .) Hence, for  $\bar{\lambda}$  to be a Lagrange multiplier for  $\bar{x} = 0$ , we calculate

$$f'(\bar{x}) + G'(\bar{x})^* \bar{\lambda} = c - 2\bar{x} + \bar{\lambda} = c + \bar{\lambda} \stackrel{!}{=} 0 \quad \iff \quad \bar{\lambda} = -c,$$

so  $-c$  is the unique candidate for an element of  $\Lambda(\bar{x})$ . Further,

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

so from  $(c, d)_{\ell^2} \geq 0$  for all  $d \in K$  it follows that  $-c \in K^\circ = T(K, G(\bar{x}))^\circ$ . Hence, in fact  $\{-c\} = \Lambda(\bar{x}) \neq \emptyset$ .

Next we investigate the proposed sufficient condition (3.24). We have  $T_\ell(G, K, \bar{x}) = K$ , hence, the critical cone  $C(\bar{x})$  is given by

$$C(\bar{x}) = \left\{ d \in K : (f'(\bar{x}), d)_{\ell^2} = 0 \right\},$$

but for all  $d \in K \setminus \{0\}$ , we obtain

$$(f'(\bar{x}), d)_{\ell^2} = (c, d)_{\ell^2} = \sum_{i=1}^{\infty} c_i d_i > 0,$$

since  $c_i > 0$  and  $0 \neq d \geq 0$ . Therefore,  $C(\bar{x}) = \{0\}$ , so the classical SOS (3.24) is indeed satisfied. In fact, we have even seen that all directions  $d \in K \setminus \{0\} = T(\mathcal{F}, \bar{x}) \setminus \{0\}$  are strict *ascent directions* (*Aufstiegsrichtungen*) for  $f$ .

Note however that (verify this!)

$$L''_{xx}(\bar{x}, \bar{\lambda})(d, d) = f''(\bar{x})(d, d) = -2\|d\|_{\ell^2}^2 \leq 0 \quad \text{for all } d \in \ell^2.$$

And indeed,  $\bar{x}$  is not a local minimum of  $f$  on  $\mathcal{F} = K$ : Define a sequence  $(x^k) \subset \ell^2$  by  $(x^k)_i := (2\delta_{ik}c_i)_{i \in \mathbb{N}} \subset K$ , so

$$(x^k)_i := \begin{cases} 2c_i & \text{if } i = k, \\ 0 & \text{otherwise.} \end{cases}$$

Then  $\|x^k - \bar{x}\|_{\ell^2} = \|x^k\|_{\ell^2} = 2c_k \rightarrow 0$  as  $k \rightarrow \infty$ , but

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

so  $\bar{x}$  cannot be a local minimum of  $f$  on  $K$ .

February 1, 2023

14 →

Now we are at a crossroads how to obtain a sufficient second-order condition for infinite-dimensional problems. There are essentially two possible ways out:

1. Strengthen the condition (3.24) but keep the general assumptions on (P) as general as possible, or
2. enforce more structure in the problem (P) but keep condition (3.24).

We will present one theorem for each alternative. But first, we state an auxiliary result which will be used in either theorem: Under the Robinson CQ, the linearizing cone  $T_\ell(G, K, \bar{x})$  at  $\bar{x}$  is approximated by feasible directions in the sense that

$$\text{dist}(x - \bar{x}, T_\ell(G, K, \bar{x})) = o(\|x - \bar{x}\|_X) \quad (3.25)$$

for  $\mathcal{F} \ni x \rightarrow \bar{x}$ . This already foreshadows that we will need to suppose that the KKT-point  $\bar{x}$  is in fact regular.

The approximation property (3.25) is implied by the following lemma:

**Lemma 3.47.** *If the Robinson constraint qualification (3.2) holds at  $\bar{x} \in \mathcal{F}$ , then there exists a map  $h: \mathcal{F} \rightarrow T_\ell(G, K, \bar{x})$  with*

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

*Proof.* Let  $x \in \mathcal{F}$  be arbitrary. Then the F-differentiability of  $G$  implies

$$G(x) = G(\bar{x}) + G'(\bar{x})(x - \bar{x}) + r(x), \quad \text{where } \|r(x)\|_Z = o(\|x - \bar{x}\|_X).$$

Then RCQ in the form of the ZKCQ (3.8) shows that there exists  $\delta > 0$  such that

$$\overline{B_{\delta,Z}(0)} \subset G'(\bar{x})\overline{B_X(0)} - ((K - G(\bar{x})) \cap \overline{B_Z(0)}).$$

Hence, we can find  $s(x) \in X$  and  $v(x) \in K - G(\bar{x})$  such that  $r(x) = G'(\bar{x})s(x) - v(x)$  and (by rescaling)

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

Setting  $h(x) := x - \bar{x} + s(x)$ , there holds

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

It remains to show that  $h(x) \in T_\ell(G, K, \bar{x})$ , that is,  $G'(\bar{x})h(x) \in T(K, G(\bar{x}))$ . So:

$$\begin{aligned} G'(\bar{x})h(x) &= G'(\bar{x})(x - \bar{x}) + G'(\bar{x})s(x) \\ &= G'(\bar{x})(x - \bar{x}) + r(x) + v(x) \\ &= G(x) - G(\bar{x}) + v(x) \end{aligned}$$

Writing  $v(x) \in K - G(\bar{x})$  in the form  $v(x) = k(x) - G(\bar{x})$  with  $k(x) \in K$ , we have

$$G'(\bar{x})h(x) = 2 \left( \underbrace{\frac{G(x) + k(x)}{2}}_{\in K} - G(\bar{x}) \right) \in \text{cone}(K, G(\bar{x})).$$

Hence,  $G'(\bar{x})h(x) \in \text{cone}(K, G(\bar{x})) \subset T(K, G(\bar{x}))$  and so  $h(x) \in T_\ell(G, K, \bar{x})$ , and  $h$  has all required properties.  $\square$

**Remark 3.48.** It is also possible to prove the approximation property (3.25) directly via metric regularity, Theorem 3.20, and to construct the function  $h$  as in Lemma 3.47 from there. See the exercises.

We next prove the theorems as announced above.

**Stronger SOSC:** In the foregoing Example 3.46, the condition (3.24) is not sufficient for optimality of a KKT point because the critical cone  $C(\bar{x})$  is too small. Thus, part of a solution to obtain a sufficient condition would be to enlarge the critical cone  $C(\bar{x})$ , so make (3.24) stronger.

**Definition 3.49** (Approximate critical cone). For  $\eta \geq 0$  we define the  $\eta$ -approximate critical cone at  $\bar{x} \in \mathcal{F}$  by

$$C_\eta(\bar{x}) = \left\{ d \in T_\ell(G, K, \bar{x}) : \langle f'(\bar{x}), d \rangle_{X^*, X} \leq \eta \|d\|_X \right\}.$$

Note that  $C_0(\bar{x}) = C(\bar{x})$  and  $C_\eta(\bar{x}) = T_\ell(G, K, \bar{x})$  for all  $\eta \geq \|f'(\bar{x})\|_{X^*}$ .

We now prove a very general theorem about second-order sufficient conditions which needs very little structural properties of (P). The proof rests on Taylor expansion for the Lagrange function and Lemma 3.47. For convenience, we will use the shorthand notation  $d^2$  for directions  $(d, d)$  in second derivatives, so for example  $f''(\bar{x})d^2$  instead of  $f''(\bar{x})(d, d)$ .

**Theorem 3.50** (Second-order sufficient conditions). *Let  $X$  and  $Z$  be Banach spaces with  $K \subset Z$  closed and convex. Further, let  $f: X \rightarrow \mathbb{R}$  and  $G: X \rightarrow Z$  be twice  $F$ -differentiable. Assume that  $\bar{x} \in \mathcal{F} = G^{-1}[K]$  satisfies the RCQ (3.2) and the*

*KKT-conditions with multiplier  $\bar{\lambda}$ :*

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

*Let in addition the following second-order condition hold:*

$$L''_{xx}(\bar{x}, \bar{\lambda})(d, d) \geq \gamma \|d\|_X^2 \quad \text{for all } d \in C_\eta(\bar{x}), \quad (3.26)$$

*where  $\gamma > 0$  and  $\eta > 0$  are fixed constants. Then  $\bar{x}$  is an isolated local solution of (P) and there exist  $\kappa, \delta > 0$  such that the quadratic growth condition*

$$f(x) \geq f(\bar{x}) + \kappa \|x - \bar{x}\|_X^2 \quad \text{for all } x \in \mathcal{F} \cap B_{\delta, X}(\bar{x})$$

*holds true.*

*Proof.* Consider  $x \in \mathcal{F} \cap B_{\delta, X}(\bar{x})$  for some  $\delta > 0$  sufficiently small to be chosen later, and set  $d(x) := x - \bar{x}$ . Note that we do not know whether  $d(x) \in T_\ell(K, G, \bar{x})$ , so (3.26) cannot be used for  $d(x)$  directly. However, Lemma 3.47 braces us with the approximation  $d(x) = h(x) + r(x)$  where  $h(x) \in T_\ell(K, G, \bar{x})$  and  $\|r(x)\|_X = o(\|d(x)\|_X)$ . The plan is to use  $h(x)$  as a surrogate for  $d(x)$  in  $T_\ell(K, G, \bar{x})$ . To this end, we first note some approximation properties of  $h(x)$ , namely (verify those!)

$$\|h(x)\|_X = \|d(x)\|_X + o(\|d(x)\|_X). \quad (3.27)$$

and the quadratic equivalent

$$\|h(x)\|_X^2 = \|d(x)\|_X^2 + o(\|d(x)\|_X^2). \quad (3.28)$$

Now the case where  $\langle f'(\bar{x}), h(x) \rangle_{X^*, X} > \eta \|h(x)\|_X$  is quite easy: Using  $\langle f'(\bar{x}), r(x) \rangle_{X^*, X} = o(\|d(x)\|_X)$ , we find

$$\begin{aligned} f(x) - f(\bar{x}) &= \langle f'(\bar{x}), d(x) \rangle_{X^*, X} + o(\|d(x)\|_X) \\ &= \langle f'(\bar{x}), h(x) \rangle_{X^*, X} + o(\|d(x)\|_X) \\ &> \eta \|h(x)\|_X + o(\|d(x)\|_X) \\ &\stackrel{(3.27)}{=} \eta \|d(x)\|_X + o(\|d(x)\|_X) \geq \eta \|d(x)\|_X^2 = \eta \|x - \bar{x}\|_X^2, \end{aligned}$$

where the last inequality holds true for  $\delta$  sufficiently small.

Next we deal with the case  $\langle f'(\bar{x}), h(x) \rangle_{X^*, X} \leq \eta \|h(x)\|_X$ , so  $h(x) \in C_\eta(\bar{x})$ . From the KKT conditions, we have  $\bar{\lambda} \in T(K, G(\bar{x}))^\circ$ , so

$$L(x, \bar{\lambda}) - L(\bar{x}, \bar{\lambda}) = f(x) - f(\bar{x}) + \langle \bar{\lambda}, G(x) - G(\bar{x}) \rangle_{Z^*, Z} \leq f(x) - f(\bar{x}), \quad (3.29)$$



since  $G(x) - G(\bar{x}) \in T(K, G(\bar{x}))$ . Moreover, by (3.27) and the construction of  $r(x)$ , we have

$$L''_{xx}(\bar{x}, \bar{\lambda})(r(x), h(x)) + L''_{xx}(\bar{x}, \bar{\lambda})(h(x), r(x)) + L''_{xx}(\bar{x}, \bar{\lambda})r(x)^2 = o(\|d(x)\|_X^2). \quad (3.30)$$

Recalling that the first KKT condition means in fact  $L'_x(\bar{x}, \bar{\lambda}) = 0$  in  $X^*$ , we thus obtain by Taylor expansion

$$\begin{aligned} f(x) - f(\bar{x}) &\stackrel{(3.29)}{\geq} L(x, \bar{\lambda}) - L(\bar{x}, \bar{\lambda}) \\ &= \langle L'_x(\bar{x}, \bar{\lambda}), d(x) \rangle_{X^*, X} + \frac{1}{2} L''_{xx}(\bar{x}, \bar{\lambda})(d(x), d(x)) + o(\|d(x)\|_X^2) \\ &= \frac{1}{2} L''_{xx}(\bar{x}, \bar{\lambda})(h(x) + r(x), h(x) + r(x)) + o(\|d(x)\|_X^2) \\ &\stackrel{(3.30)}{=} \frac{1}{2} L''_{xx}(\bar{x}, \bar{\lambda})h(x)^2 + o(\|d(x)\|_X^2) \\ &\geq \frac{\gamma}{2} \|h(x)\|_X^2 + o(\|d(x)\|_X^2) \\ &= \frac{\gamma}{2} \|d(x)\|_X^2 + o(\|d(x)\|_X^2) \geq \frac{\gamma}{4} \|d(x)\|_X^2 = \frac{\gamma}{4} \|x - \bar{x}\|_X^2, \end{aligned}$$

the last inequality again for  $\delta$  sufficiently small.

The assertion then follows with  $\kappa := \min(\frac{\gamma}{4}, \eta)$ . □

**Remark 3.51.** One can dispose of the assumption that  $\bar{x}$  satisfies the RCQ if the feasible set is “flat” enough around  $\bar{x}$  in the sense that there is a  $\delta > 0$  such that for all  $x \in B_{\delta, X}(\bar{x})$  we have  $x - \bar{x} \in T_\ell(G, K, \bar{x})$ . Indeed, we have only used RCQ to invoke Lemma 3.47 which provides the substitute  $h(x) \in T_\ell(G, K, \bar{x})$  for  $x - \bar{x}$ , because for the latter we do not know in general that it is an element of the linearizing cone. For example, box constraints are “flat”.

**More structural assumptions for (P):** Theorem 3.50 works in a very general setting, which is very nice, but makes it hard to actually verify the conditions in practical situations. On the other hand, we have posed no further assumptions on  $f$  and  $G$  or even the geometry of the Banach space  $X$  at all, so there might be room to strengthen assumptions there and weaken condition (3.26). Indeed, it is possible to derive second-order sufficient conditions in the classical form (3.24) under more specific assumptions on  $f$  and  $G$ , if  $X$  is reflexive. These conditions are met by many problems subject to PDEs or, generally, by optimal control problems.

**Theorem 3.52.** *Let  $X$  and  $Z$  be Banach spaces with  $K \subset Z$  closed and convex and  $X$  reflexive. Further, let  $f: X \rightarrow \mathbb{R}$  be given in the form  $f = f_1 + f_2$  with  $f_1, f_2: X \rightarrow \mathbb{R}$ , and let  $f_1, f_2$  and  $G: X \rightarrow Z$  be twice  $F$ -differentiable. Assume that*

$\bar{x} \in \mathcal{F} = G^{-1}[K]$  satisfies the RCQ (3.2) and the KKT-conditions with multiplier  $\bar{\lambda}$ :

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

Further, suppose that the second derivatives exhibit the following weak continuity properties: If  $d_k \rightharpoonup d$  in  $X$ , then, for  $i = 1, 2$ ,

$$f_i''(\bar{x})d^2 \leq \liminf_{k \rightarrow \infty} f_i''(\bar{x})d_k^2 \quad \text{and} \quad G''(\bar{x})d_k^2 \rightharpoonup G''(\bar{x})d \quad \text{in } Z.$$

Let in addition  $d \mapsto f_2''(\bar{x})d^2$  be coercive on the linearizing cone, that is, there exists a constant  $\alpha > 0$  such that

$$f_2''(\bar{x})d^2 \geq \alpha \|d\|_X^2 \quad \text{for all } d \in T_\ell(G, K, \bar{x}),$$

and suppose that the classical SOSOC is satisfied:

$$L''_{xx}(\bar{x}, \bar{\lambda})d^2 > 0 \quad \text{for all } d \in C(\bar{x}) \setminus \{0\}. \quad (3.24)$$

Then  $\bar{x}$  is an isolated local solution of (P) and there exist  $\kappa, \delta > 0$  such that the quadratic growth condition

$$f(x) \geq f(\bar{x}) + \kappa \|x - \bar{x}\|_X^2 \quad \text{for all } x \in \mathcal{F} \cap B_{\delta, X}(\bar{x})$$

holds true.

*Proof.* Let  $\bar{x}$  be as in the theorem and let  $\bar{\lambda} \in \Lambda(\bar{x})$  its associated Lagrange multiplier. We argue by contradiction, so suppose that the quadratic growth condition is wrong. Then there exists a sequence  $(x_k) \subseteq \mathcal{F}$  such that  $x_k \rightarrow \bar{x}$  and

$$\frac{1}{k} \|x_k - \bar{x}\|_X^2 \geq f(x_k) - f(\bar{x}).$$

Set

$$d_k := \frac{x_k - \bar{x}}{\|x_k - \bar{x}\|_X}.$$

Since  $X$  is assumed to be reflexive,  $(d_k)$  admits a weakly convergent subsequence, which we do not relabel, with the weak limit  $d_k \rightharpoonup d \in X$ . We show that  $d \in C(\bar{x})$  and use the SOSOC to derive that  $d = 0$ . This will give a contradiction with the coercivity of  $f_2''(\bar{x})$ .

Let  $h: \mathcal{F} \rightarrow T_\ell(G, K, \bar{x})$  be the map established in Lemma 3.47 (here we use the assumption that  $\bar{x}$  is regular) and set  $h_k := h(x_k)$ . Then

$$\langle x', h_k \rangle_{X^*, X} = \langle x', x_k - \bar{x} \rangle_{X^*, X} + o(\|x_k - \bar{x}\|_X),$$

for any  $x' \in X^*$ , so

$$g_k := \frac{h_k}{\|x_k - \bar{x}\|_X} \rightharpoonup d.$$

The linearizing cone  $T_\ell(G, K, \bar{x})$  is closed and convex (Remark 3.12 and Theorem 3.22), hence weakly closed, and we find  $d \in T_\ell(G, K, \bar{x})$  (Proposition 2.3). Moreover, from the above contradiction assumption,

$$\begin{aligned} \frac{1}{k} \|x_k - \bar{x}\|_X^2 &\geq f(x_k) - f(\bar{x}) = \langle f'(\bar{x}), x_k - \bar{x} \rangle_{X^*, X} + o(\|x_k - \bar{x}\|_X) \\ &= \langle f'(\bar{x}), h_k \rangle_{X^*, X} + o(\|x_k - \bar{x}\|_X), \end{aligned}$$

so by dividing by  $\|x_k - \bar{x}\|_X$  and letting  $k \rightarrow \infty$ , we obtain that  $\langle f'(\bar{x}), d \rangle_{X^*, X} \leq 0$  and  $d \in C(\bar{x})$ .

Recall that  $\bar{\lambda} \in T(K, G(\bar{x}))^\circ$ , so with  $G(x_k) - G(\bar{x}) \in \text{cone}(K, G(\bar{x})) \subseteq T(K, G(\bar{x}))$  (Lemma 3.8), there holds

$$\langle \bar{\lambda}, G(x_k) - G(\bar{x}) \rangle_{Z^*, Z} \leq 0.$$

In particular, again from the contradiction assumption,

$$\begin{aligned} \frac{1}{k} \|x_k - \bar{x}\|_X^2 &\geq f(x_k) - f(\bar{x}) \geq f(x_k) - f(\bar{x}) + \langle \bar{\lambda}, G(x_k) - G(\bar{x}) \rangle_{Z^*, Z} \\ &= L(x_k, \bar{\lambda}) - L(\bar{x}, \bar{\lambda}). \end{aligned}$$

Of course we will use Taylor expansion for the latter term to get higher derivatives involved:

$$L(x_k, \bar{\lambda}) - L(\bar{x}, \bar{\lambda}) = \langle L'_x(\bar{x}, \bar{\lambda}), x_k - \bar{x} \rangle_{X^*, X} + L''_{xx}(\bar{x}, \bar{\lambda})(x_k - \bar{x})^2 + o(\|x_k - \bar{x}\|_X^2).$$

The first derivative vanishes since  $\bar{x}$  is a KKT point and  $\bar{\lambda}$  is the associated multiplier. In the second derivative, we substitute with  $h_k$  instead of  $x_k - \bar{x}$  at the cost of an error of  $o(\|x_k - \bar{x}\|_X^2)$ :

$$L''_{xx}(\bar{x}, \bar{\lambda})(x_k - \bar{x})^2 = L''_{xx}(\bar{x}, \bar{\lambda})h_k^2 + o(\|x_k - \bar{x}\|_X^2)$$

to altogether obtain

$$\frac{1}{k} \geq L''_{xx}(\bar{x}, \bar{\lambda})g_k^2 + \frac{o(\|x_k - \bar{x}\|_X^2)}{\|x_k - \bar{x}\|_X^2}.$$

From the assumptions, the map  $h \mapsto L''_{xx}(\bar{x}, \bar{\lambda})h^2 = f''(\bar{x})h^2 + \langle \bar{\lambda}, G''(\bar{x})h^2 \rangle$  is weakly lower semicontinuous, thus (recall  $g_k \rightharpoonup d$ )

$$0 \geq \liminf_{k \rightarrow \infty} \left[ L''_{xx}(\bar{x}, \bar{\lambda})g_k^2 + \frac{o(\|x_k - \bar{x}\|_X^2)}{\|x_k - \bar{x}\|_X^2} \right] = \liminf_{k \rightarrow \infty} L''_{xx}(\bar{x}, \bar{\lambda})g_k^2 \geq L''_{xx}(\bar{x}, \bar{\lambda})d^2.$$

But then the SOSC (3.24) implies that  $d = 0$ .

Now finally, let  $\alpha > 0$  be the coercivity parameter for  $f''_2(\bar{x})$ . Note that also  $h \mapsto L''_{xx}(\bar{x}, \bar{\lambda})h^2 - f''_2(\bar{x})h^2$  is weakly lower semicontinuous. Thus, with  $g_k \rightharpoonup d = 0$ ,

$$0 = L''_{xx}(\bar{x}, \bar{\lambda})d^2 - f''_2(\bar{x})d^2 \leq \liminf_{k \rightarrow \infty} \left[ L''_{xx}(\bar{x}, \bar{\lambda})g_k^2 - f''_2(\bar{x})g_k^2 \right]$$

and so

$$\begin{aligned} \alpha &\leq \alpha \liminf_{k \rightarrow \infty} \|d_k\|_X^2 = \alpha \liminf_{k \rightarrow \infty} \|g_k\|_X^2 \leq \liminf_{k \rightarrow \infty} f_2''(\bar{x})g_k^2 \\ &\leq \liminf_{k \rightarrow \infty} \left[ L_{xx}''(\bar{x}, \bar{\lambda})g_k^2 - f_2''(\bar{x})g_k^2 \right] + \liminf_{k \rightarrow \infty} f_2''(\bar{x})g_k^2 \\ &\leq \liminf_{k \rightarrow \infty} L_{xx}''(\bar{x}, \bar{\lambda})g_k^2 = 0. \end{aligned}$$

But this is a contradiction to  $\alpha > 0$ . □

[Remark 3.51](#) about disposing of the RCQ assumption for  $\bar{x}$  applies also to the foregoing [Theorem 3.52](#).

February 8, 2023

15 →

**Remark 3.53.** The assumptions on the coercive part  $f_2$  in the objective function in [Theorem 3.52](#) typically follow from a “control cost” term in the optimal control formulation which usually takes the form of a (squared and scaled) norm. For example, in the running example [Example 2.9](#), we had  $\frac{\alpha}{2}\|u\|_{L^2(\partial\Omega)}^2$ . (Verifying the coercivity assumption in the case of  $x = (y, u)$  still requires some reasoning and identification of the linearizing cone, though.) The assumptions on  $f_2$  and  $G$  will usually require to employ compactness methods. See [Example 3.55](#) below.

**Remark 3.54.** Due to the particular structural assumptions in [Theorem 3.52](#), in this setting the classical SOSC ([3.24](#)) is in fact equivalent to the uniform SOSC which requires that there exists  $\gamma > 0$  such that

$$L_{xx}''(\bar{x}, \bar{\lambda})d^2 \geq \gamma\|d\|_X^2 \quad \text{for all } d \in C(\bar{x}). \quad (3.31)$$

Of course, [\(3.31\)](#) always implies [\(3.24\)](#), but the reverse is in general *false* in infinite-dimensional spaces. (In finite-dimensional spaces,  $\gamma$  in [\(3.31\)](#) is related to the finite set of eigenvalues of  $L_{xx}''(\bar{x}, \bar{\lambda})$  of which there must exist a smallest one; in infinite dimensions that need not be true at all.) But here the classical SOSC ([3.24](#)) is in fact sufficient for [\(3.31\)](#). Indeed, set

$$\gamma := \inf_{\substack{\|d\|_X=1 \\ d \in C(\bar{x})}} L_{xx}''(\bar{x}, \bar{\lambda})d^2$$

and consider an infimal sequence  $(d_k)$  with  $d_k \in C(\bar{x})$  and  $\|d_k\|_X = 1$  such that  $\gamma = \lim_{k \rightarrow \infty} L_{xx}''(\bar{x}, \bar{\lambda})d_k^2$ . Then, after passing to a subsequence,  $d_k \rightharpoonup d$  in  $X$ . We show that  $\gamma > 0$ .

1. If  $d = 0$ , then at the end of the proof of [Theorem 3.52](#)) we have seen that since  $d_k \rightharpoonup d = 0$ ,

$$\alpha \leq \lim_{k \rightarrow \infty} L_{xx}''(\bar{x}, \bar{\lambda})d_k^2.$$

But the limit is exactly  $\gamma$ , by construction, so  $\gamma \geq \alpha > 0$ .

2. If  $d \neq 0$ , then already from weak lower semicontinuity of  $L''_{xx}(\bar{x}, \bar{\lambda})$  we have

$$0 < L''_{xx}(\bar{x}, \bar{\lambda})d^2 \leq \liminf_{k \rightarrow \infty} L''_{xx}(\bar{x}, \bar{\lambda})d_k^2 = \gamma.$$

**Example 3.55.** Consider the setting of [Example 2.9](#) and [Example 3.39](#) for the semilinear optimal control problem. We verify the assumptions on second derivatives in [Theorem 3.52](#). In the example we had

$$G(\bar{x}) = \begin{pmatrix} E(\bar{y}, \bar{u}) \\ \bar{u} \end{pmatrix} \quad \text{and} \quad G'(\bar{x})d = \begin{pmatrix} E'_y(\bar{y}, \bar{u})z + E'_u(\bar{y}, \bar{u})h \\ h \end{pmatrix}$$

where  $d = (z, v) \in Y \times U = H^1(\Omega) \times L^2(\partial\Omega)$  and

$$E'_y(\bar{y}, \bar{u})z: \left[ v \mapsto (\nabla z, \nabla v)_{L^2(\Omega)^n} + (z, v)_{L^2(\partial\Omega)} + 3(\bar{y}^2 z, v)_{L^2(\Omega)} \right] \in H^1(\Omega)^*,$$

and  $E'_u(\bar{y}, \bar{u})h = -(\cdot, h)_{L^2(\partial\Omega)}$ . Note that neither derivative depends on  $\bar{u}$  any more, since the control was entering the original problem *linearly*. Thus, when taking another derivative with respect to  $x = (y, u)$  in direction  $d = (z, h)$ , all derivatives in  $u$  direction vanish and we obtain

$$G''(\bar{x})d^2 = \begin{pmatrix} E''_{yy}(\bar{y}, \bar{u})z^2 \\ 0 \end{pmatrix}$$

with

$$E''_{yy}(\bar{y}, \bar{u})z^2: \left[ v \mapsto 6(\bar{y}z^2, v)_{L^2(\Omega)} \right] \in H^1(\Omega)^*.$$

We show that  $z \mapsto E''_{yy}(\bar{y}, \bar{u})z^2$  is weakly continuous as a mapping  $H^1(\Omega) \rightarrow H^1(\Omega)^*$ . Note that for every  $v \in H^1(\Omega) \hookrightarrow L^6(\Omega)$  (because  $n \leq 3$ ), we find that  $\bar{y}v \in L^3(\Omega)$  (Hölder's inequality!), thus by duality

$$\left[ w \mapsto \int_{\Omega} \bar{y}wv \, dx \right] \quad \text{is an element of} \quad \mathcal{L}(L^{3/2}(\Omega), \mathbb{R}).$$

Now, let  $(z_k) \subseteq H^1(\Omega)$  be a sequence with  $z_k \rightharpoonup z$  in  $H^1(\Omega)$ . Then, since the embedding  $H^1(\Omega) \hookrightarrow L^3(\Omega)$  is compact,  $z_k \rightarrow z$  in  $L^3(\Omega)$  and  $z_k^2 \rightarrow z^2$  in  $L^{3/2}(\Omega)$ . (The latter follows from continuity of the  $t \mapsto |t|^{3/2}$  superposition operator, recall the exercises.) Hence, from the above,

$$\langle E''_{yy}(\bar{y}, \bar{u})z_k^2, v \rangle = \int_{\Omega} \bar{y}z_k^2 v \, dx \quad \longrightarrow \quad \int_{\Omega} \bar{y}z^2 v \, dx = \langle E''_{yy}(\bar{y}, \bar{u})z^2, v \rangle$$

for every  $v \in H^1(\Omega)$  and this is the desired weak continuity. It follows that  $d \mapsto G''(\bar{x})d^2$  is weakly continuous as in the assumptions of [Theorem 3.52](#).

We next turn to the second derivatives of  $f$ . Set

$$f_1(x) = \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 \quad \text{and} \quad f_2(x) = \frac{\alpha}{2} \|u\|_{L^2(\partial\Omega)}^2.$$

The first-order derivatives of  $f(x) = J(y, u)$  were already calculated in [Example 3.39](#). From there, we find, again with  $d = (z, h)$ ,

$$f_1''(\bar{x})d^2 = J''_{yy}(\bar{y}, \bar{u})z^2 = (z, z)_{L^2(\Omega)}^2 = \|z\|_{L^2(\Omega)}^2$$

and

$$f_2''(\bar{x})d^2 = J''_{uu}(\bar{y}, \bar{u})h^2 = \alpha(h, h)_{L^2(\partial\Omega)}^2 = \alpha\|h\|_{L^2(\partial\Omega)}^2.$$

There is no  $z$  component in  $f_2''(\bar{x})d^2$ , so we have to recover the coercivity assumption on  $f_2''(\bar{x})$  from the linearizing cone. One might be tempted to conjecture that setting  $f_1(\bar{x}) = 0$  and  $f_2(\bar{x}) = f(\bar{x})$  would have been more easy since then a squared norm of  $z$  would have occurred in  $f_2''(\bar{x})$ . But  $\|z\|_{L^2(\Omega)}^2$  is *not* coercive with respect to  $Y = H^1(\Omega)$ , so there is nothing gained. In fact, we can make do with  $f_2$  as defined.

Indeed, suppose that  $d \in T_\ell(G, K, \bar{x})$ . In [Example 3.38](#), we had seen that with  $K = \{0\}_{H^1(\Omega)^*} \times U_{\text{ad}}$ , there was  $T(K, G(\bar{x})) = \{0\}_{H^1(\Omega)^*} \times T(U_{\text{ad}}, \bar{u})$  and  $T(U_{\text{ad}}, \bar{u}) = \overline{\text{cone}(U_{\text{ad}}, \bar{u})}$ . By definition, the condition that  $d = (z, h) \in T_\ell(G, K, \bar{x})$  means that  $G'(\bar{x})d \in T(K, G(\bar{x}))$ , so with the above formula for  $G'(\bar{x})d$ , we obtain the conditions that

$$h \in \overline{\text{cone}(U_{\text{ad}}, \bar{u})} \quad \text{and} \quad E'_y(\bar{y}, \bar{u})z = -E'_u(\bar{y}, \bar{u})h.$$

The latter means (recall [Example 3.4](#)) that  $z$  is the weak solution to

$$\left. \begin{aligned} -\Delta z + 3\bar{y}^2 z &= 0 && \text{on } \Omega, \\ \frac{\partial z}{\partial \nu} + z &= h && \text{on } \partial\Omega. \end{aligned} \right\} \quad (3.32)$$

In [Example 3.39](#) it was already mentioned that this problem has a unique solution  $z \in H^1(\Omega)$  for every  $h \in L^2(\partial\Omega) \hookrightarrow H^1(\Omega)^*$ , and in fact,  $E'_y(\bar{y}, \bar{u})$  is surjective  $H^1(\Omega) \rightarrow H^1(\Omega)^*$ . It is thus in particular continuously invertible—see the open mapping theorem reminder from the beginning of [Section 3.2](#)—, so if  $d = (z, h) \in T_\ell(G, K, \bar{x})$ , then there exists a constant  $C > 0$  such that

$$\|z\|_{H^1(\Omega)} \leq C\|h\|_{H^1(\Omega)^*} \leq C\|h\|_{L^2(\partial\Omega)}.$$

But then, going back to the second derivative of  $f_2$ , we find

$$f_2''(\bar{x})d^2 = \alpha \|h\|_{L^2(\partial\Omega)}^2 \geq \frac{\alpha}{2} \|h\|_{L^2(\partial\Omega)}^2 + \frac{\alpha}{2C^2} \|z\|_{H^1(\Omega)}^2 \geq \bar{\alpha} \|d\|_X^2,$$

for some  $\bar{\alpha} > 0$ , since  $X = H^1(\Omega) \times L^2(\partial\Omega)$ . Thus  $d \mapsto f_2''(\bar{x})d^2$  is indeed coercive on  $T_\ell(G, K, \bar{x})$ .

Now finally it remains to observe that if  $d_k = (z_k, h_k) \rightharpoonup (z, h)$  in  $H^1(\Omega) \times L^2(\partial\Omega)$ , then  $z_k \rightarrow z$  in  $L^2(\Omega)$  due to compact embedding, hence

$$f_1(\bar{x})d_k^2 = \|z_k\|_{L^2(\Omega)}^2 \longrightarrow \|z\|_{L^2(\Omega)}^2 = f_1(\bar{x})d^2$$

and

$$\liminf f_2(\bar{x})d_k^2 = \alpha \|h_k\|_{L^2(\partial\Omega)}^2 \geq \alpha \|h\|_{L^2(\partial\Omega)}^2$$

since norms are the prime example of weakly lower semicontinuous functions.

Now assume that  $(\bar{y}, \bar{u})$  is a KKT-point. Then, as derived in [Example 3.39](#),

$$\bar{u}(x) = \text{proj}_{[a(x), b(x)]} \left( \frac{1}{\alpha} \bar{p}(x) \right) = \begin{cases} a(x) & \text{if } \frac{1}{\alpha} \bar{p}(x) \leq a(x), \\ \frac{1}{\alpha} \bar{p}(x) & \text{if } a(x) < \frac{1}{\alpha} \bar{p}(x) < b(x), \\ b(x) & \text{if } \frac{1}{\alpha} \bar{p}(x) \geq b(x). \end{cases}$$

with the unique weak solution  $\bar{p} \in H^1(\Omega)$  to

$$\begin{aligned} -\Delta \bar{p} + 3\bar{y}^2 \bar{p} &= y_d - \bar{y} && \text{on } \Omega, \\ \frac{\partial \bar{p}}{\partial \nu} + \bar{p} &= 0 && \text{on } \partial\Omega. \end{aligned}$$

We calculate the critical cone  $C(\bar{x}) = C(\bar{y}, \bar{u})$ . It was already derived that  $d = (z, h)$  is in the linearizing cone  $T_\ell(G, K, \bar{x})$  if and only if  $h \in \overline{\text{cone}(U_{\text{ad}}, \bar{u})}$ , that is,  $h|_{[\bar{u}=a]} \geq 0$  and  $h|_{[\bar{u}=b]} \leq 0$ , and  $z \in H^1(\Omega)$  is the weak solution to [\(3.32\)](#). Further,

$$f'(\bar{x})d = J_y'(\bar{y}, \bar{u})z + J_u(\bar{y}, \bar{u})h = (\bar{y} - y_d, z)_{L^2(\Omega)} + \alpha(\bar{u}, h)_{L^2(\partial\Omega)} \stackrel{!}{=} 0.$$

But, since  $\bar{p}$  is the weak solution to the problem with  $y_d - \bar{y}$  on the right-hand side as above, and  $z$  is the weak solution with boundary data  $h$ ,

$$(\bar{y} - y_d, z)_{L^2(\Omega)} = -(\nabla \bar{p}, \nabla z)_{L^2(\Omega)} - 3(\bar{y}^2 \bar{p}, z)_{L^2(\Omega)} = -(\bar{p}, h)_{L^2(\partial\Omega)}.$$

So  $f'(\bar{x})d = 0$  for  $d = (z, h) \in T_\ell(G, K, \bar{x})$  precisely when

$$(\bar{p}, h)_{L^2(\partial\Omega)} = \alpha(\bar{u}, h)_{L^2(\partial\Omega)}.$$

Since  $\bar{u}$  was given by the pointwise projection of  $\bar{p}/\alpha$  onto  $[a, b]$ , we have  $\bar{p} = \alpha\bar{u}$  on  $[a \leq \bar{p}/\alpha \leq b]$ . Due to the sign conditions of  $h$  on  $[\bar{u} = a] \supset [\bar{p}/\alpha < a]$  and  $[\bar{u} = b] \supset [\bar{p}/\alpha > b]$ , we obtain the condition that

$$h|_{[\bar{p}/\alpha < a]} = h|_{[\bar{p}/\alpha > b]} = 0.$$

So,  $(z, h) \in C(\bar{x})$  precisely when  $z$  is the unique solution to (3.32) with data  $h \in L^2(\partial\Omega)$  satisfying

$$h|_{[\bar{u}=a]} \geq 0, \quad h|_{[\bar{u}=b]} \leq 0 \quad \text{and} \quad h|_{[\bar{p}/\alpha < a]} = h|_{[\bar{p}/\alpha > b]} = 0.$$

The resulting SOSOC is then given by the requirement that

$$\int_{\Omega} (1 + 6\bar{y}\bar{p})z^2 \, dx + \alpha \int_{\partial\Omega} h^2 \, d\omega > 0$$

for all pairs  $(z, h) \neq 0$  in the critical cone as before.



---

## References

- [Br10] H. Brezis: Functional Analysis, Sobolev Spaces and Partial Differential Equations. Springer, New York, 2010.
- [BS00] J.F. Bonnans, A. Shapiro: Perturbation Analysis of Optimization Problems. Springer Science + Business Media, 2000.
- [GT01] D. Gilbarg and N.S. Trudinger: Elliptic Partial Differential Equations of Second Order. Springer Science + Business Media, 2001.
- [HPUU09] M. Hinze, R. Pinnau, M. Ulbrich, S. Ulbrich: Optimization with PDE Constraints. Springer Netherlands, 2009.
- [La93] S. Lang: Real and Functional Analysis. Springer, New York, 1993.
- [NW06] J. Nocedal, S. Wright: Numerical Optimization. Springer, New York, 2006.
- [Ro72] S.M Robinson: Normed convex processes, *Trans. Amer. Math. Soc.* 174, 127–140, 1972.
- [Ro76] S.M Robinson: Stability theory for systems of inequalities in nonlinear programming, part II: differentiable nonlinear systems. *SIAM J. Num. Anal.* 13, 497–513, 1976.
- [Ro76b] S.M Robinson: Regularity and stability for convex multivalued functions, *Math. Oper. Res.* 1, 130–143, 1976.
- [Sc07] W. Schirotzek: Nonsmooth analysis. Springer, Berlin Heidelberg, 2007.
- [Tr10] F. Tröltzsch: Optimal Control of Partial Differential Equations: Theory, Methods, Applications. American Mathematical Society, 2010.
- [UU12] M. Ulbrich, S. Ulbrich: Nichtlineare Optimierung. Birkhäuser, Basel, 2012.
- [ZK79] J. Zowe, S. Kurcyusz: Regularity and stability for the mathematical programming problem in Banach spaces. *Appl. Math. Optimization* 5, 49–62, 1979.