# Optimal control and directional differentiability for elliptic quasi-variational inequalities

Amal Alphonse\* Michael Hintermüller<sup>†</sup> Carlos N. Rautenberg<sup>‡</sup>
March 28, 2022

#### Abstract

We focus on elliptic quasi-variational inequalities (QVIs) of obstacle type and prove a number of results on the existence of solutions, directional differentiability and optimal control of such QVIs. We give three existence theorems based on an order approach, an iteration scheme and a sequential regularisation through partial differential equations. We show that the solution map taking the source term into the set of solutions of the QVI is directionally differentiable for general data and locally Hadamard differentiable obstacle mappings, thereby extending in particular the results of our previous work which provided the first differentiability result for QVIs in infinite dimensions. Optimal control problems with QVI constraints are also considered and we derive various forms of stationarity conditions for control problems, thus supplying among the first such results in this area.

#### **Contents**

I	Intr	oduction	2
	1.1	Contributions of the paper	3
	1.2	Basic assumptions and notations	
2	Exis	tence for QVIs	4
	2.1	Iteration scheme	4
	2.2	Birkhoff–Tartar order approach	
	2.3	Sequential regularisation by PDEs	
3	Dire	ectional differentiability 1	0
	3.1	Iteration scheme and expansion formulae	2
	3.2	Passage to the limit	
	3.3	Continuity properties of the directional derivative	
	3.4	Complementarity characterisation of the directional derivative	
	3.5	Examples of QVIs with multiple solutions	
4	Exis	tence of optimal controls	7
	4.1	The penalised optimal control problem	8
5	Stat	ionarity 1	-
	5.1	Bouligand stationarity	(
	5.2	Weak C-stationarity	
		5.2.1 Stationarity for the penalised optimal control problem	2
		5.2.2 Passage to the limit $\rho \to 0$	
	5.3	E-almost C-stationarity	
	5.4	From $\mathcal{E}$ -almost to C-stationarity	
	5.5	Strong stationarity	

The authors extend their gratitude to the two referees for their careful reading and excellent comments which helped to greatly improve some of the results and presentation. AA and MH were partially supported by the DFG through the DFG SPP 1962 Priority Programme *Non-smooth and Complementarity-based Distributed Parameter Systems: Simulation and Hierarchical Optimization* within project 10. MH and CNR acknowledge the support of Germany's Excellence Strategy - The Berlin Mathematics Research Center MATH+ (EXC-2046/1, project ID: 390685689) within project AA4-3. In addition, MH acknowledges the support of SFB-TRR154 within subproject B02, and CNR was supported by NSF grant DMS-2012391.

<sup>\*</sup>Weierstrass Institute, Mohrenstrasse 39, 10117 Berlin, Germany (alphonse@wias-berlin.de)

<sup>†</sup>Weierstrass Institute, Mohrenstrasse 39, 10117 Berlin, Germany (hintermueller@wias-berlin.de)

<sup>†</sup>Department of Mathematical Sciences and the Center for Mathematics and Artificial Intelligence (CMAI), George Mason University, Fairfax, VA 22030, USA (crautenb@gmu.edu)

#### 31

## 1 Introduction

Quasi-variational inequalities (QVIs) are generalisations of variational inequalities (VIs) where the constraint set in which the solution is sought depends on the unknown solution itself. The very nature of the dependency of the constraint set on the solution intrinsically leads to a complicated and challenging mathematical structure since it significantly amplifies the nonlinear and nonsmooth nature of VIs. Another attribute that fundamentally distinguishes QVIs from VIs is the lack of uniqueness of solutions (in general) which then necessitates the consideration of multi-valued or set-valued solution mappings. QVIs arise in a multitude of models describing phenomena in fields such as biology, physics, economics and social sciences amongst others. First introduced by Bensoussan and Lions [17, 48] in the study of stochastic impulse controls, specific applications involving QVIs are thermoforming processes [4, 7], the formation and growth of lakes, rivers and sandpiles [59, 15, 58, 56, 16], games in the context of generalised Nash equilibrium problems [34, 25, 55], and magnetisation of superconductors [44, 14, 57, 62]. See [5, 12] for additional details and references.

In this paper, we focus on elliptic QVIs of obstacle type or compliant obstacle problems. These have the form

find 
$$y \in \mathbf{K}(y) : \langle Ay - f, y - v \rangle \le 0 \quad \forall v \in \mathbf{K}(y) \text{ where } \mathbf{K}(y) := \{ v \in V : v \le \Phi(y) \}.$$
 (1)

Here  $f \in V^*$  is data,  $\Phi \colon V \to V$  is a given obstacle map, and V is a reflexive Banach space possessing an ordering  $\leq$  which is used in the definition of the constraint set (we shall be more precise below). Let us define  $\mathbf{Q}$  to be the solution map associated to the QVI in (1) so that it reads  $y \in \mathbf{Q}(f)$ . We develop in this paper theory addressing the matters of existence for (1), directional differentiability of  $\mathbf{Q}$  and stationarity conditions for optimal control problems with QVI constraints of the form

$$\min_{\substack{u \in U_{ad} \\ y \in \mathbf{Q}(u)}} \frac{1}{2} \|y - y_d\|_H^2 + \frac{\nu}{2} \|u\|_U^2.$$
(2)

Different methodologies exist for the mathematical treatment of existence for QVIs. There is an approach based on order that was pioneered by Tartar [67] which relies on the existence of subsolutions and supersolutions to guarantee existence of solutions (typically, one takes 0 as a subsolution which would hold under sign conditions on the source term). In certain cases, the QVI can be expressed as a generalized equation and it therefore belongs to a more general problem class [40, 41, 26, 39, 27]. In problems involving constraints on derivatives (which is not the case under consideration in this paper), special forms of regularisation of the constraint that modify the partial differential operator may be suitable, see [62, 52, 10, 11]. For more details, we refer the reader to the survey paper [5]. We discuss in §2 appropriate conditions on the function spaces and the obstacle map  $\Phi$  for  $\mathbf{Q}(f)$  to be non-empty. One approach relies on an iteration argument where a contraction-type property of  $\Phi$  is used. Another existence result is given for source terms bounded from below by using the aforementioned Birkhoff–Tartar theory, and we also study a sequential regularisation approach of the QVI by PDEs where the QVI constraint is handled by a penalty term.

Literature on the differentiability and sensitivity analysis for solution maps associated to QVIs in infinite dimensions is almost non-existent: our contributions [4, 6] appear to be the first ones that address these issues. In [4], we give a first directional differentiability result for the solution map taking the source term into the set of solutions for non-negative sources and directions whilst in [6] we studied continuity properties related to minimal and maximal solution mappings of QVIs. In §3, we derive directional differentiability results for Q. We extend and improve here our previous work [4] which provided differentiability results for source and direction terms that are non-negative; in this paper we shall remove this restriction in our Theorem 3.2, which requires minimal (and locally formulated) assumptions to apply. We give a characterisations of the QVI that is satisfied by the directional derivative of Q as a complementarity system and in §3.3 we also prove a continuity result that shows that the derivative depends continuously on the direction under some assumptions. This gives a comprehensive answer to the question of sensitivity analysis of QVIs under rather general conditions.

The scarcity of work done on the optimal control of QVIs in infinite dimensions is unsurprisingly even more pronounced; see [2, 6, 23, 24, 54] for some of the very few contributions. In our work [6], in addition to stability properties we also provided results on the optimal control of minimal and maximal solutions of QVIs. While this article was under preparation, we note that [72] has appeared wherein the author considers elliptic QVIs and their differential sensitivity and strong stationarity conditions for the optimal control problem but for Frèchet differentiable obstacle maps  $\Phi$ ; we assume only Hadamard differentiability of  $\Phi$  for the differentiability result and we furthermore provide other forms of stationarity as well as existence/approximation results. For QVIs in the finite dimensional setting, see [53] and the references therein. In sharp contrast, control problems with VI constraints have attracted wide attention: see for example [13, 51, 19, 18, 38, 37, 36, 65, 43, 33, 69] and the references therein. We shall consider in §4 the optimal control problem (2) where existence of the optimal control will be shown using a standard calculus of variations argument. Then we turn our attention to the derivation of stationarity conditions for the optimal control and state. There are a number of concepts

of stationarity for these types of control problems, see [37] for a discussion. We first work on obtaining Bouligand stationarity in §5.1, then a form of weak C-stationarity in §5.2, moving on to  $\mathcal{E}$ -almost C-stationarity conditions in §5.3 by approximating the QVI control-to-state map through PDEs (as done in §2.3) and then passing to the limit. We discuss in §5.4 how to upgrade to C-stationarity from  $\mathcal{E}$ -almost C-stationarity and finally, in §5.5, we provide a strong stationarity result.

## 1.1 Contributions of the paper

We summarise the main results of this work.

#### · QVI:

- Theorems 2.18 and 2.19: existence for (1) via a penalty approach,
- Theorem 3.2: directional differentiability for QVIs for locally Hadamard maps  $\Phi$  under local Lipschitz conditions,
- Proposition 2.1: complementarity characterisations of the QVI in (1),
- Proposition 3.12: continuity properties of the QVI satisfied by directional derivative,
- Proposition 3.13: complementarity characterisation of the QVI satisfied by the directional derivative of the solution map.

#### • Optimal control:

- Theorem 4.1: existence of optimal controls for (2).

#### • Stationarity conditions for (2):

- Proposition 5.2: Bouligand stationarity,
- Theorem 5.5: weak C-stationarity,
- Theorem 5.11:  $\mathcal{E}$ -almost C-stationarity,
- Proposition 5.15: C-stationarity,
- Theorem 5.16: strong stationarity.

#### 1.2 Basic assumptions and notations

We make some standing assumptions that are necessary throughout the paper, except where mentioned otherwise.

We always work with real Banach or Hilbert spaces. Let V be a Banach space and denote the standard duality pairing on  $V^* \times V$  by  $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_{V^*, V}$ . Take  $A \colon V \to V^*$  to be a linear operator that satisfies the following properties for all  $u, v \in V$ :

$$\begin{split} \langle Au,v\rangle &\leq C_b \left\|u\right\|_V \left\|v\right\|_V, \\ \langle Au,u\rangle &\geq C_a \left\|u\right\|_V^2, \\ \langle Au^+,u^-\rangle &\leq 0, \end{split} \tag{boundedness}$$

where  $C_a, C_b > 0$  are constants. We will frequently suppose that the Banach space V is a vector lattice for a partial ordering  $\leq$ . This means that for all  $u, v \in V$ , the following holds:

- (i)  $u \le u$  (reflexivity),
- (ii)  $u \le v$  and  $v \le u$  implies u = v (anti-symmetry),
- (iii)  $u \le v$  and  $v \le w$  implies  $u \le w$  (transitivity),
- (iv)  $u \le v$  implies that  $u + w \le v + w$  and  $\lambda u \le \lambda v$  for  $\lambda \ge 0$ ,
- (v) there exists a greatest lower bound  $\inf(u, v)$  and a least upper bound  $\sup(u, v)$  belonging to V.

See for example [3, 49] or [61, §4:5] for more details. It should be emphasised that in the context of function spaces over a bounded Lipschitz domain  $\Omega$ , with  $\leq$  chosen as the usual a.e. ordering, (v) allows for  $V = L^p(\Omega)$  and  $V = W^{1,p}(\Omega)$  for  $1 \leq p < \infty$  but not  $V = W^{2,p}(\Omega)$  in general. We write the positive cone of V as

$$V_{+} := \{ v \in V : v \ge 0 \}$$

(this is convex but not necessarily closed). If V is a Banach lattice, the projection onto  $V_+$  (assuming this is well defined) of an element  $v \in V$  agrees with  $\sup(0,v)$ , but this is not necessarily the case for a general vector lattice. Note that the dual space  $V^*$  inherits an ordering: we say  $f \leq g$  in  $V^*$  if and only if  $\langle g - f, v \rangle \geq 0$  for all  $v \in V_+$ .

Regarding the obstacle map, we take  $\Phi \colon V \to V$  to be given.

The identity operator will be denoted by I. We denote continuous, dense, and compact embeddings of spaces by  $\hookrightarrow$ ,  $\stackrel{d}{\hookrightarrow}$ , and  $\stackrel{c}{\hookrightarrow}$  respectively. The notation  $B_R(u)$  will be used to mean the closed ball in V of radius R centred at u.

# 2 Existence for QVIs

We begin by discussing three existence results for the QVI in (1), reproduced here:

$$y \le \Phi(y) : \langle Ay - f, y - v \rangle \le 0 \quad \forall v \in V : v \le \Phi(y),$$

involving different approaches. We start by obtaining existence through iteration by solutions of VIs. Then we consider a translation of the theory by Birkhoff–Tartar for source terms that are bounded from below and we finish by considering a sequential regularisation approach through PDEs. These existence results entail different assumptions. The third approach is useful for purposes of numerical realisation. The second approach requires  $\Phi$  to be increasing and bounded below in a certain sense.

Before we proceed, let us give the following characterisation involving (1).

**Proposition 2.1.** The QVI in (1) is equivalent to the complementarity system

$$\xi := f - Ay, \tag{3a}$$

$$\xi \ge 0,$$
 (3b)

$$\langle \xi, \Phi(y) - y \rangle = 0,$$
 (3c)

$$0 \le \Phi(y) - y. \tag{3d}$$

*Proof.* The proof is standard. By definition,  $\xi$  satisfies  $\langle \xi, y - v \rangle \geq 0$  for all feasible v. Setting  $v = \Phi(y)$  and then  $v = 2y - \Phi(y)$ , we obtain the orthogonality condition (3c) for  $\xi$ . Testing with  $v = y - \varphi$  for  $\varphi \in V$  with  $\varphi \geq 0$  gives the stated non-negativity. The reverse direction follows from writing  $\langle Ay - f, y - v \rangle = \langle Ay - f, y - \Phi(y) \rangle + \langle Ay - f, \Phi(y) - v \rangle$  (where v is a feasible test function) and using the second and third lines in the system.

#### 2.1 Iteration scheme

We need the following assumption for this section (as an example,  $V = L^2(\Omega)$  or  $H^1(\Omega)$  on a bounded Lipschitz domain are valid).

**Assumption 2.2.** Let V be a Hilbert space and a vector lattice with  $V_+$  closed and suppose that  $\Phi \colon V \to V$  is increasing.

The lattice and increasing properties are necessary to apply the comparison principle for VIs [61, §4:5]. This assumption also implies the following useful property (whose proof is in Appendix A), which can be thought of as a weak monotone convergence theorem (in fact, it suffices for V to be a reflexive Banach space rather than Hilbert for the result).

**Lemma 2.3.** If  $\{v_n\} \subset V$  is a bounded sequence which is either increasing or decreasing (i.e., either  $v_n \leq v_{n+1}$  for all n, or  $v_n \geq v_{n+1}$  for all n), then there exists a  $v \in V$  such that  $v_n \rightharpoonup v$  in V (for the full sequence).

Let  $S \colon V^* \times V \to V$  be the solution mapping of the VI associated to the class of QVIs under consideration, i.e.  $y = S(f, \psi)$  solves

$$y \le \Phi(\psi) : \langle Ay - f, y - v \rangle \le 0 \quad \forall v \in V : v \le \Phi(\psi).$$

Take a source term  $f \in V^*$  and set

$$y_0 := A^{-1}f$$

to be the solution of the unconstrained problem. The function  $y_1 := S(f, y_0)$  satisfies  $y_1 \le y_0$  by the comparison principle [61, §4:5, Theorem 5.1], and defining

$$y_n := S(f, y_{n-1}),$$

we see that  $y_n \leq y_{n-1}$  by repeated applications of the comparison principle. Hence  $\{y_n\}$  is monotonically decreasing and each  $y_n$  satisfies

$$y_n \in V, y_n \le \Phi(y_{n-1}) : \langle Ay_n - f, y_n - v \rangle \le 0 \quad \forall v \in V : v \le \Phi(y_{n-1}). \tag{4}$$

We look for a uniform bound on  $\{y_n\}$ . When the obstacle map is such that it always dominates some given function  $v_0 \in V$ , this is easy since we may test with  $v = v_0$ . Otherwise, we need the following.

#### **Lemma 2.4.** *If*

$$\|\Phi(v)\|_{V} \le C_X \|v\|_{V} \quad \forall v \in V \text{ where } C_X < \frac{C_a}{C_b},\tag{5}$$

then  $\{y_n\}$  is bounded in V.

<sup>&</sup>lt;sup>1</sup>Heuristically,  $y_0$  is considered as a solution of the VI with source f and obstacle equal to  $\infty$ .

*Proof.* Since  $y_n \leq y_{n-1}$  and  $\Phi$  is increasing,  $\Phi(y_n) \leq \Phi(y_{n-1})$  and so  $\Phi(y_n)$  is a valid test function in (4) and we obtain

$$C_{a} \|y_{n}\|_{V}^{2} \leq \langle Ay_{n}, \Phi(y_{n}) \rangle + \langle f, y_{n} - \Phi(y_{n}) \rangle$$

$$\leq C_{b} \|y_{n}\|_{V} \|\Phi(y_{n})\|_{V} + \|f\|_{V^{*}} \|y_{n} - \Phi(y_{n})\|_{V}$$

$$\leq C_{b}C_{X} \|y_{n}\|_{V}^{2} + (1 + C_{X}) \|f\|_{V^{*}} \|y_{n}\|_{V}.$$

From this, we deduce that  $y_n$  is bounded in V under the condition on  $C_X$  in (5).

Now we pass to the limit and show that  $\mathbf{Q} \colon V^* \rightrightarrows V$  is such that  $\mathbf{Q}(f) \neq \emptyset$  under certain circumstances.

**Theorem 2.5.** Let Assumption 2.2 hold and suppose that

either there exists 
$$v_0 \in V$$
 such that  $v_0 \leq \Phi(v)$  for all  $v \in V$ , or (5),

if 
$$\{v_n\} \subset V$$
 is decreasing with  $v_n \rightharpoonup v$  in  $V$  and  $v \leq \Phi(v_n)$ , then  $v \leq \Phi(v)$ . (7)

For any  $f \in V^*$ , there exists a solution  $y \in \mathbf{Q}(f) \cap (-\infty, A^{-1}f]$  which is the weak limit of the sequence  $\{y_n\}$  defined above.

*Proof.* We obtain, thanks to monotonicity and Lemma 2.3 that  $y_n \rightharpoonup y$  in V (for the full sequence) for some y. Since  $\{y_n\}$  is decreasing,  $y_m - y_n \in V_+$  where  $n \geq m$  for m fixed. As  $V_+$  is closed and convex, it is weakly sequentially closed, giving  $y_m \geq y$ . This implies that for arbitrary  $v^* \in V$  with  $v^* \leq \Phi(y)$ , we have  $v^* \leq \Phi(y_m)$ . We take such a  $v^*$  as the test function in the VI for  $y_m$  and then pass to the limit to obtain that y satisfies the inequality in (1) and it remains to be seen that  $y \leq \Phi(y)$ . This follows from passing to the limit in  $y \leq y_m \leq \Phi(y_{m-1})$  by making use of (7).

The assumption (7) is rather weak and it is satisfied if, for example,  $\Phi \colon V \to V$  is weakly sequentially continuous.

**Remark 2.6.** For QVIs with more general or different types of constraints one might need to assume Mosco convergence (see [61, §4:4]) properties of the underlying constraint sets.

**Example 2.7.** The prototypical example for  $\Phi$  to have in mind is a map given by the inverse of a partial differential operator such as

$$\Phi(w) := L^{-1}w + f_0,$$

for example with  $L\colon V\to V^*$  a second-order linear elliptic operator on a bounded Lipschitz domain  $\Omega$  and  $f_0\in V$ . The validity of elliptic regularity and continuous dependence estimates for L would give compactness properties for  $\Phi$  (and weak maximum principles would also yield the increasing property for  $\Phi$ ). See [4, §1.2] for more details on this and on an application to fluid flow.

#### 2.2 Birkhoff-Tartar order approach

In this section, we extend Birkhoff–Tartar-type existence results typically used for QVIs with non-negative source terms to QVIs with source terms that are allowed to be negative. This leads to different assumptions than those made in §2.1. The bedrock of this technique, as detailed in the introduction, is the result of Tartar [67] that gives existence of fixed points for increasing maps that possess subsolutions and supersolutions, see also [9, Chapter 15, §15.2]. We need the following functional setup in this section.

**Assumption 2.8.** Let  $V \stackrel{d}{\hookrightarrow} H$  be a continuous and dense embedding of Hilbert spaces and let  $C \subset H$  be a closed convex cone satisfying

$$C = \{ h \in H : (h, g)_H \ge 0 \text{ for all } g \in C \}.$$
 (8)

This induces an ordering defined by

$$h_1 \leq h_2$$
 if and only if  $h_2 - h_1 \in H_+$ .

Note that  $H_+ \equiv C$ . We write  $h^+ = P_{H_+}h$  to denote the orthogonal projection of  $h \in H$  onto  $H_+$  and define  $h^- := h^+ - h$ . We assume that  $v \in V$  implies  $v^+ \in V$  and that there exists a C > 0 with  $||v^+||_V \leq C||v||_V$  for all  $v \in V$ . Finally, suppose that

$$\Phi \colon H \to V$$
 is increasing.

Note that H is a vector lattice (induced by C) and that the ordering induces an ordering for V in the obvious way and also an ordering for  $V^*$  as elucidated in §1.2. Also,  $-h^- \in P_{-H_+}h$  because C satisfies (8).

Let us recall the Birkhoff–Tartar result (see [9, §15.2, Proposition 2]) for increasing maps under the assumptions on the function spaces in Assumption 2.8.

**Theorem 2.9** (Birkhoff–Tartar). Suppose that  $T: H \to H$  is an increasing map and let  $\underline{h}$  be a subsolution and  $\overline{h}$  be a supersolution of the map T, i.e.,

$$\underline{h} \le T(\underline{h})$$
 and  $T(\overline{h}) \le \overline{h}$ .

If  $h \leq \overline{h}$ , then the set of fixed points of T in the interval  $[h, \overline{h}]$  is non-empty and has a minimal and a maximal element.

With this at hand, we can study existence for (1).

**Theorem 2.10.** Let Assumption 2.8 hold and suppose that

there exists 
$$v_0 \in V$$
 such that  $v_0 \le \Phi(v_0)$ . (9)

Given  $f \in V^*$  with  $Av_0 \leq f \leq F$  for some  $F \in V^*$ , there exist solutions  $y \in \mathbf{Q}(f) \cap [v_0, A^{-1}F]$ . Furthermore, there exists a minimal and a maximal solution on this interval.

*Proof.* By the comparison principle,  $S(f,v_0) \geq S(Av_0,v_0) = v_0$ , hence  $v_0$  is a subsolution for  $S(f,\cdot)$ . Since  $\Phi$  is increasing,  $A^{-1}F = S(F,\Phi^{-1}(\infty))^2 \geq S(F,A^{-1}F) \geq S(f,A^{-1}F)$  so that  $A^{-1}F$  is a supersolution. We also have  $v_0 = S(Av_0,v_0) \leq S(F,v_0) \leq S(F,\Phi^{-1}(\infty)) = A^{-1}F$ , i.e., the subsolution lies below the supersolution. Finally,  $S(f,\cdot)$  is increasing due to  $\Phi$  being increasing. The result follows from the Birkhoff–Tartar theorem.

A typical situation in examples is when  $\Phi(0) \geq 0$  and  $f \in V_+^*$ . While the assumption (9) of the existence of such a  $v_0$  may appear to be restrictive, note that choosing  $v_0 \equiv 0$  recovers the setting of [4] which has been successfully applied to an application in thermoforming. The next example illustrates the existence of such a function  $v_0$  to a map  $\Phi$  related to solution maps of elliptic PDEs.

**Example 2.11.** Let  $\Omega \subset \mathbb{R}^n$  be a bounded Lipschitz domain and set  $H := L^2(\Omega)$ . Suppose  $(V, H, V^*)$  is a Gelfand triple<sup>3</sup> with V a reflexive Banach space. Given a linear, bounded, coercive and T-monotone operator  $B : V \to V^*$  and a source term  $g \in V^*$ , let  $\Phi(u) = \varphi$  be defined<sup>4</sup> as the solution of

$$B\varphi = g + u$$
.

Take any  $v_0 \in V$ . We claim that if g is such that

$$g \geq Bv_0 - v_0$$

in  $V^*$ , then (9) is satisfied. To see this, set  $v := \Phi(v_0)$  so that  $Bv = g + v_0$ . Adding the same term to both sides, we obtain  $B(v - v_0) = g + v_0 - Bv_0$ . Test this with the function  $(v - v_0)^-$  to obtain

$$\langle B(v-v_0)^-, (v-v_0)^- \rangle < -\langle q+v_0-Bv_0, (v-v_0)^- \rangle < 0.$$

## 2.3 Sequential regularisation by PDEs

In this section, we obtain existence results for (1) by regularising the QVI by PDEs by a penalty approach similar to [47, §3.5.2, p. 370]. There has been considerable effort on various aspects and methods of regularisation of VIs by PDEs; see for example [30, §3.2] for an approach similar to what we consider here and [46] and [42, §IV] for a penalisation involving approximations to the Heaviside graph (see also [61, §5:3] on this).

**Assumption 2.12.** Let V be a reflexive Banach space and a vector lattice such that  $V_+$  is closed.

Recall that an operator  $T: X \to X^*$  is hemicontinuous [63, Definition 2.3] if  $s \mapsto \langle T(x+sy), z \rangle_{X^*,X}$  is continuous for all  $x, y, z \in X$ . For each  $\rho > 0$ , let  $m_{\rho}: V \to V^*$  be a hemicontinuous map such that

$$m_{\rho}(v) = 0 \text{ if } v \le 0 \tag{10}$$

$$\langle m_{\rho}(u) - m_{\rho}(v), u - v \rangle \ge 0 \tag{11}$$

$$z_{\rho} \rightharpoonup z \text{ in } V \text{ and } m_{\rho}(z_{\rho}) \to 0 \text{ in } V^* \text{ (as } \rho \to 0) \implies z \le 0$$
 (12)

**Remark 2.13.** The last condition precludes the possibility of having 'bad' choices of  $m_{\rho}$  such as  $\rho(\cdot)^+$ . It is also worth pointing out that if  $m_{\rho} \equiv m$  for some map m, then (12) implies that

$$m(z) = 0 \implies z \le 0,$$

which is the converse of (10), so (12) can be thought of a strengthening of the classical kernel or penalty condition that one finds in penalty approaches for VIs.

 $<sup>^2</sup>$ By  $S(F,\Phi^{-1}(\infty))$  we simply mean the solution of the unconstrained problem with source F. No invertibility of  $\Phi$  is necessary.

<sup>&</sup>lt;sup>3</sup>Recall that  $V \subset H \equiv H^* \subset V^*$  is called a *Gelfand triple* if V is a reflexive Banach space continuously and densely embedded into the Hilbert space H and H has been identified with its dual through the Riesz map.

<sup>&</sup>lt;sup>4</sup>The interest in such obstacle mappings is not merely academic, see [4] for some applications.

It is always possible to find such a sequence of maps  $\{m_{\rho}\}$ , see the next example as well as Example 2.17 for the Gelfand triple case. For this reason, we will not usually explicitly refer to (10)–(12) in statements of theorems.

**Example 2.14** (Existence of  $m_{\rho}$ ). Let V and  $V^*$  be strictly convex<sup>5</sup>. Indeed, with  $\mathcal{J}: V \to V^*$  denoting the duality mapping<sup>6</sup> the choice

$$m_{\rho}(u) := \mathcal{J}(u - P_{V_{-}}(u))$$

furnishes such an example where  $P_{V_-}: V \to V_-$  is the metric projection onto the set of non-positive elements  $V_-$ . Properties (10) and (11) as well as hemicontinuity follow as in [47, §3.5.2, Theorem 5.1, p. 370]. In fact, note that  $m_{\rho}(u) = 0$  implies that  $u \leq 0$  (because  $\mathcal{J}$  is bijective and passes through the origin [73, Proposition 32.22 (a), (b)]). For (12), denoting  $m_{\rho} \equiv m$ , by monotonicity, we have for every  $\lambda > 0$  and  $v \in V$  that

$$\langle m(z_{\rho}) - m(z + \lambda v), z_{\rho} - z - \lambda v \rangle \ge 0$$

whence passing to the limit  $\rho \to 0$ , using  $m(z_{\rho}) \to 0$  in  $V^*$  (by hypothesis),  $\langle m(z + \lambda v), \lambda v \rangle \geq 0$  and then dividing through by  $\lambda$  and sending  $\lambda \to 0$ , by hemicontinuity, we obtain that m(z) = 0 in  $V^*$  and thus  $z \leq 0$ .

We consider the penalisation<sup>8</sup>

$$Ay_{\rho} + \frac{1}{\rho}m_{\rho}(y_{\rho} - \Phi(y_{\rho})) = f \tag{13}$$

of (1) and study the convergence properties of its solution as  $\rho \to 0$ . First, we discuss existence. We recall that a map  $T\colon X\to X^*$  is said to be *radially continuous* [63, Definition 2.3] if  $s\mapsto \langle T(x+sy),y\rangle_{X^*,X}$  is continuous for all  $x,y\in X$ , and a map  $R\colon X\to Y$  between Banach spaces is said to be *completely continuous* [66, §2] if  $x_n\rightharpoonup x$  in X implies that  $R(x_n)\to R(x)$  in Y.

Proposition 2.15 (Existence for the penalised equation). Under Assumption 2.12, assume

there exists 
$$v_0 \in V$$
 such that  $v_0 \leq \Phi(v)$  for all  $v \in V$  (14)

and one of the following:

$$m_{\rho}(I - \Phi) \colon V \to V^*$$
 is completely continuous, (15a)

$$m_o(I - \Phi) \colon V \to V^*$$
 is monotone, radially continuous and bounded. (15b)

Given  $f \in V^*$ , there exists a solution  $y_{\rho} \in V$  of (13). Furthermore, every solution satisfies

$$||y_o||_V \le C (||f||_{V^*} + ||v_0||_V),$$

where C is independent of  $\rho$ .

*Proof.* We have that  $A + (1/\rho)m_{\rho}(I - \Phi)$  is a bounded operator (under (15a), recall that completely continuous maps are bounded). Let us show that it is also coercive. First, by adding and subtracting the same term, observe the formula

$$\langle m_{\rho}(y_{\rho} - \Phi(y_{\rho})), y_{\rho} - v_{0} \rangle = \langle m_{\rho}(y_{\rho} - \Phi(y_{\rho})) - m_{\rho}(v_{0} - \Phi(y_{\rho})), y_{\rho} - v_{0} \rangle$$

$$> 0$$

(by monotonicity (11) and because  $m_{\rho} \equiv 0$  on  $(-\infty, 0]$  from (10)). Now, using this, we have

$$\langle Ay_{\rho}, y_{\rho} - v_{0} \rangle + \langle m_{\rho}(y_{\rho} - \Phi(y_{\rho})), y_{\rho} - v_{0} \rangle \ge C_{a} \|y_{\rho}\|_{V}^{2} - C_{b} \|y_{\rho}\|_{V} \|v_{0}\|_{V},$$

which yields coercivity of the full elliptic operator.

Suppose that (15a) is available. By [66, §2, Lemma 2.1], A is a type M operator. Since the sum of a type M operator and a completely continuous operator is type M [66, §2, Example 2.B], we get that the full elliptic operator is of type M. Then [66, §2, Corollary 2.2] yields existence. Under (15b), the full elliptic operator is pseudomonotone by [63, Lemma 2.9 and Lemma 2.11] giving existence via [63, Theorem 2.6].

 $<sup>^5</sup>$ All Hilbert spaces (and thus their duals) are strictly convex. In fact, the strict convexity requirement in the assumption is no issue in the setting of reflexive Banach spaces: by Asplund's theorem (see e.g., [47, §2.2.2, Theorem 2.5]), V can be renormed via an equivalent norm making V and  $V^*$  strictly convex.

 $<sup>^6</sup>$ The assumption of strict convexity gives appropriate properties of  $\mathcal{J}$  (such as single-valuedness), see [47, §2.2.2, p. 174] and [73, §32.3d] for more details.

<sup>&</sup>lt;sup>7</sup>This is well defined since we assumed  $V_+$  (and hence  $V_-$ ) is closed and because V is a reflexive and strictly convex space.

<sup>&</sup>lt;sup>8</sup>For the results of this section, it would be sufficient to simply consider the case where each  $m_{\rho} \equiv m$ , but in anticipation of the optimal control problem that we shall later study (in particular when we derive optimality conditions), it becomes useful to consider this generality now.

Regarding the estimate on the solution, we test the equation with  $y_{\rho} - v_0$  and use the above coercivity estimate to find

$$\begin{aligned} C_a \left\| y_\rho \right\|_V^2 &\leq C_b \left\| y_\rho \right\|_V \left\| v_0 \right\|_V + \left\| f \right\|_{V^*} \left\| y_\rho \right\|_V + \left\| f \right\|_{V^*} \left\| v_0 \right\|_V \\ &\leq \frac{C_a}{3} \left\| y_\rho \right\|_V^2 + \frac{3C_b^2}{4C_a} \left\| v_0 \right\|_V^2 + \frac{3}{4C_a} \left\| f \right\|_{V^*}^2 + \frac{C_a}{3} \left\| y_\rho \right\|_V^2 + \frac{1}{2} \left\| f \right\|_{V^*}^2 + \frac{1}{2} \left\| v_0 \right\|_V^2. \end{aligned}$$

This gives the uniform bound

$$\frac{C_a}{3} \|y_\rho\|_V^2 \le \left(\frac{3C_b^2}{4C_a} + \frac{1}{2}\right) \|v_0\|_V^2 + \left(\frac{3}{4C_a} + \frac{1}{2}\right) \|f\|_{V^*}^2.$$

**Remark 2.16.** The assumptions of the previous lemma are by no means necessary. One could, for example, ask for  $(I - \Phi) : V \to V$  to be invertible and  $A(I - \Phi)^{-1} : V \to V^*$  to be pseudomonotone and coercive instead of (15a) or (15b) and then apply [63, Theorem 2.6] to obtain the same result.

Let us point out a very common setting.

**Example 2.17** (Gelfand triple case). Suppose that

$$V \subset H \equiv H^* \subset V^*$$
 is a Gelfand triple with  $V \stackrel{c}{\hookrightarrow} H$  and  $H$  is a vector lattice defined via (8) with  $H_+$  closed, (16)  $\Phi \colon V \to H$  is completely continuous.

Set  $h^+ = P_{H_+}h$  to be the orthogonal projection in H. We assume that  $(\cdot)^+ : V \to V$ . We can take  $m_\rho : V \to H^* \equiv H$  defined by

$$m_{\rho}(v) := (v^+, \cdot)_H$$

and this satisfies (10), (11), (12) and (15a). Indeed, (12) follows because  $P_{H_+}: H \to H_+$  is Lipschitz continuous and the compact embedding and complete continuity imply (15a) (using the fact that the projection operator is continuous in H).

We write the possibly multivalued solution mapping associated to the equation under study as  $\mathbf{P}_{\rho} \colon V^* \rightrightarrows V$ , so (13) reads  $y_{\rho} \in \mathbf{P}_{\rho}(f)$ . Now, thanks to the lemma, for every source term  $f_{\rho} \in V^*$ , the following equation has a solution  $y_{\rho}$ :

$$Ay_{\rho} + \frac{1}{\rho}m_{\rho}(y_{\rho} - \Phi(y_{\rho})) = f_{\rho}.$$
 (18)

The next two theorems show that solutions of the regularised problem (13) converge to solutions of the QVI under varying assumptions.

**Theorem 2.18** (Existence and approximation of solutions to the QVI). Let Assumption 2.12, (14), either (15a) or (15b) and

$$\Phi \colon V \to V \text{ is completely continuous}$$
 (19)

hold. Take a sequence  $f_{\rho} \to f$  in  $V^*$ . Then there exists a subsequence  $\{\rho_n\}_n$  and elements  $y_{\rho_n} \in \mathbf{P}_{\rho_n}(f_{\rho_n})$  such that  $y_{\rho_n} \to y$  in V where  $y \in \mathbf{Q}(f)$ .

*Proof.* The proof is in four steps and is similar to the proof of Theorem 2.3 of [36].

1. Uniform estimates and feasibility of limit. For each  $\rho$ , let  $y_{\rho}$  be a solution of (18) (such a selection is possible due to the axiom of choice). By Proposition 2.15, it satisfies the estimate

$$||y_{\rho}||_{V} \le C \left(||f_{\rho}||_{V^*} + ||v_{0}||_{V}\right),$$

and this is bounded, hence for a subsequence (which we do not attempt to differentiate for ease of reading),  $y_{\rho} \rightharpoonup y$  in V to some y. Rearranging the equality (18),

$$||m_{\rho}(y_{\rho} - \Phi(y_{\rho}))||_{V^*} = \rho ||f_{\rho} - Ay_{\rho}||_{V^*} \le C\rho$$

and therefore  $m_{\rho}(y_{\rho} - \Phi(y_{\rho})) \to 0$  in  $V^*$  as  $\rho \to 0$ . Then (12) implies that  $y \le \Phi(y)$ .

2. Monotonicity formula. For  $v \in V$ , we get by adding and subtracting the same term and using the monotonicity of  $m_{\rho}$ ,

$$\langle m_{\rho}(y_{\rho} - \Phi(y_{\rho})), y_{\rho} - v \rangle = \langle m_{\rho}(y_{\rho} - \Phi(y_{\rho})) - m_{\rho}(v - \Phi(y_{\rho})), y_{\rho} - \Phi(y_{\rho}) + \Phi(y_{\rho}) - v \rangle + \langle m_{\rho}(v - \Phi(y_{\rho})), y_{\rho} - v \rangle \geq \langle m_{\rho}(v - \Phi(y_{\rho})), y_{\rho} - v \rangle.$$

$$(20)$$

3. Passage to the limit. Test the equation (18) with  $y_{\rho} - v$  for  $v \in V$  and use (20) to find

$$\langle Ay_{\rho}, y_{\rho} \rangle + \frac{1}{\rho} \langle m_{\rho}(v - \Phi(y_{\rho})), y_{\rho} - v \rangle \le \langle f_{\rho}, y_{\rho} - v \rangle + \langle Ay_{\rho}, v \rangle. \tag{21}$$

Now, choose an arbitrary  $v^* \in V$  with  $v^* \leq \Phi(y)$  and select the test function to be

$$v_o = v^* - \Phi(y) + \Phi(y_o).$$

With this choice, the second term on the left-hand side of the above inequality (21) is equal to zero by (10) and we find

$$\langle Ay_o, y_o \rangle \leq \langle f_o, y_o - v_o \rangle + \langle Ay_o, v_o \rangle.$$

Noting that  $v_{\rho} \to v^*$  in V (thanks to the complete continuity (19)) and  $v_{\rho} \le \Phi(y_{\rho})$ , take the limit inferior as  $\rho \to 0$  above and use weak lower semicontinuity to get  $y \in \mathbf{Q}(f)$ .

4. Strong convergence. Define  $v_{
ho}:=y+\Phi(y_{
ho})-\Phi(y)$  which has the properties

$$\begin{split} &v_\rho \to y \text{ in } V,\\ &v_\rho \le \Phi(y_\rho),\\ &y_\rho - v_\rho = (y_\rho - y) + (\Phi(y) - \Phi(y_\rho)) \rightharpoonup 0 \text{ in } V, \end{split}$$

the first holding since we already have  $y_{\rho} \rightharpoonup y$  in V. Testing (18) appropriately, we have

$$\langle A(y_{\rho} - v_{\rho}), y_{\rho} - v_{\rho} \rangle = \langle f_{\rho}, y_{\rho} - v_{\rho} \rangle - \frac{1}{\rho} \langle m_{\rho}(y_{\rho} - \Phi(y_{\rho})), y_{\rho} - v_{\rho} \rangle - \langle Av_{\rho}, y_{\rho} - v_{\rho} \rangle$$

and to this we apply the monotonicity formula and coercivity of A to find

$$C_{a} \|y_{\rho} - v_{\rho}\|_{V}^{2} \leq \langle f_{\rho}, y_{\rho} - v_{\rho} \rangle - \frac{1}{\rho} \langle m_{\rho}(v_{\rho} - \Phi(y_{\rho})), y_{\rho} - v_{\rho} \rangle - \langle Av_{\rho}, y_{\rho} - v_{\rho} \rangle$$

$$= \langle f_{\rho}, y_{\rho} - v_{\rho} \rangle - \langle Av_{\rho}, y_{\rho} - v_{\rho} \rangle.$$
 (since  $v_{\rho} \leq \Phi(y_{\rho})$ )

The right-hand side converges to zero, hence  $y_{\rho} - v_{\rho} \to 0$  strongly in V, implying  $y_{\rho} \to y$ .

Theorem 2.18 requires the complete continuity condition (19) on  $\Phi$ . Let us consider how this assumption can be weakened or substituted.

**Theorem 2.19.** Assume the conditions of Theorem 2.18, except replace the assumption (19) with

$$\langle A(\cdot), (I - \Phi)(\cdot) \rangle \colon V \to \mathbb{R}$$
 is weakly lower semicontinuous (22)

and assume one of the following:

$$\Phi \colon V \to V$$
 is weakly sequentially continuous, (23) (16), (17) and  $f_{\rho} \rightharpoonup f$  in  $H$ .

Then there exists a subsequence  $\{\rho_n\}_n$  and elements  $y_{\rho_n} \in \mathbf{P}_{\rho_n}(f_{\rho_n})$  such that  $y_{\rho_n} \rightharpoonup y \in \mathbf{Q}(f)$  in V.

*Proof.* We modify the third step of the proof of Theorem 2.18 (and, like before, we do not distinguish subsequences of  $\{\rho\}$ ). We write the final inequality of step 3 as  $\langle Ay_{\rho}, y_{\rho} - v_{\rho} \rangle \leq \langle f_{\rho}, y_{\rho} - v_{\rho} \rangle$ , which, recalling  $v_{\rho} = v^* - \Phi(y) + \Phi(y_{\rho})$ , is

$$\langle Ay_{\rho}, y_{\rho} - \Phi(y_{\rho}) \rangle + \langle Ay_{\rho}, \Phi(y) - v^* \rangle \leq \langle f_{\rho}, y_{\rho} - v_{\rho} \rangle.$$

By (22), we can take the limit inferior on the left-hand side. Regarding the right-hand side, let us consider the two cases separately.

- (1) Under (23),  $y_{\rho} v_{\rho} \rightharpoonup y v^*$  in V, and since  $f_{\rho} \to f$  in  $V^*$ , we can pass to the limit on the right-hand side and we obtain  $y \in \mathbf{Q}(f)$ , hence  $y_{\rho} \rightharpoonup y$  in V.
- (2) In the Gelfand triple case, we write the final term in the inequality above as the inner product  $(f_{\rho}, y_{\rho} v_{\rho})_H$  and pass to the limit easily.

**Remark 2.20.** It is not difficult to see that (22) and (23) are weaker assumptions than (19).

The theorem provides only weak convergence but strong convergence can be attained under additional assumptions as the next remark shows.

**Remark 2.21.** If, in addition to the conditions of Theorem 2.19 under the weak sequential continuity condition (23), we also have

$$\langle A\Phi(\cdot), (I-\Phi)(\cdot) \rangle \colon V \to \mathbb{R} \text{ is weakly lower semicontinuous}^9,$$
 (24)

then  $y_{\rho_n} - \Phi(y_{\rho_n}) \to y - \Phi(y)$  in V. To see this, returning to step 4 of the proof of Theorem 2.18 where we recall  $v_{\rho} = y + \Phi(y_{\rho}) - \Phi(y)$ , we start with the calculation

$$\begin{split} & \liminf_{\rho \to 0} \langle Av_{\rho}, y_{\rho} - v_{\rho} \rangle \\ & = \liminf_{\rho \to 0} \left( \langle A(y - \Phi(y)), (\mathbf{I} - \Phi)(y_{\rho}) + \Phi(y) - y \rangle + \langle A\Phi(y_{\rho}), (\mathbf{I} - \Phi)(y_{\rho}) \rangle + \langle A\Phi(y_{\rho}), \Phi(y) - y \rangle \right) \\ & \geq \langle A(y - \Phi(y)), (\mathbf{I} - \Phi)(y) + \Phi(y) - y \rangle + \langle A\Phi(y), (\mathbf{I} - \Phi)(y) \rangle + \langle A\Phi(y), \Phi(y) - y \rangle \\ & = 0, \end{split}$$

where for the inequality we used weak continuity for the first and last terms and (24) for the middle term. Now, taking the limit superior in the final inequality of the proof of Theorem 2.18, using the identity  $\limsup(a_n) + \liminf(b_n) \le \limsup(a_n + b_n)$  and the above calculation, we get

$$\begin{split} \limsup_{\rho \to 0} \langle f_{\rho}, y_{\rho} - v_{\rho} \rangle &\geq \limsup_{\rho \to 0} \left( C_{a} \left\| y_{\rho} - v_{\rho} \right\|_{V}^{2} + \langle A v_{\rho}, y_{\rho} - v_{\rho} \rangle \right) \\ &\geq \limsup_{\rho \to 0} C_{a} \left\| y_{\rho} - v_{\rho} \right\|_{V}^{2}. \end{split}$$

Since the left-hand side is zero (by (23)), we deduce that  $y_{\rho}-v_{\rho}\to 0$  and hence  $y_{\rho}-\Phi(y_{\rho})\to y-\Phi(y)$  in V. We see then that if for example

$$(I - \Phi)^{-1} \colon V \to V$$
 exists and is continuous,

we would also get the strong convergence  $y_{\rho} \to y$ .

**Remark 2.22.** If  $\mathbf{Q}(f)$  is a singleton, then the convergence results of the previous theorems hold for the entire sequence and not just a subsequence because the limit  $y = \mathbf{Q}(f)$  is unique.

# 3 Directional differentiability

In this section, we extend the results of our previous work [4] which dealt with directional differentiability of the source-to-solution map  $\mathbf{Q}$  associated to (1) for non-negative source terms and directions. Formally, the goal is to show that the following limit exists:

$$\lim_{s \to 0^+} \frac{\mathbf{Q}(f + sd) - \mathbf{Q}(f)}{s}.$$

This is merely a formal limit since  $\mathbf{Q}: V \rightrightarrows V$  is set valued in general, however in case  $\mathbf{Q}: V \to V$  is single valued, it is precise. It is important to obtain such a sensitivity result not only for applications but also for the procurement of certain types of stationarity conditions for optimal control problems with QVI constraints, a topic that we will address in §5.

We will follow closely the approach of our earlier work [4] where we combined an iteration (by VIs) argument with the directional differentiability result for VIs in Dirichlet space case provided by Mignot [50] but here, we make two refinements: instead of the order approach for the iterations employed in [4], we shall use a contraction technique similar to that in §2.1, and secondly, we shall use the VI differentiability result in [71] given under a general vector lattice setting, which generalises the result in [50]. For this, we begin with the following assumption on the ordering.

**Assumption 3.1.** Let V be a reflexive Banach space which is a vector lattice induced by a closed convex cone C satisfying  $C \cap -C = \{0\}$  and suppose that  $v_n \to v$  in V implies  $\sup(0, v_n) \rightharpoonup \sup(0, v)$  in V.

As before, we will identify C with  $V_+$  and note that the strong-weak convergence part of the above assumption is satisfied if there exists a constant M>0 such that  $\|\sup(0,v)\|_V\leq C\|v\|_V$  for all  $v\in V$ . To state the main result, we need to introduce some notation. Recall from (1) the constraint set mapping  $\mathbf{K}\colon V\rightrightarrows V$  defined by

$$\mathbf{K}(w) := \{ v \in V : v \le \Phi(w) \}.$$

This is convex and closed (since  $V_+$  is closed), and associated to this, we define the *radial cone* of  $\mathbf{K}(w)$  at a point  $u \in \mathbf{K}(w)$  by

$$\mathcal{R}_{\mathbf{K}(w)}(u) := \{ h \in V : \exists s^* > 0 \text{ such that } u + sh \in \mathbf{K}(w) \ \forall s \in [0, s^*] \}$$

and the corresponding tangent cone  $\mathcal{T}_{\mathbf{K}(w)}(u) := \overline{\mathcal{R}_{\mathbf{K}(w)}(u)}$ . Finally, recall the notation  $B_R(y)$  to stand for the closed ball in V of radius R centred at u.

<sup>&</sup>lt;sup>9</sup>Note that (23) and (24) imply (22). Indeed, taking the limit inferior of  $\langle Au_n, (I-\Phi)(u_n) \rangle = \langle A(I-\Phi)u_n, (I-\Phi)u_n \rangle + \langle A\Phi(u_n), (I-\Phi)(u_n) \rangle$ , using superadditivity and weak sequential continuity on the first term and (24) on the second term allows us to deduce the claim.

**Theorem 3.2.** Let Assumption 3.1 hold and given  $f \in V^*$  and  $d \in V^*$ , take  $y \in \mathbf{Q}(f)$  satisfying the local assumptions

there exists 
$$\epsilon > 0$$
 such that  $\Phi \colon B_{\epsilon}(y) \to V$  is Lipschitz with Lipschitz constant  $C_{\Phi} < C_a/(C_a + C_b)$ , (25)

$$\Phi \colon V \to V$$
 is Hadamard directionally differentiable at y. (26)

Then, for s > 0 sufficiently small, there exists  $y^s \in \mathbf{Q}(f + sd) \cap B_R(y)$  (where  $0 < R \le \epsilon$ ) and  $\alpha = \alpha(d) \in V$  such that

$$y^s = y + s\alpha + o(s)$$

where  $s^{-1}o(s) \to 0$  in V as  $s \to 0^+$  and  $\alpha$  satisfies the QVI

$$\alpha \in \mathcal{K}^{y}(\alpha) : \langle A\alpha - d, \alpha - v \rangle \leq 0 \quad \forall v \in \mathcal{K}^{y}(\alpha),$$
  
$$\mathcal{K}^{y}(\alpha) := \Phi'(y)(\alpha) + \mathcal{T}_{\mathbf{K}(y)}(y) \cap [f - Ay]^{\perp}.$$
 (27)

The directional derivative  $\alpha = \alpha(d)$  is positively homogeneous in d.

The proof of this theorem will be given in the next subsections. For now, let us make some observations.

- **Remark 3.3.** (i) The stated assumptions do **not** force solutions of the QVI to be unique. We will construct examples demonstrating this fact in §3.5.
  - (ii) If there exists an  $\epsilon$  such that  $\Phi$  is Hadamard differentiable on  $B_{\epsilon}(y)$  and

$$\forall z \in B_{\epsilon}(y), \forall v \in V, \quad \|\Phi'(z)(v)\|_{V} \le C_{\Phi} \|v\|_{V} \quad \text{where } C_{\Phi} < C_{a}/(C_{a} + C_{b}), \tag{28}$$

then (25) holds. This is immediate: take  $u, v \in B_{\epsilon}(y)$  and use the mean value theorem to find

$$\|\Phi(u) - \Phi(v)\|_{V} \le \sup_{\lambda \in (0,1)} \|\Phi'(\lambda u + (1-\lambda)v)(u-v)\|_{V} \le C_{\Phi} \|u-v\|_{V},$$

where we utilised the fact that  $\lambda u + (1 - \lambda)v \in B_{\epsilon}(y)$ . It can sometimes be easier to verify (28) than (25) depending on the problem at hand.

- (iii) The derivative  $\alpha$  is the unique solution of the QVI (27), see Proposition 3.9.
- (iv) All of the required assumptions on Φ are local, i.e., they are based at or around a neighbourhood of the chosen point y and we do not ask for them to hold globally on the whole of V. We may introduce more local assumptions in the course of the paper and one should bear in mind that such conditions are stated in terms of a fixed element y which, in later sections, need to be modified appropriately (for example in §5 such assumptions should be evaluated at the function that we call y\*). This should become apparent from the context.
- (v) In the theorem, the existence of a particular  $y \in \mathbf{Q}(f)$  is assumed; conditions under which  $\mathbf{Q}(f)$  is non-empty were given in the existence results of §2.
- (vi) This theorem generalises and improves the result of Theorem 1.6 in our earlier paper [4]. In particular, the case  $f, d \in V_+^*$  corresponds to the main result of [4] (which also requires additional assumptions).
- (vii) A differentiability result for QVIs also appears in [72, Theorem 5.5]. There, in particular, the author requires Fréchet differentiability for  $\Phi$  at y. In contrast, we require only Hadamard differentiability. In [72], A can be nonlinear of Fréchet type; we have taken A to be linear in this paper for simplicity but this can be generalised: see Remark 3.4.
- **Remark 3.4.** We have taken A to be linear for technical simplicity but an examination of the proofs that follow show that it would be possible for us to consider nonlinear A that are Hadamard differentiable in this section (a key point would be to generalise [4, Proposition 1], as we shall come to see in the proceeding). For the stationarity results of section §5.2, A would need to be continuously Fréchet differentiable. The details and the resulting changes are left to the reader.

Let us give an example of the functional setup which is typical for many applications.

**Example 3.5** (The case of a Dirichlet space). Suppose that  $H := L^2(X; \mu)$  where X is a locally compact, separable metric space and  $\mu$  is a positive Radon measure on X with full support  $^{10}$ , and let  $V \subset H$  be a dense subspace. The ordering on these spaces is given by the usual a.e. ordering of functions.

Assume that there exists a symmetric, positive semidefinite bilinear form  $\xi \colon V \times V \to \mathbb{R}$  such that endowing V with

$$(\cdot,\cdot)_V := (\cdot,\cdot)_H + \xi(\cdot,\cdot)$$

 $<sup>^{10}</sup>$ That is,  $\mu$  is a non-negative Borel measure which is finite on compact sets and strictly positive on non-empty open sets.

makes it a Hilbert space. Furthermore, we assume the Markov property<sup>11</sup>

if 
$$u \in V$$
 then  $\hat{u} := \min(u^+, 1) \in V$  and  $\xi(\hat{u}, \hat{u}) \leq \xi(u, u)$ 

and the density

$$V \cap C_c(X) \stackrel{d}{\hookrightarrow} C_c(X)$$
 and  $V \cap C_c(X) \stackrel{d}{\hookrightarrow} V$ .

The pair  $(V, \xi)$  is known as a regular Dirichlet form and V is the so-called Dirichlet space. This framework allows us to define the notions of capacity, quasi-continuity and related concepts, see [28, §2.1] and [32, §3] for more details.

In this setting, Mignot proved<sup>12</sup>the polyhedricity of sets of obstacle type in [50, Theorem 3.2] and the differentiability of VI solution maps associated to such constraint sets in [50, Theorem 3.3]. We also have an explicit expression for the critical cone appearing in (27) via [50, Lemma 3.2]:

$$\mathcal{K}^y(w) := \{ \varphi \in V : \varphi \leq \Phi'(y)(w) \text{ q.e. on } \mathcal{A}(y) \text{ and } \langle Ay - f, \varphi - \Phi'(y)(w) \rangle = 0 \}.$$

Here, 'q.e.' stands for quasi-everywhere and a statement holds quasi-everywhere if it holds everywhere except on a set of capacity zero, and A(y) refers to the active or coincidence set of the solution y to the QVI related to an obstacle map  $\Phi$ , i.e.,

$$\mathcal{A}(y) := \{ x \in X : y(x) = \Phi(y)(x) \} \quad \text{for } y \in V.$$

We in fact take the quasi-continuous representatives of the functions appearing in the definition so that A(y) is quasiclosed and defined up to sets of capacity zero. It is important to note that the set of points defining the active set is taken over X; in the context of some Sobolev spaces over a domain  $\Omega$ , this can sometimes be  $X = \overline{\Omega}$  and not merely  $\Omega$ , see [4, §1.2] for more details.

Before we proceed, let us provide some notation. Define the critical cone

$$\mathcal{K}^y := \mathcal{T}_{\mathbf{K}(y)}(y) \cap [f - Ay]^{\perp},\tag{29}$$

and observe the relation

$$\mathcal{K}^{y}(w) = \Phi'(y)(w) + \mathcal{K}^{y}.$$

Recall that the  $polar\ cone$  of a set  $M\subset V$  is defined as

$$M^{\circ} = \{ g \in V^* : \langle g, v \rangle \le 0 \quad \forall v \in M \}.$$

#### 3.1 Iteration scheme and expansion formulae

To prove Theorem 3.2, we employ an iteration and passage to the limit approach like in our previous work [4]. We fix an arbitrary  $f \in V^*$  and take an arbitrary but fixed  $y \in \mathbf{Q}(f)^{13}$ . Pick a direction  $d \in V^*$  and construct, similarly to §2.1, the sequence

$$y_0^s := y,$$
  
 $y_n^s := S(f + sd, y_{n-1}^s).$  (30)

The idea here is to expand each  $y_n^s$  in terms of y, a directional derivative and a remainder term (both of these would depend on n) and then to pass to the limit in such an expansion. The natural way to proceed would be to obtain a uniform bound on  $\{y_n^s\}$  which would result in the existence of a weakly convergent *subsequence*  $\{y_{n_j}^s\}$ . This is not enough to identify the limit of  $\{y_{n_j}^s\}$  due to the (n-1) index in the definition of  $y_n^s$ , so one would need convergence of the whole sequence which holds true when, for example, one has monotonicity. However, in contrast to the sequence considered in §2.1, we do not obtain any monotonicity of  $\{y_n^s\}$  since we do not assume a sign on d nor do we assume monotonicity of  $\Phi$ . Therefore, for convergence of the full sequence, we instead look for a contraction of the map associated to  $\{y_n^s\}$  on some small ball.

**Lemma 3.6.** Assume the Lipschitz property (25). Then for any  $0 < R \le \epsilon$ ,  $S(f+sd,\cdot)$ :  $B_R(y) \to B_R(y)$  is a contraction whenever

$$s \le C_a \|d\|_{V^*}^{-1} R(1 - (1 + C_b C_a^{-1}) C_{\Phi}).$$

*Proof.* Let  $v \in B_R(y)$ ; we want to show that  $S(f + sd, v) \in B_R(y)$ . Observe that, using y = S(f, y) and continuous dependence (e.g. [4, Equation (21)]),

$$\begin{split} \|S(f+sd,v)-y\|_{V} &\leq (1+C_{b}C_{a}^{-1}) \, \|\Phi(v)-\Phi(y)\|_{V} + C_{a}^{-1}s \, \|d\|_{V^{*}} \\ &\leq (1+C_{b}C_{a}^{-1})C_{\Phi} \, \|v-y\|_{V} + C_{a}^{-1}s \, \|d\|_{V^{*}} \\ &\leq (1+C_{b}C_{a}^{-1})C_{\Phi}R + C_{a}^{-1}s \, \|d\|_{V^{*}} \,, \end{split} \tag{since } v,y \in B_{R}(y) \subset B_{\epsilon}(y) )$$

<sup>&</sup>lt;sup>11</sup>This is also known as the *unit contraction property*.

<sup>&</sup>lt;sup>12</sup>In fact, Mignot uses a weaker setting of *positivity-preserving forms* rather than the Dirichlet form setting described here with also some other weaker conditions.

<sup>&</sup>lt;sup>13</sup>Again, see §2 for existence of such y.

and, using the fact that  $(1 + C_b C_a^{-1})C_{\Phi}$  equals a constant strictly less than 1, the right-hand side is bounded above by R. This shows that  $S(f+sd,\cdot)$  maps  $B_R(y)$  into itself. To see that the map is a contraction, take  $v,w\in B_R(y)$  and observe that

$$||S(f+sd,v) - S(f+sd,w)||_{V} \le (1 + C_{a}^{-1}C_{b}) ||\Phi(v) - \Phi(w)||_{V} \le C_{\Phi}(1 + C_{a}^{-1}C_{b}) ||v - w||_{V}.$$

Hence, under (25), we have that each  $y_n^s \in B_R(y)$ . By applying the Banach fixed point theorem, we obtain the following existence and convergence result.

**Proposition 3.7.** Given  $f, d \in V^*$  and  $y \in \mathbf{Q}(f)$ , under (25) and sufficiently small s > 0, there exists  $y^s \in \mathbf{Q}(f + sd) \cap B_R(y)$  such that  $y_n^s \to y^s$  in V (where  $y_n^s$  is defined in (30)).

Since we want to study differentiability of QVIs, we need some differentiability for the constraint set mapping. We will henceforth assume the Hadamard differentiability at y condition (26). Now, making use of [71, Theorems 4.18 and 5.2] we can expand  $y_1^s = S(f + sd, y)$  as follows:

$$y_1^s = y + s\alpha_1 + o_1(s),$$

where  $s^{-1}o_1(s) \to 0$  as  $s \to 0^+$  and  $\alpha_1 = \partial S(f, y)(d)$  is the directional derivative of  $S(\cdot, y)$  in the direction d, and this satisfies the VI (recall  $\mathcal{K}^y$  from (29))

$$\alpha_1 \in \mathcal{K}^y : \langle A\alpha_1 - d, \alpha_1 - v \rangle \le 0 \quad \forall v \in \mathcal{K}^y.$$

To acquire an expansion formula for a general  $y_n^s$ , define for n>1,

$$\alpha_n := \Phi'(y)(\alpha_{n-1}) + \partial S(f, y)(d - A\Phi'(y)(\alpha_{n-1})).$$

In exactly the same way as in [4, Proposition 2], we obtain the following result.

**Proposition 3.8.** *Under* (25) *and* (26), *for*  $n \ge 1$ ,

$$y_n^s = y + s\alpha_n + o_n(s) \tag{31}$$

where  $s^{-1}o_n(s) \to 0$  as  $s \to 0^+$  and  $\alpha_n = \alpha_n(d)$  is positively homogeneous in the direction d and satisfies the VI

$$\alpha_n \in \mathcal{K}^y(\alpha_{n-1}) : \langle A\alpha_n - d, \alpha_n - \varphi \rangle \le 0 \quad \forall \varphi \in \mathcal{K}^y(\alpha_{n-1}),$$
  
$$\mathcal{K}^y(\alpha_{n-1}) := \mathcal{K}^y + \Phi'(y)(\alpha_{n-1}).$$

See (35) for the precise definition of  $o_n$ . The proof of this proposition, which we omit here, is by induction and makes use of the expansion formula of [4, Proposition 1], which tells us that

$$y_{n+1}^{s} = S(f + sd, y + s\alpha_{n} + o_{n}(s)) = y + s(\Phi'(y)(\alpha_{n}) + \partial S(f, y)(d - A\Phi'(y)(\alpha_{n}))) + o_{n+1}(s).$$

It remains then to pass to the limit in (31) and to identify the corresponding limits.

#### 3.2 Passage to the limit

Observe that the conditions (25) and (26) imply that

$$\Phi'(y): V \to V$$
 is Lipschitz with Lipschitz constant  $C_L < C_a/C_b$ , (32)

which is precisely what is needed for the coming intermediary results. In particular, it allows for the Banach fixed point theorem to be amenable to show the convergence of  $\{\alpha_n\}$  as the next proposition demonstrates. But first, let us prove that (32) is indeed a consequence. From the expansion formula  $\Phi(y+sh)=\Phi(y)+s\Phi'(y)(h)+o(s;h)$  where  $o(\cdot,h)$  is a remainder term, we find

$$\|\Phi'(y)(h) - \Phi'(y)(d)\|_{V} \le \frac{1}{s} \|\Phi(y + sh) - \Phi(y + sd)\|_{V} + \frac{1}{s} \|o(s; d) - o(s; h)\|_{V}.$$

Without loss of generality, we may assume that at least one of h and d is non-zero. We see that if  $s \le \epsilon/(\|h\|_V + \|d\|_V)$ , we have  $y + sh, y + sd \in B_{\epsilon}(y)$  and therefore, by (25),

$$\|\Phi'(y)(h) - \Phi'(y)(d)\|_{V} \le C_{\Phi} \|h - d\|_{V} + \frac{1}{s} \|o(s; d) - o(s; h)\|_{V}.$$

Taking  $s \to 0^+$  we obtain the statement after noting that  $C_{\Phi} < C_L$ .

**Proposition 3.9.** Under (26) and (32),  $\alpha_n \to \alpha$  in V where  $\alpha$  is the unique solution of the QVI (27).

*Proof.* Denote by  $T \colon V \to V$  the solution map  $\gamma \mapsto \beta$  of the inequality

$$\beta \in \mathcal{K}^y(\gamma) : \langle A\beta - d, \beta - \varphi \rangle \le 0 \quad \forall \varphi \in \mathcal{K}^y(\gamma).$$

This has a unique solution by the Lions-Stampacchia theorem [61, §4:3, Theorem 3.1], hence T is well defined.

Consider  $\gamma_1, \gamma_2 \in V$  with  $\beta_1 := T(\gamma_1)$  and  $\beta_2 := T(\gamma_2)$ . Testing the inequality for  $\beta_1$  with the feasible element  $\beta_2 - \Phi'(y)(\gamma_2) + \Phi'(y)(\gamma_1)$  and vice versa and then combining both of the resulting inequalities, we find

$$\langle A(\beta_1 - \beta_2), \beta_1 - \beta_2 + \Phi'(y)(\gamma_2) - \Phi'(y)(\gamma_1) \rangle \le 0,$$

which implies, using (32),

$$\|\beta_1 - \beta_2\|_V \le \frac{C_b}{C_a} \|\Phi'(y)(\gamma_2) - \Phi'(y)(\gamma_1)\|_V < \|\gamma_2 - \gamma_1\|_V.$$

This shows that  $T\colon V\to V$  is a contraction. Therefore, thanks to the Banach fixed point theorem, the iterative sequence  $\beta_n:=T(\beta_{n-1}),\,\beta_1:=\alpha_1$ , is such that  $\beta_n\equiv\alpha_n$  (by uniqueness of solutions) and  $\alpha_n\to\alpha$  strongly in V where  $\alpha$  is the fixed point of T.

Thanks to this result, it follows that  $o_n(s) \to o^*(s)$  in V for some  $o^*(s)$ . We can send  $n \to \infty$  in (31) to obtain

$$y^s = y + s\alpha + o^*(s),$$

and it is left for us to show that  $o^*$  is a remainder term. The idea in [4] was to show that the convergence  $s^{-1}o_n(s) \to 0$  as  $s \to 0^+$  is *uniform* in n, which is sufficient to commute the limits  $s \to 0^+$  and  $n \to \infty$  for  $s^{-1}o_n(s)$ , giving the desired behaviour  $s^{-1}o^*(s) \to 0$  as  $s \to 0^+$ . This was done in [4, Lemma 14], the proof of which we will now adapt under the context of our current (more general) setting. For this, we need some more notation. For  $v \in V$  and  $h_s \in V$ , we define the remainder term associated to  $\Phi$ 

$$\hat{l}(s, h, h_s; v) := \Phi(v + sh_s) - \Phi(v) - s\Phi'(v)(h), \tag{33}$$

and since  $\Phi$  is Hadamard differentiable at y, if  $h_s \to h$  in V as  $s \to 0^+$ , then  $s^{-1}\hat{l}(s,h,h_s;y) \to 0$  as  $s \to 0^+$ . We write  $\hat{l}(s,h,h;v) = l(s,h;v)$  when  $h_s \equiv h$ . Now let  $S_0 \colon V^* \to V$  be the map  $f \mapsto u$  of the following VI with zero lower obstacle:

$$u \in V_+ : \langle Au - f, u - v \rangle \le 0 \quad \forall v \in V_+.$$

In a similar fashion to  $\hat{l}$ , we denote the remainder term associated to the expansion formula of  $S_0$  by  $\hat{o}$ :

$$\hat{o}(s, h, h_s; f) := S_0(f + sh_s) - S_0(f) - sS_0'(f)(h).$$

**Proposition 3.10.** Under (25) and (26),  $s^{-1}o^*(s) \to 0$  as  $s \to 0^+$ .

*Proof.* Since  $y + s\alpha_n + o_n(s) = y_n^s \in B_{\epsilon}(y)$  and  $\{\alpha_n\}$  is bounded, let us say by M, if  $s \leq M^{-1}\epsilon$ , then  $y + s\alpha_n \in B_{\epsilon}(y)$  too. Hence, from (33) and the Lipschitz property (25), we have

$$\begin{aligned} \left\| \hat{l}(s, \alpha_{n}, \alpha_{n} + s^{-1}o_{n}(s); y) \right\|_{V} &\leq \left\| \hat{l}(s, \alpha_{n}, \alpha_{n} + s^{-1}o_{n}(s); y) - l(s, \alpha_{n}; y) \right\|_{V} + \left\| l(s, \alpha_{n}; y) \right\|_{V} \\ &= \left\| \Phi(y + s(\alpha_{n} + s^{-1}o_{n}(s))) - \Phi(y + s\alpha_{n}) \right\|_{V} + \left\| l(s, \alpha_{n}; y) \right\|_{V} \\ &\leq C_{\Phi} \left\| o_{n}(s) \right\|_{V} + \left\| l(s, \alpha_{n}; y) \right\|_{V}. \end{aligned}$$
(34)

We see from [4, Equation (34) and Proposition 1] that  $o_n$  has the definition

$$o_n(s) := \hat{l}(s, \alpha_{n-1}, \alpha_{n-1} + s^{-1}o_{n-1}(s); y) - \hat{o}(s, A\Phi'(y)(\alpha_{n-1}) - d, A\Phi'(y)(\alpha_{n-1}) - d + As^{-1}\hat{l}(s, \alpha_{n-1}, \alpha_{n-1} + s^{-1}o_{n-1}(s)); A\Phi(y) - f).$$
 (35)

For ease of reading, let us omit the base point from the expressions for  $\hat{l}, l, \hat{o}$  and o from now on. That is, we write  $\hat{l}(\cdot, \cdot, \cdot)$  instead of  $\hat{l}(\cdot, \cdot, \cdot; y)$  and likewise for the other terms. In the above equality, taking norms and, on the right-hand side, using (34) on the first term and the corresponding estimate

$$\|\hat{o}(s, h, h + s^{-1}h_s)\|_{V} \le C_a^{-1} \|h_s\|_{V^*} + \|o(s, h)\|_{V}$$

for  $S_0$  and its remainder term (see [4, Lemma 1]) on the second term, we find

$$\begin{aligned} \|o_n(s)\|_V &\leq C_{\Phi} \|o_{n-1}(s)\|_V + \|l(s,\alpha_{n-1})\|_V + C_a^{-1}C_b \|\hat{l}(s,\alpha_{n-1},\alpha_{n-1}+s^{-1}o_{n-1}(s))\|_V + \|o(s,A\Phi'(y)(\alpha_{n-1})-d)\|_V \\ &\leq C_{\Phi} (1 + C_a^{-1}C_b) \|o_{n-1}(s)\|_V + (1 + C_a^{-1}C_b) \|l(s,\alpha_{n-1})\|_V + \|o(s,A\Phi'(y)(\alpha_{n-1})-d)\|_V, \end{aligned}$$

where we again used (34) on the penultimate term in the first line to obtain the second inequality. Defining

$$a_n(s) := \|o_n(s)\|_V$$
 and  $b_n(s) := (1 + C_a^{-1}C_b) \|l(s, \alpha_n)\|_V + \|o(s, A\Phi'(y)(\alpha_n) - d)\|_V$ ,

the above can be recast as

$$a_n(s) \le Ca_{n-1}(s) + b_{n-1}(s)$$

for some C<1 by the assumption on  $C_\Phi$  in (25). Solving this recurrence inequality gives

$$a_n(s) \le C^{n-1}a_1(s) + C^{n-2}b_1(s) + C^{n-3}b_2(s) + \dots + Cb_{n-2}(s) + b_{n-1}(s).$$
 (36)

Now, consider

$$\frac{b_{n-1}(s)}{s} = \frac{(1 + C_a^{-1}C_b) \|l(s, \alpha_{n-1})\|_V}{s} + \frac{\|o(s, A\Phi'(y)(\alpha_{n-1}) - d)\|_V}{s}.$$

By Proposition 3.9,  $\alpha_n \to \alpha$  strongly in V, thus  $\overline{\{\alpha_{n-1}\}}$  and  $\overline{\{A\Phi'(y)(\alpha_{n-1})-d\}}$  are compact sets in V and  $V^*$  respectively. Since the remainder terms l and o appearing in the displayed equality above arise from the Hadamard (and hence compact) differentiability of  $\Phi$  and the solution map  $S_0$  associated to VIs, it follows that  $l(s,\gamma)/s$  and o(s,h)/s both converge to zero uniformly for  $\gamma$  and h belonging to  $\overline{\{\alpha_{n-1}\}}$  and  $\overline{\{A\Phi'(y)(\alpha_{n-1})-d\}}$  respectively. Because  $\{\alpha_{n-1}\}\subset\overline{\{\alpha_{n-1}\}}$  and  $\{A\Phi'(y)(\alpha_{n-1})-d\}\subset\overline{\{A\Phi'(y)(\alpha_{n-1})-d\}}$ , we have that

$$\frac{l(s,\gamma)}{s} \to 0 \text{ uniformly in } \gamma \in \{\alpha_{n-1}\} \qquad \text{and} \qquad \frac{o(s,h)}{s} \to 0 \text{ uniformly in } h \in \{A\Phi'(y)(\alpha_{n-1}) - d\},$$

which then gives

$$\frac{b_{n-1}(s)}{s} \to 0 \quad \text{uniformly in } n.$$

This, along with (36) and the geometric series estimate  $C^{n-2}+C^{n-3}+\ldots+C+1=(1-C^{n-1})/(1-C)\leq 1/(1-C)$  implies that for every  $\epsilon>0$ , there exists an  $s_0$  independent of n such that

$$\frac{\|o_n(s)\|_V}{s} \le \epsilon \quad \text{when } s \le s_0$$

which means precisely that  $s^{-1}o_n(s) \to 0$  as  $s \to 0^+$  uniformly in n. Finally, recalling that  $o_n(s)$  converges in V, taking the limit as  $n \to \infty$  in the above inequality, we deduce that  $s^{-1}o^*(s) \to 0$  as  $s \to 0^+$ .

This concludes the proof of Theorem 3.2.

**Remark 3.11.** It is worth noting that the complete continuity assumption (19) is not needed for the result (the strong convergence of  $\{y_n^s\}$  assured by the application of the Banach fixed point theorem allowed us to circumvent complete continuity). Furthermore, complete continuity of  $\Phi'(y)$  is not needed for the characterisation of the directional derivative; continuity suffices (which is guaranteed since Hadamard derivatives are continuous with respect to the direction), unlike in §5.1 and §5.2 of [4].

#### 3.3 Continuity properties of the directional derivative

We now study the conditions under which continuity of the map taking the direction d into the directional derivative  $\alpha$  in (27) is assured. We recall (27) for convenience:

$$\alpha \in \mathcal{K}^y(\alpha) : \langle A\alpha - d, \alpha - v \rangle \le 0 \quad \forall v \in \mathcal{K}^y(\alpha),$$
  
 $\mathcal{K}^y(w) := \mathcal{K}^y + \Phi'(y)(w).$ 

**Proposition 3.12.** Under (32),  $d \mapsto \alpha(d)$  is continuous from  $V^*$  to V. That is, if  $d_j \to d$  in  $V^*$ , then

$$\alpha_j \to \alpha$$
 in  $V$ 

where  $\alpha_j$  and  $\alpha$  are the solutions of (27) with source terms  $d_j$  and d respectively.

*Proof.* The element  $\alpha_j$  associated to  $d_j$  satisfies

$$\alpha_i \in \mathcal{K}^y(\alpha_i) : \langle A\alpha_i - d_i, \alpha_i - v \rangle \leq 0 \quad \forall v \in \mathcal{K}^y(\alpha_i).$$

Take  $j, k \in \mathbb{N}$  and in the inequality for  $\alpha_j$ , take the test function  $v = \alpha_k - \Phi'(y)(\alpha_k) + \Phi'(y)(\alpha_j)$  which is clearly feasible, whilst in the inequality for  $\alpha_k$ , set  $v = \alpha_j - \Phi'(y)(\alpha_j) + \Phi'(y)(\alpha_k)$  to obtain

$$\langle A\alpha_j - d_j, \alpha_j - \alpha_k + \Phi'(y)(\alpha_k) - \Phi'(y)(\alpha_j) \rangle \le 0,$$
  
 $\langle A\alpha_k - d_k, \alpha_k - \alpha_n + \Phi'(y)(\alpha_j) - \Phi'(y)(\alpha_k) \rangle \le 0.$ 

Adding these, we find

$$\langle A(\alpha_i - \alpha_k) - (d_i - d_k), \alpha_i - \alpha_k + \Phi'(y)(\alpha_k) - \Phi'(y)(\alpha_n) \rangle \le 0,$$

which implies, using (32),

$$C_{a} \|\alpha_{j} - \alpha_{k}\|_{V}^{2} \leq \|d_{j} - d_{k}\|_{V^{*}} \|\alpha_{j} - \alpha_{k}\|_{V} + C_{b} \|\alpha_{k} - \alpha_{j}\|_{V} \|\Phi'(y)(\alpha_{k}) - \Phi'(y)(\alpha_{j})\|_{V}$$

$$+ \|d_{k} - d_{j}\|_{V^{*}} \|\Phi'(y)(\alpha_{k}) - \Phi'(y)(\alpha_{j})\|_{V}$$

$$\leq \|d_{j} - d_{k}\|_{V^{*}} \|\alpha_{j} - \alpha_{k}\|_{V} + C_{b}C_{L} \|\alpha_{k} - \alpha_{j}\|_{V}^{2} + C_{L} \|d_{k} - d_{j}\|_{V^{*}} \|\alpha_{k} - \alpha_{j}\|_{V}.$$

Manipulating, we find that  $\{\alpha_i\}$  is a Cauchy sequence and thus there exists an  $\alpha \in V$  with

$$\alpha_i \to \alpha$$
 in  $V$ .

Now, in the inequality for  $\alpha_j$ , choose the test function  $v_j := v - \Phi'(y)(\alpha) + \Phi'(y)(\alpha_j)$  where v is such that  $v \in \mathcal{K}^y(\alpha)$ . It follows that  $v_j \to v$  in V. This allows us to pass to the limit and we get

$$\langle A\alpha - d, \alpha - v \rangle \le 0 \quad \forall v \in \mathcal{K}^y(\alpha)$$

and it remains to be seen that  $\alpha \in \mathcal{K}^y(\alpha)$ , which is evident since the critical cone is closed.

#### 3.4 Complementarity characterisation of the directional derivative

We now look for an analogue of the complementarity characterisation of Proposition 2.1 for the QVI (27) satisfied by the directional derivative.

**Proposition 3.13.** The QVI (27) is equivalent to the complementarity system

$$\alpha - \Phi'(y)(\alpha) \in \mathcal{K}^y, \tag{37a}$$

$$\xi_d = d - A\alpha, \tag{37b}$$

$$\xi_d \in (\mathcal{K}^y)^\circ,$$
 (37c)

$$\langle \xi_d, \Phi'(y)(\alpha) - \alpha \rangle = 0.$$
 (37d)

*Proof.* As noted above,  $\alpha - \Phi'(y)(\alpha)$  belongs to the set  $\mathcal{K}^y$ . Define  $\xi_d := d - A\alpha$  which by definition satisfies

$$\alpha - \Phi'(y)(\alpha) \in \mathcal{K}^y : \langle \xi_d, \alpha - v \rangle \ge 0 \quad \forall v \in V : v - \Phi'(y)(\alpha) \in \mathcal{K}^y.$$

Taking  $v = \Phi'(y)(\alpha)$  here and then  $v = 2\alpha - \Phi'(y)(\alpha)$  (which is feasible since  $v - \Phi'(y)(\alpha)$  is twice a function that belongs to  $\mathcal{K}^y$ ) shows the orthogonality condition (37d).

Let  $w \in \mathcal{K}^y$  and select  $v = \alpha + w$  (this is feasible since  $v - \Phi'(y)(\alpha) = \alpha - \Phi'(y)(\alpha) + w \in \mathcal{K}^y + \mathcal{K}^y$  and the tangent cone, being a convex cone, is closed under addition). With this choice, we obtain

$$\langle \xi_d, w \rangle \leq 0 \quad \forall w \in \mathcal{K}^y,$$

meaning precisely that  $\xi_d \in (\mathcal{K}^y)^\circ$ . The reverse direction holds by the same trick as in the proof of Proposition 2.1.

## 3.5 Examples of QVIs with multiple solutions

In this section, we construct explicit examples of QVIs with non-unique solutions such that the assumptions of Theorem 3.2 are satisfied, thus verifying that multiplicity of solutions is not lost under our assumptions.

**Example 1** Let  $\Omega \subset \mathbb{R}^n$  be a bounded Lipschitz domain and set  $V := H^k(\Omega)$  with  $H = L^2(\Omega)$  forming a Gelfand triple. Below, all norms and inner products that appear are over H.

Pick  $\delta>0$  and select a sequence  $\{y_n\}_{n=1}^N$  of smooth functions satisfying  $\|y_n-y_m\|^2>4\delta^2$  for each  $m,n\in\{1,\ldots,N\}$  with  $m\neq n$  and  $N\geq 2$  fixed. Take a smooth cutoff function  $\nu\in C^\infty(\mathbb{R})$  with  $0\leq \nu\leq 1$  and

$$\nu(t) = \begin{cases} 1 & : \text{if } t \in (-\delta^2, \delta^2), \\ 0 & : \text{if } |t| \ge 2\delta^2. \end{cases}$$

For a parameter  $y \in V$ , define the map  $\Phi_y \colon V \to V$  by

$$\Phi_{u}(u) := \nu(\|u - y\|^{2})y$$

and set

$$\Phi(u) := \sum_{n=1}^{N} \Phi_{y_n}(u).$$

Note that  $\Phi \colon V \to V$  and  $\Phi(y_n) = y_n$  (because  $\Phi_{y_n}(y_m) = y_n \delta_{nm}$ ). Let the elliptic operator  $A \colon V \to V^*$  have the property that  $Ay_n \in H$  for each n and define the pointwise a.e. maximum  $f := \max(Ay_1, \cdots, Ay_N) \in H$ . Then the QVI

find 
$$u \le \Phi(u) : \langle Au - f, u - v \rangle \le 0 \quad \forall v \in V : v \le \Phi(u)$$

has multiple solutions and indeed each  $y_n \in \mathbf{Q}(f)$  is a solution. To see this, simply observe that  $Ay_n - f \leq 0$  and  $y_n - v = \Phi(y_n) - v \geq 0$  for all  $v \in V$  with  $v \leq \Phi(y_n)$ .

It follows from the expression  $\Phi_y'(u)(h) = 2y\nu'(\|u-y\|^2)(h,u-y)$  that

$$\Phi'(u)(h) = \sum_{n=1}^{N} 2y_n \nu'(\|u - y_n\|^2)(h, u - y_n)$$

and hence  $\Phi'(B_{\delta}(y_n)) \equiv 0$  and thus (28) is trivially satisfied (hence also (25) and (32) by Remark 3.3 (ii) and the digression at the start of §3.2). Hence, all the requirements of Theorem 3.2 have been met and we obtain for every  $d \in V^*$  the existence of of  $y_m^s \in \mathbf{Q}(f+sd)$  and  $\alpha_m \in V$  such that

$$\lim_{s \to 0^+} \frac{y_m^s - y_m}{s} = \alpha_m.$$

Let us also note that in addition,  $\Phi'(y_n) \colon V \to V$  is completely continuous.

**Example 2** A second example, without the need for the source term f to be defined in terms of  $\{y_n\}$ , can be given under the same initial setting as above. For  $n=1,\ldots,N$ , take  $\psi_n\in V$  to be given distinct obstacles such that the associated solutions  $y_n\in V$  of the VIs

$$y_n \le \psi_n : \langle Ay_n - f, y_n - v \rangle \le 0 \quad \forall v \in V : v \le \psi_n$$

are distinct too. We suppose that  $\delta$  is chosen such that  $\|y_n - y_m\|^2 > 4\delta^2$ , which is possible since the  $y_n$  are distinct functions. With  $\nu$  as above, define now  $\Phi_n \colon V \to V$  by

$$\Phi_n(u) := \nu(\|u - y_n\|^2)\psi_n$$

and set

$$\Phi(u) := \sum_{n=1}^{N} \Phi_n(u).$$

We have  $\Phi(y_n) = \psi_n$  and each  $y_n$  is again a solution of the QVI associated to  $\Phi$  with source term f, i.e.,  $y_n \in \mathbf{Q}(f)$ . Furthermore,

$$\Phi'(u)(h) = \sum_{n=1}^{N} 2\psi_n \nu'(\|u - y_n\|^2)(h, u - y_n)$$

and we can argue as before to derive the other properties and results.

# 4 Existence of optimal controls

We now address the optimal control problem (2). Regarding the function space context in this section, we take

- (i)  $V \hookrightarrow H$  to be a continuous embedding of reflexive Banach spaces,
- (ii) U to be a reflexive Banach space with  $U \stackrel{c}{\hookrightarrow} V^*$ ,
- (iii)  $U_{ad} \subseteq U$  to be a non-empty and weakly sequentially closed<sup>14</sup> set.

Given  $\nu > 0$  and a desired state  $y_d \in H$ , define  $J \colon H \times U \to \mathbb{R}$  by

$$J(y, u) := \frac{1}{2} \|y - y_d\|_H^2 + \frac{\nu}{2} \|u\|_U^2,$$

and consider the problem (2) which we recall here:

$$\min_{\substack{u \in U_{ad} \\ y \in \mathbf{Q}(u)}} J(y, u).$$

<sup>&</sup>lt;sup>14</sup>That is, if  $u_n \rightharpoonup u$  in U with  $u_n \in U_{ad}$ , then  $u \in U_{ad}$ .

**Theorem 4.1.** Let Assumption 2.12 hold, suppose that  $\mathbf{Q}(u)$  is non-empty<sup>15</sup> for every  $u \in U_{ad}$  and let the feasibility condition (6) and the complete continuity (19) hold. Then there exists an optimal control  $u^* \in U_{ad}$  and associated state  $y^* \in \mathbf{Q}(u^*)$  to the problem (2).

*Proof.* Let  $u_n \in U_{ad}$  be an infimising sequence with  $y_n \in \mathbf{Q}(u_n)$ , i.e.,

$$J(y_n, u_n) \to \inf_{\substack{u \in U_{ad}, \\ y \in \mathbf{Q}(u)}} J(y, u).$$

Then  $\{u_n\}$  and  $\{y_n\}$  are bounded in U and V respectively (the latter arises from (6)) and therefore, there exists a subsequence such that

$$u_{n_j} \rightharpoonup u^* \text{ in } U$$
 and  $y_{n_j} \rightharpoonup y^* \text{ in } V$ .

By assumption,  $u^*$  also belongs to  $U_{ad}$ . Since the  $y_n$  are solutions of QVIs, we have the following estimate

$$||y_{n_j} - y_{n_k}||_V \le C(||u_{n_j} - u_{n_k}||_{V^*} + ||\Phi(y_{n_j}) - \Phi(y_{n_k})||_V).$$

In the limit, the first term on the right-hand side vanishes due to the compact embedding, and the second term vanishes too because  $\Phi$  is completely continuous due to (19). Thus  $\{y_{n_j}\}$  is Cauchy in V and  $y_{n_j} \to y^*$  in V. Taking an arbitrary  $v \in V$  such that  $v \leq \Phi(y^*)$ , we set  $v_{n_j} := v - \Phi(y^*) + \Phi(y_{n_j})$  and use this as a test function in the QVI for  $y_{n_j}$  in which we can pass to the limit to find  $y^* \in \mathbf{Q}(u^*)$ . To see that this pair is optimal, we observe that (dispensing with the subsequence notation now), using the continuity of the embedding  $V \hookrightarrow H$ ,

$$J(y^*, u^*) \le \liminf_{n \to \infty} J(y_n, u_n) \le \lim_{n \to \infty} J(y_n, u_n) = \min_{\substack{u \in U_{ad} \\ y \in \mathbf{Q}(u)}} J(y, u).$$

Regarding regularity of the optimal control, see Theorem 5.11. In general there is no uniqueness for the optimal control and state regardless of whether  $\mathbf{Q}$  is single valued or not.

## 4.1 The penalised optimal control problem

Let us return to the context of §2.3 and consider for each  $\rho > 0$  the penalisation of (2):

$$\min_{u \in U_{ad}} J(y_{\rho}, u) \quad \text{such that} \quad Ay_{\rho} + \frac{1}{\rho} m_{\rho} (y_{\rho} - \Phi(y_{\rho})) = u. \tag{38}$$

We remind the reader that  $m_{\rho}$  is taken to satisfy (10)–(12). Recalling the map  $\mathbf{P}_{\rho}$  from §2.3, we can write the equation above as  $y_{\rho} \in \mathbf{P}_{\rho}(u)$ . The reason for considering this problem is because we will use this to derive stationarity conditions in the next section but first let us check that this minimisation problem suitably approximates (2).

**Proposition 4.2.** Let Assumption 2.12, (14), (15a) and (19) hold and suppose that  $\mathbf{Q}$  is single valued. Then there exist optimal pairs  $(y_{\rho}^*, u_{\rho}^*)$  of (38) and an optimal pair  $(y^*, u^*)$  of (2) such that

$$(y_{\rho}^*, u_{\rho}^*) \rightarrow (y^*, u^*)$$
 in  $V \times U$ .

*Proof.* First, observe that  $\mathbf{P}_{\rho}(u)$  is non-empty for all  $u \in U_{ad}$  by Proposition 2.15 (after possibly renorming V, see Example 2.14). Now, let  $(y_{\rho}^*, u_{\rho}^*)$  denote an optimal pair of (38), which exists by standard arguments (like in the proof of Theorem 4.1) making use of (15a) (to show weak continuity of the solution map). By definition,

$$J(y_o^*, u_o^*) \le J(w_o, u) \quad \forall u \in U_{ad}, \quad \forall w_o \in \mathbf{P}_o(u). \tag{39}$$

Given any  $\tilde{u} \in U_{ad}$ , we pick a subsequence  $\{\tilde{y}_{\rho_n}\}$  such that  $\mathbf{P}_{\rho_n}(\tilde{u}) \ni \tilde{y}_{\rho_n} \to \tilde{y}$  where  $\tilde{y} \in \mathbf{Q}(\tilde{u})$ ; this is possible by Theorem 2.18. The inequality (39) implies that  $J(y_{\rho_n}^*, u_{\rho_n}^*)$  is bounded above by  $J(\tilde{y}_{\rho_n}, \tilde{u})$  which in turn is bounded uniformly in  $\rho_n$  because  $\tilde{y}_{\rho_n}$  is bounded in V by the estimate of Proposition 2.15:

$$\|y_{\rho_n}^*\|_V \le C \left(\|u_{\rho_n}^*\|_{V^*} + \|v_0\|_V\right).$$

Hence for another subsequence (which we shall relabel)

$$u_{\rho_n}^* \rightharpoonup u^* \quad \text{in } U_{ad},$$
  
 $y_{\rho_n}^* \rightharpoonup y^* \quad \text{in } V,$ 

<sup>&</sup>lt;sup>15</sup>See §2.

for some  $(u^*, y^*)$  that we need to show is an optimal pair. By following steps 3 and 4 in the proof of Theorem 2.18,  $y_{\rho_n}^* \to y^* = \mathbf{Q}(u^*)$  in V (since  $u_{\rho_n}^* \to u^*$  in  $V^*$ ). Hence  $(y^*, u^*)$  is a feasible point of (2). Then observe that for  $(\hat{y}, \hat{u})$  being any optimal point of (2),

$$J(\hat{y},\hat{u}) \leq J(y^*,u^*) \leq \liminf_{n \to \infty} J(y^*_{\rho_n},u^*_{\rho_n}) \leq \limsup_{n \to \infty} J(y^*_{\rho_n},u^*_{\rho_n}) \leq \limsup_{n \to \infty} J(w^*_{\rho_n},\hat{u}) \quad \forall w^*_{\rho_n} \in \mathbf{P}_{\rho_n}(\hat{u})$$

with the last inequality by (39). Now it becomes necessary for  $\mathbf{Q}$  to be single-valued since then,  $\hat{y} = \mathbf{Q}(\hat{u})$  and it must be the case that we can select a sequence  $\{w_{\rho_n}^*\}$  such that  $w_{\rho_n}^* \in \mathbf{P}_{\rho_n}(\hat{u})$  and  $w_{\rho_n}^* \to \hat{y}$  in V (by Theorem 2.18), and we find

$$J(\hat{y}, \hat{u}) \le J(y^*, u^*) \le \lim_{n \to \infty} J(y_{\rho_n}^*, u_{\rho_n}^*) \le J(\hat{y}, \hat{u}).$$

Because  $J(\hat{y},\hat{u})$  is the minimal value and hence is either independent of  $(\hat{y},\hat{u})$  or uniquely determined by  $(\hat{y},\hat{u})$ , the subsequence principle shows that  $J(y_{\rho}^*,u_{\rho}^*)\to J(\hat{y},\hat{u})$  (for the entire sequence). Furthermore, the above inequality shows that  $(y^*,u^*)$  is optimal and we get  $u_{\rho}^*\to u^*$  in H since we have weak convergence and convergence of the norm.

Regarding the assumption in this lemma that  $\mathbf{Q}$  is single valued, this is the case if, for example,  $\Phi$  is (globally) Lipschitz with Lipschitz constant strictly smaller than  $C_a/(C_a+C_b)$ , see the discussion around [4, Equation (21)]. An alternative condition for uniqueness for QVIs in a specific setting is given in [45].

Let us see how the results of this section change if we do not assume complete continuity of  $\Phi \colon V \to V$ .

**Remark 4.3.** (1) We can drop (19) from Theorem 4.1 in favour of the conditions in Theorem 2.19 as long as in the Gelfand triple regime (16) we assume  $U \hookrightarrow H$ . Examining the proof of Theorem 4.1, the feasibility of the limit of the infimising sequence follows exactly as in the proof of Theorem 2.19. The Cauchy estimate is not necessary. Weak lower semicontinuity of the norm allows us to retain the final line in the proof.

(2) If we drop (19) from Proposition 4.2 in favour of  $V \stackrel{c}{\hookrightarrow} H$  and the conditions in Theorem 2.19 as long as in the Gelfand triple regime (16) we assume  $U \hookrightarrow H$ , we would get  $y_{\rho}^* \rightharpoonup y^*$  in V (i.e., a weak convergence). To see this, we simply need to modify the proof to use Theorem 2.19 instead of Theorem 2.18. The compact embedding into H is needed to bound from above the term  $\limsup_{n\to\infty} J(w_{\rho_n}^*, \hat{u})$  by  $J(\hat{y}, \hat{u})$ .

# 5 Stationarity

In this section, we shall derive various forms of necessary conditions satisfied by optimal controls and states. Let us first formally define some concepts of stationarity which are motivated by analogous concepts from the VI case and also by the results that we shall obtain later.

Let  $(y,u) \in V \times H$  be a solution of the optimal control problem (2) where  $V \hookrightarrow H$  with V a reflexive Banach space and H a Hilbert space,  $U_{ad} \subset H$  is non-empty and weakly sequentially closed (in the context of the previous section, we have assumed  $U \equiv H$ ).

Inspired by the results we obtain in §5.2 in a general function space setting, we say that (y, u) is a *weak C-stationarity* point of (2) if there exists  $(p, \xi, \lambda) \in V \times V^* \times V^*$  such that

$$\begin{split} y + (\mathbf{I} - \Phi'(y))^*\lambda + A^*p &= y_d, \\ Ay - u + \xi &= 0, \\ \xi &\geq 0 \text{ in } V^*, \quad y \leq \Phi(y), \quad \langle \xi, y - \Phi(y) \rangle &= 0, \\ u &\in U_{ad} : (\nu u - p, u - v)_H \leq 0 \quad \forall v \in U_{ad}, \\ \langle \lambda, p \rangle &\geq 0. \end{split}$$

The function p is said to be the *adjoint state* and  $\lambda$  is the *Lagrange multiplier* associated to the adjoint state equation (the first equation above).

Let us now restrict the discussion to when  $H = L^2(\Omega)$  on a domain  $\Omega \subset \mathbb{R}^n$ . Certain sets associated to the lower-level QVI problem in (2) are important in stating the following stationarity conditions. Denoting  $\xi := u - Ay$  (see Proposition 2.1), let us formally define then the following sets:

$$\mathcal{A} := \{y = \Phi(y)\}$$
 is the *active* (or coincidence) set,  $\mathcal{I} := \{y < \Phi(y)\}$  is the *inactive* set,  $\mathcal{A}_s := \{\xi > 0\}$  is the *strongly active* set,  $\mathcal{B} := \{y = \Phi(y)\} \cap \{\xi = 0\}$  is the *biactive* set.

These definitions are merely heuristic due to the (in general) low regularity of  $\xi$ , see for example [69, §3 and Appendix A] or [33] for a rigorous approach to define these objects.

We say that  $(y,u) \in V \times H$  is a *C-stationarity* point of (2) if (y,u) is a solution of (2) and there exists  $(p,\xi,\lambda) \in V \times V^* \times V^*$  such that

$$y + (I - \Phi'(y))^* \lambda + A^* p = y_d,$$
 (40a)

$$Ay - u + \xi = 0, (40b)$$

$$\xi \ge 0 \text{ in } V^*, \quad y \le \Phi(y), \quad \langle \xi, y - \Phi(y) \rangle = 0,$$
 (40c)

$$u \in U_{ad} : (\nu u - p, u - v)_H \le 0 \quad \forall v \in U_{ad}, \tag{40d}$$

$$\langle \xi, p^+ \rangle = \langle \xi, p^- \rangle = 0 \tag{40e}$$

$$\langle \lambda, p \rangle \ge 0, \quad \langle \lambda, y - \Phi(y) \rangle = 0,$$
 (40f)

$$\langle \lambda, v \rangle = 0 \quad \forall v \in V : v = 0 \text{ a.e. on } \Omega \setminus \mathcal{I}.$$
 (40g)

Note that we use the condition (40e) in lieu of the more commonly seen condition p = 0 a.e. in  $\{\xi > 0\}$  due to the low regularity of  $\xi$ .

**Remark 5.1.** It is worth remarking that in certain works [65], rather than the inequality constraint in (40f), the stronger condition

$$\langle \lambda, \psi p \rangle \ge 0$$
 for all sufficiently smooth and non-negative  $\psi$  (41)

is required in order to satisfy C-stationarity; this is a direct analogy of the corresponding (element-wise) condition in the finite dimensional setting in [64]. We will also consider the obtainment of (41) in Proposition 5.13.

The condition (40g) is in practice difficult to check due to the fact that in general,  $\lambda$  possesses only the low  $V^*$  regularity. Therefore, one looks for a weaker concept. In the first instance, for an *almost C-stationarity* point, (40g) is replaced by

$$\langle \lambda, v \rangle = 0 \quad \forall v \in V : v = 0 \text{ a.e. on } \Omega \setminus \mathcal{I}, \ v|_{\mathcal{I}} \in H_0^1(\mathcal{I}).$$

More generally, an  $\mathcal{E}$ -almost C-stationarity point, the concept of which was introduced by Hintermüller and Kopacka in [36, 35], satisfies (40a)–(40f) but now (40g) is replaced with

$$\forall \tau > 0, \exists E^{\tau} \subset \mathcal{I} \text{ with } |\mathcal{I} \setminus E^{\tau}| \leq \tau : \langle \lambda, v \rangle = 0 \quad \forall v \in V : v = 0 \text{ a.e. on } \Omega \setminus E^{\tau}.$$

This is a condition that arises from an application of Egorov's theorem as we shall see later.

Now, in the other direction, a point which satisfies (40a)–(40c) and additionally

$$p\geq 0$$
 q.e. on  $\mathcal B$  and  $p=0$  q.e. on  $\mathcal A_s$ ,  $\langle \lambda,v \rangle \geq 0$   $\forall v\in V:v\geq 0$  q.e. on  $\mathcal B$  and  $v=0$  q.e. on  $\mathcal A_s$ ,

is called a *strong stationarity* point, which is typically the most stringent notion of stationarity possible and requires differentiability of the control-to-state map to be obtainable.

In the proceeding sections, we will show that there exist weak C-stationarity, ( $\mathcal{E}$ -almost) C-stationarity and strong stationarity points under various assumptions. We will, however, first start in §5.1 with the so-called *Bouligand stationarity* which is a primal condition and is defined below. It also requires differentiability of  $\mathbf{Q}$ .

## 5.1 Bouligand stationarity

In the case where  $\mathbf{Q}$  is directionally differentiable from the results of §3, we have the following Bouligand stationarity (or B-stationarity) characterisation of the optimal control, see [50, §5] and [51, Lemma 3.1] for the VI case. To start, define the radial cone of  $U_{ad}$  at  $u^*$  and the tangent cone respectively by

$$\mathcal{R}_{U_{ad}}(u^*) = \{h \in H : \exists s^* > 0 \text{ such that } u^* + sh \in U_{ad} \quad \forall s \in [0, s^*] \} \quad \text{and} \quad \mathcal{T}_{U_{ad}}(u^*) := \overline{\mathcal{R}_{U_{ad}}(u^*)}.$$

**Proposition 5.2** (Bouligand stationarity). Let  $U_{ad}$  be non-empty and  $(y^*, u^*)$  be a local minimiser of (2) and let the assumptions of Theorem 3.2 hold. Then

$$(\alpha_h, y^* - y_d)_H + \nu(u^*, h)_H \ge 0 \quad \forall h \in \mathcal{T}_{U_{ad}}(u^*),$$
 (42)

where  $\alpha_h$  is the directional derivative given uniquely through Theorem 3.2 as the solution of (27) with source h.

*Proof.* Take h in the radial cone of  $U_{ad}$  at  $u^*$  so that it is an admissible direction. Using this direction term, we define  $y_s$  as given by Theorem 3.2 after having initially selected  $y^* \in \mathbf{Q}(u^*)$ . This satisfies  $y_s = y^* + s\alpha_h + o(s)$  where  $\alpha_h$  is the directional derivative (uniquely determined thanks to Proposition 3.9) and o is a remainder term. It follows that  $(u^* + sh, y_s)$  can be made arbitrarily close to  $(u^*, y^*)$  if s is sufficiently small (since  $y_s - y^* = s\alpha_h + o(s)$  and the

 $<sup>^{16}</sup>$ These assumptions should be evaluated locally at  $y^*$ , of course.

right-hand side tends to zero in V). Hence, by definition of local minimiser, we have  $J(y_s, u^* + sh) \ge J(y^*, u^*)$  for s sufficiently small. Writing this inequality out, we get

$$0 \le \|y_s - y_d\|_H^2 + \nu \|u^* + sh\|_H^2 - \|y^* - y_d\|_H^2 - \nu \|u^*\|_H^2$$
  
=  $\|y_s\|_H^2 - \|y^*\|_H^2 + 2(y^* - y_s, y_d)_H + \nu s^2 \|h\|_H^2 + 2\nu s(u^*, h)_H.$ 

This leads to

$$0 \leq \|y^* + s\alpha_h + o(s)\|_H^2 - \|y^*\|_H^2 - 2(s\alpha_h + o(s), y_d)_H + \nu s^2 \|h\|_H^2 + 2\nu s(u^*, h)_H$$
$$= \|s\alpha_h + o(s)\|_H^2 + 2(s\alpha_h + o(s), y^* - y_d)_H + \nu s^2 \|h\|_H^2 + 2\nu s(u^*, h)_H$$
$$= s^2 \|\alpha_h + s^{-1}o(s)\|_H^2 + 2(s\alpha_h + o(s), y^* - y_d)_H + \nu s^2 \|h\|_H^2 + 2\nu s(u^*, h)_H.$$

Dividing by s and sending to zero, the above yields

$$0 \le 2(\alpha_h, y^* - y_d)_H + 2\nu(u^*, h)_H \quad \forall h \in \mathcal{R}_{U_{ad}}(u^*),$$

and by density and the continuity result of Proposition 3.12, also for  $h \in \mathcal{T}_{U_{ad}}(u^*)$ .

# 5.2 Weak C-stationarity

In this section we will show a type of weak C-stationarity for the optimal pair by passing to the limit in the stationarity system satisfied by the optimal pair of the PDE regularisation of the QVI. Recall the notations and framework of §2.3 and §4.1 where we studied the convergence of solutions of certain PDEs to a solution of the associated QVI and the associated optimal control problems. In this section, we again take

 $(y^*, u^*)$  to be an arbitrary local minimiser of (2).

In addition to the basic setup of Assumption 2.12, we need the following fundamental conditions related to  $\Phi$ , in which we also recall two assumptions that were stated earlier for the convenience of the reader.

#### Assumption 5.3. Assume that

there exists 
$$v_0 \in V$$
 such that  $v_0 \le \Phi(v)$  for all  $v \in V$ , (14)

$$\Phi \colon V \to V \text{ is completely continuous},$$
 (19)

there exists 
$$\epsilon > 0$$
 such that  $\Phi \colon V \to V$  is continuously Fréchet differentiable on  $B_{\epsilon}(y^*)$ , (43)  $\mathbf{Q}$  is single valued.

We also introduce the following invertibility assumptions; these are stated separately from above since they will come in use later in another section. Note that these types of conditions are also needed in [72].

#### Assumption 5.4. Assume that

$$(I - \Phi'(z)): V \to V \text{ is invertible for } z \in B_{\epsilon}(y^*),$$
 (44)

$$A(I - \Phi'(z))^{-1}: V \to V^*$$
 is uniformly bounded and uniformly coercive in  $z \in B_{\epsilon}(y^*)$ . (45)

The main result of this section is the following theorem which shows that local minimisers are weak C-stationarity points.

Theorem 5.5 (Weak C-stationarity). Suppose that

$$U_{ad}$$
 is non-empty, closed and convex and  $V \stackrel{c}{\hookrightarrow} H \hookrightarrow V^*$  is a Gelfand triple. (46)

In addition to Assumptions 2.12, 5.3 and 5.4, suppose that  $m_{\rho}$  satisfies along with (10)–(12) the conditions

$$m_o: H \to V^* \text{ is continuous}$$
 (47)

$$m_{\rho} \colon V \to V^*$$
 is continuously Fréchet differentiable. (48)

Then there exist multipliers  $(p^*, \xi^*, \lambda^*) \in V \times V^* \times V^*$  satisfying the weak C-stationarity system

$$y^* + (I - \Phi'(y^*))^* \lambda^* + A^* p^* = y_d, \tag{49a}$$

$$Ay^* - u^* + \xi^* = 0, (49b)$$

$$\xi^* \ge 0 \text{ in } V^*, \quad y^* \le \Phi(y^*), \quad \langle \xi^*, y^* - \Phi(y^*) \rangle = 0,$$
 (49c)

$$u^* \in U_{ad} : (\nu u^* - p^*, u^* - v)_H \le 0 \quad \forall v \in U_{ad},$$
 (49d)

$$\langle \lambda^*, p^* \rangle \ge 0. \tag{49e}$$

Here, we have assumed the existence of  $C^1$  maps  $m_\rho$  — this will have to be verified on a case-by-case basis (we leave the possibility of being able to define such maps satisfying the conditions (10)–(12) and (47)–(48) in the general setting to the interested reader, who may find [31] useful for this purpose). However, let us note that in the most common case of interest where the function spaces involve functions over domains in  $\mathbb{R}^n$  with the usual ordering, it is usually possible to construct sufficiently smooth  $m_\rho$ , see for example §5.3.

- **Remark 5.6.** (i) We assumed the complete continuity (19) to utilise the strong convergence result of Theorem 2.18. It would be interesting to see how the calculations below can be adapted in the case where (we do not have complete continuity and) we only have weak convergence from Theorem 2.19.
  - (ii) Due to the Gelfand triple setup and complete continuity of  $\Phi$  here, we find from (47) that the complete continuity of  $m_{\rho}$  condition (15a) is satisfied.
  - (iii) The meaning of (45) is that for all  $z \in B_{\epsilon}(y^*)$ , the operator  $A(I \Phi'(z))^{-1}$  has a boundedness constant  $C'_b$  and a coercivity constant  $C'_a$  both of which are independent of z. A consequence is that

$$(I - \Phi'(z))^{-1} : V \to V \text{ is bounded uniformly for all } z \in B_{\epsilon}(y^*).$$
 (50)

Since  $\Phi$  is  $C^1$ , we automatically have that  $(I - \Phi'(z))^{-1}$  is bounded; (50) clarifies that the bound is uniform.

Let us proceed with proving this result.

#### 5.2.1 Stationarity for the penalised optimal control problem

Recall the penalised problem (38) that approximates (2):

$$\min_{u \in U_{ad}} J(y_{\rho}, u) \quad \text{such that} \quad Ay_{\rho} + \frac{1}{\rho} m_{\rho} (y_{\rho} - \Phi(y_{\rho})) = u. \tag{38}$$

Under Assumption 5.3 and (15a), Proposition 4.2 is applicable. For the moment and for purposes of a simpler exposition, let us assume that

$$(y^*, u^*)$$
 is the optimal point of (2) given in Proposition 4.2 (51)

(we will discard this later on). Via the proposition, we obtain the existence of minimisers  $(y_o^*, u_o^*)$  of (38) such that

$$(y_\rho^*,u_\rho^*)\to (y^*,u^*) \text{ in } V\times H.$$

Thus, for any  $\epsilon > 0$ , we can find a  $\rho_0$  such that  $\rho \leq \rho_0$  implies

$$y_{\rho}^* \in B_{\epsilon}(y^*)$$

(this is why it has been possible to formulate most assumptions on  $\Phi$  only locally). To derive stationarity conditions for the penalised problem (38), we check the Zowe–Kurcyusz constraint qualification [74] (see also the Robinson condition [60]). To do so, we make the necessary surjectivity assumption (52) below regarding existence for the linearised equation — we discuss instances where it holds in Remark 5.8.

**Lemma 5.7.** Assume (43), (48), (46) and suppose that

$$\forall \rho \le \rho_0, \ \forall f \in V^*, \ \exists z \in V : Az + \frac{1}{\rho} m_\rho' (y_\rho^* - \Phi(y_\rho^*)) (I - \Phi'(y_\rho^*)) (z) = f. \tag{52}$$

Then, for such  $\rho$  and any optimal point  $(y_{\rho}^*, u_{\rho}^*)$  of (38), there exists  $p_{\rho}^* \in V$  such that

$$A^* p_{\rho}^* + \frac{1}{\rho} (\mathbf{I} - \Phi'(y_{\rho}^*))^* m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^* = y_d - y_{\rho}^*,$$

$$(\nu u_{\rho}^* - p_{\rho}^*, u_{\rho}^* - v)_H \le 0 \qquad \forall v \in U_{ad}.$$
(53)

*Proof.* We introduce the following notation:

$$X := V \times H, \quad g(x) = g(y, u) := Ay + \frac{1}{\rho} m_{\rho} (y - \Phi(y)) - u,$$

$$x_{\rho} = (y_{\rho}^*, u_{\rho}^*), \quad C(x_{\rho}) := \{ k(v - y_{\rho}^*, h - u_{\rho}^*) : v \in V, h \in U_{ad}, k \ge 0 \}.$$

The map  $g\colon X\to V^*$ , being a composition of  $C^1$  maps, is continuously Fréchet differentiable at  $x_\rho$  and we must check that  $g'(x_\rho)C(x_\rho)=V^*$ , but since  $\tilde C:=V\times\{0\}\subset C(x_\rho)$ , it suffices to verify  $g'(x_\rho)\tilde C=V^*$ . Observing that

$$g'(x_{\rho})(y,0) = Ay + \frac{1}{\rho}m'_{\rho}(y_{\rho}^* - \Phi(y_{\rho}^*))(y - \Phi'(y_{\rho}^*)(y)),$$

it follows that we need existence for the PDE in (52) and this is guaranteed by assumption for  $\rho$  sufficiently small. Calculating the adjoint  $g'(x_{\rho})^*: V \to X^*$  of g' via

$$\langle g'(x_{\rho})(y,u),v\rangle = \langle Ay,v\rangle + \frac{1}{\rho} \langle m'_{\rho}(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))(y - \Phi'(y_{\rho}^{*})(y)),v\rangle - (u,v)$$

$$= \langle y, A^{*}v\rangle + \frac{1}{\rho} \langle y, (I - \Phi'(y_{\rho}^{*}))^{*}m'_{\rho}(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{*}v\rangle - (v,u),$$

we find

$$g'(x_{\rho})^{*}(v) = \left(A^{*}v + \frac{1}{\rho}(I - \Phi'(y_{\rho}^{*}))^{*}m'_{\rho}(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{*}v, -v\right).$$

Applying, for example, [68, Theorem 6.3], we get the existence of  $p_{\rho}^* \in V$  such that  $J'(x_{\rho}) - g'(x_{\rho})^* p_{\rho}^* \in C(x_{\rho})^{\circ}$ , i.e., for all  $k \geq 0$ ,

$$\langle y_{\rho}^{*} - y_{d} + \frac{1}{\rho} (I - \Phi'(y_{\rho}^{*}))^{*} m_{\rho}' (y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{*} p_{\rho}^{*} + A^{*} p_{\rho}^{*}, k(c_{1} - y_{\rho}^{*}) \rangle \ge 0 \qquad \forall c_{1} \in V,$$

$$(\nu u_{\rho}^{*} - p_{\rho}^{*}, k(c_{2} - u_{\rho}^{*}))_{H} \ge 0 \qquad \forall c_{2} \in U_{ad}$$

As  $c_1 \in V$  can be chosen arbitrarily, we find the stated result.

**Remark 5.8.** The conditions of Assumption 5.4 are clearly sufficient to guarantee the surjectivity condition (52); and in fact (45) can be replaced with asking for  $A(I - \Phi'(z))^{-1} : V \to V^*$  to be coercive for all  $z \in B_{\epsilon}(y^*)$ . Indeed, first observe that the bounded inverse theorem guarantees that  $A(I - \Phi'(z))^{-1} : V \to V^*$  is bounded for  $z \in B_{\epsilon}(y^*)$ . Now, the equation

$$A(I - \Phi'(y_{\rho}^*))^{-1}w + \frac{1}{\rho}m_{\rho}'(y_{\rho}^* - \Phi(y_{\rho}^*))w = f$$

has a unique solution  $w \in V$  by the Lax–Milgram theorem, leading to existence of  $z := (I - \Phi'(y_{\rho}^*))^{-1}w \in V$  satisfying the equation in (52).

#### **5.2.2** Passage to the limit $\rho \rightarrow 0$

Now the objective is to pass to the limit in (53) as  $\rho \to 0$  for which we shall need some technical results.

**Lemma 5.9.** Under Assumption 5.4, if  $z_n \to z$  and  $q_n \rightharpoonup q$  in V with  $z_n, z \in B_{\epsilon}(y^*)$ , then

$$(I - \Phi'(z_n))^{-1}q_n \rightharpoonup (I - \Phi'(z))^{-1}q \quad in V,$$
 (54)

$$\langle A(\mathbf{I} - \Phi'(z))^{-1} q, q \rangle \le \liminf_{n \to \infty} \langle A(\mathbf{I} - \Phi'(z_n))^{-1} q_n, q_n \rangle.$$
 (55)

The convergence in (54) is strong if  $q_n \to q$  in V.

In order to not disturb the flow of the paper, the proof of this lemma has been placed in Appendix A. As an immediate corollary to Lemma 5.9, for sequences  $w_{\rho} \to w$  and  $q_{\rho} \rightharpoonup q$  in V, we have

$$\lim_{n \to \infty} (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1} w_{\rho} = (\mathbf{I} - \Phi'(y^*))^{-1} w \text{ in } V,$$
(56)

$$(y^*, (\mathbf{I} - \Phi'(y^*))^{-1}q)_H \le \liminf_{n \to \infty} (y_\rho^*, (\mathbf{I} - \Phi'(y_\rho^*))^{-1}q_\rho)_H, \tag{57}$$

$$(y_d, (I - \Phi'(y^*))^{-1}q)_H \ge \limsup_{n \to \infty} (y_d, (I - \Phi'(y_\rho^*))^{-1}q_\rho)_H.$$
(58)

We are now ready to conclude.

*Proof of Theorem 5.5.* First, note that Proposition 2.1 directly gives (49c). Assumption 5.4 implies the surjectivity condition (52) (see Remark 5.8), therefore the stationarity conditions in (53) for the penalised problem are available. Now, the weak form of the equation for  $p_a^*$  is

$$\langle A^* p_{\rho}^*, \varphi \rangle + \frac{1}{\rho} \langle m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*)) \varphi \rangle = (y_d - y_{\rho}^*, \varphi)_H \qquad \forall \varphi \in V.$$

By defining  $v := (I - \Phi'(y_{\rho}^*))\varphi$ , thanks to the invertibility assumption (44), this can be transformed to

$$\langle A(\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}v, p_{\rho}^* \rangle + \frac{1}{\rho} \langle m_{\rho}'(y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^*, v \rangle = (y_d - y_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}v)_H \qquad \forall v \in V.$$

Selecting  $v=p_{\rho}^*$ , using the coercivity (45), the monotonicity of  $m_{\rho}$  (which implies that  $\langle m_{\rho}'(v)(h), h \rangle \geq 0$  for all  $v,h\in V$ ), Young's inequality with  $\gamma>0$  and the uniform boundedness of  $(I-\Phi'(y_{\rho}^*))^{-1}$  assured by (50), we obtain

$$C_a' \|p_{\rho}^*\|_V^2 \le C_{\gamma} \|y_d - y_{\rho}^*\|_H^2 + \gamma \|p_{\rho}^*\|_V^2$$

Selecting  $\gamma$  sufficiently small so that the right-most term is absorbed onto the left, we obtain a bound on  $\{p_{\rho}^*\}$  independent of  $\rho$ . This gives rise to the convergence (for a subsequence that has been relabelled)

$$p_{\rho}^* \rightharpoonup p^* \quad \text{in } V.$$

Define

$$\begin{split} \lambda_{\rho}^* &:= \frac{1}{\rho} m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^*, \\ \mu_{\rho}^* &:= \frac{1}{\rho} (\mathbf{I} - \Phi'(y_{\rho}^*))^* m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^* = y_d - y_{\rho}^* - A^* p_{\rho}^*, \\ \xi_{\rho}^* &:= \frac{1}{\rho} m_{\rho} (y_{\rho}^* - \Phi(y_{\rho}^*)) = u_{\rho}^* - A y_{\rho}^*, \end{split}$$

the latter two of which, since their right-hand sides converge, satisfy the following convergences both in  $V^*$ :

$$\mu_{\rho}^* \rightharpoonup \mu^* := y_d - y^* - A^* p^* \quad \text{and} \quad \xi_{\rho}^* \to \xi^* := u^* - A y^*.$$
 (59)

Again using monotonicity of  $m_o$ ,

$$\langle \mu_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1} p_{\rho}^* \rangle = \frac{1}{\rho} \langle m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^*, p_{\rho}^* \rangle \ge 0,$$

and taking the limit superior of this, recalling the definition of  $\mu^*$ , we obtain

$$\begin{split} 0 &= \limsup_{\rho \to 0} \langle y_d, (\mathbf{I} - \Phi'(y_\rho^*))^{-1} p_\rho^* \rangle - \liminf_{\rho \to 0} \langle y_\rho^*, (\mathbf{I} - \Phi'(y_\rho^*))^{-1} p_\rho^* \rangle - \liminf_{\rho \to 0} \langle A(\mathbf{I} - \Phi'(y_\rho^*))^{-1} p_\rho^*, p_\rho^* \rangle \\ &\leq \langle y_d - y^*, (\mathbf{I} - \Phi'(y^*))^{-1} p^* \rangle - \langle A(\mathbf{I} - \Phi'(y^*))^{-1} p^*, p^* \rangle \\ &\qquad \qquad \text{(using the weak semicontinuity results (55), (57) and (58))} \\ &= \langle \mu^*, (\mathbf{I} - \Phi'(y^*))^{-1} p^* \rangle. \end{split}$$

Finally, writing the VI relating  $u_o^*$  and  $p_o^*$  in (53) as

$$(\nu u_{\rho}^*, u_{\rho}^* - v)_H - \langle u_{\rho}^* - v, p_{\rho}^* \rangle \le 0 \qquad \forall v \in U_{ad},$$

using the strong convergence of  $u_{\rho}^*$  in H (and hence also in  $V^*$ ) and the weak convergence of  $p_{\rho}^*$  in V, we can pass to the limit.

Collecting the results (and recalling that the inverses and adjoints of bounded linear operators commute), we have shown the satisfaction of (49b)–(49d) and

$$y^* + \mu^* + A^* p^* = y_d,$$
$$\langle (I - \Phi'(y^*)^*)^{-1} \mu^*, p^* \rangle \ge 0,$$

Setting  $\lambda^* := (I - \Phi'(y^*)^*)^{-1} \mu^*$  we get the system (49).

Thus far, we have only shown the existence of a stationarity point and not that every local minimiser is such a point since we assumed (51). Suppose now that  $(y^*, u^*)$  is an arbitrary local minimiser (instead of (51)) as claimed in the statement of the theorem. Denote by  $\gamma$  the radius such that  $u^*$  is the minimiser on  $U_{ad} \cap B^H_{\gamma}(u^*)$  (the latter object is the closed ball in H of radius  $\gamma$  with centre  $u^*$ ). Consider for  $\bar{J}(y_\rho, u) := J(y_\rho, u) + \|u - u^*\|_H^2$  the problem

$$\min_{u \in U_{ad} \cap B_{\gamma}^{H}(u^{*})} \bar{J}(y_{\rho}, u) \quad \text{such that} \quad Ay_{\rho} + \frac{1}{\rho} m_{\rho}(y_{\rho} - \Phi(y_{\rho})) = u. \tag{60}$$

Denote by  $(\bar{y}_{\rho}, \bar{u}_{\rho})$  a minimiser of this problem. It follows from  $\bar{J}(\bar{y}_{\rho}, \bar{u}_{\rho}) \leq \bar{J}(y_{\rho}(u^*), u^*)$  and  $\mathbf{P}_{\rho}(u^*) \ni y_{\rho}(u^*) \to y^*$  that

$$\limsup_{\rho \to 0} \bar{J}(\bar{y}_{\rho}, \bar{u}_{\rho}) \le J(y^*, u^*).$$

On the other hand, from uniform bounds, we obtain the existence of  $\hat{u}$  such that  $\bar{u}_{\rho} \rightharpoonup \hat{u}$  in H and  $\bar{y}_{\rho} \to \mathbf{Q}(\hat{u}) =: \hat{y}$  in V, giving (by the identity  $\lim \sup(a_n) + \lim \inf(b_n) \le \lim \sup(a_n + b_n)$  and using weak lower semicontinuity)

$$\limsup_{\rho \to 0} \bar{J}(\bar{y}_{\rho}, \bar{u}_{\rho}) \ge J(\hat{y}, \hat{u}) + \limsup_{\rho \to 0} \|\bar{u}_{\rho} - u^*\|_{H}^{2} \ge J(y^*, u^*) + \limsup_{\rho \to 0} \|\bar{u}_{\rho} - u^*\|_{H}^{2}$$

with the last inequality because  $(y^*, u^*)$  is a local minimiser and  $\hat{u} \in B_{\gamma}^H(u^*)$ . Combining these two inequalities shows that  $\hat{u} = u^*$  and  $\bar{u}_{\rho} \to u^*$  in H. The latter fact implies that for  $\rho$  sufficiently small,  $\bar{u}_{\rho} \in B_{\gamma}^H(u^*)$  automatically and hence the feasible set in (60) can be taken to be just  $U_{ad}$ . For such  $\rho$  (assuming of course that the local conditions in Assumptions 5.3 and 5.4 hold around  $y^*$ ), the same arguments as above can be used to derive stationarity conditions for (60) and in passing to the limit in those conditions, we will find that  $(y^*, u^*)$  satisfies the same conditions as above.  $\square$ 

The proof reveals that the stationarity point satisfying (51) can be characterised as a limit of the following subsequences (which we have relabelled):

$$\begin{split} y_{\rho}^* \to y^* &\quad \text{in } V, \\ u_{\rho}^* \to u^* &\quad \text{in } H, \\ p_{\rho}^* &\rightharpoonup p^* &\quad \text{in } V, \\ \rho^{-1} m_{\rho} (y_{\rho}^* - \Phi(y_{\rho}^*)) \to \xi^* &\quad \text{in } V^*, \\ \rho^{-1} m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*)) p_{\rho}^* &\rightharpoonup \lambda^* &\quad \text{in } V^*, \end{split}$$

where  $(y_{\rho}^*, u_{\rho}^*, p_{\rho}^*)$  are as in Lemma 5.7.

#### 5.3 $\mathcal{E}$ -almost C-stationarity

We specialise to the case where H is an  $L^2$  space on a bounded domain with box constraints, which allows us to improve the weak C-stationarity system.

**Assumption 5.10.** Let  $\Omega \subset \mathbb{R}^n$  be a bounded Lipschitz domain, set  $H := L^2(\Omega)$  and take  $V \in \{H^1(\Omega), H^1_0(\Omega)\}$  and assume the Gelfand triple  $(V, H, V^*)$  structure. Finally, we take  $U_{ad}$  to be of the box constraint type

$$U_{ad} = \{ u \in H : u_a \le u \le u_b \text{ a.e. in } \Omega \}$$

$$\tag{61}$$

for given functions  $u_a, u_b \in H$ .

The assumption can be generalised, see Remark 5.12.

As before, we denote by

 $(y^*, u^*)$  an arbitrary local minimiser of (2).

**Theorem 5.11** ( $\mathcal{E}$ -almost C-stationarity). Let Assumptions 5.3, 5.4 and 5.10 hold. Then there exist multipliers  $(p^*, \xi^*, \lambda^*) \in V \times V^* \times V^*$  satisfying the  $\mathcal{E}$ -almost C-stationarity system

$$y^* + (I - \Phi'(y^*))^* \lambda^* + A^* p^* = y_d, \tag{62a}$$

$$Ay^* - u^* + \xi = 0, (62b)$$

$$\xi^* \ge 0 \text{ in } V^*, \quad y^* \le \Phi(y^*), \quad \langle \xi^*, y^* - \Phi(y^*) \rangle = 0,$$
 (62c)

$$u^* \in U_{ad} : (\nu u^* - p^*, u^* - v) \le 0 \quad \forall v \in U_{ad},$$
 (62d)

$$\langle \xi^*, (p^*)^+ \rangle = \langle \xi^*, (p^*)^- \rangle = 0 \tag{62e}$$

$$\langle \lambda^*, p^* \rangle \ge 0, \quad \langle \lambda^*, y^* - \Phi(y^*) \rangle = 0,$$
 (62f)

$$\forall \tau > 0, \exists E^{\tau} \subset \mathcal{I} \text{ with } |\mathcal{I} \setminus E^{\tau}| \leq \tau : \langle \lambda^*, v \rangle = 0 \quad \forall v \in V : v = 0 \text{ a.e. on } \Omega \setminus E^{\tau}. \tag{62g}$$

In addition, if  $u_a, u_b \in V$  then the optimal control has the regularity  $u^* \in V$ .

To prove the theorem, we choose a particular  $m_{\rho}$  (that appeared in the work of Hintermüller and Kopacka [36] for VIs), namely the superposition operator defined through the real-valued function

$$m_{\rho}(r) \equiv \max_{\epsilon(\rho)}(0, \cdot) := \begin{cases} 0 & : r \le 0\\ \frac{r^2}{2\epsilon} & : 0 < r < \epsilon\\ r - \frac{\epsilon}{2} & : r \ge \epsilon; \end{cases}$$

$$(63)$$

here,  $\epsilon = \epsilon(\rho) > 0$  is chosen such that  $\{\epsilon(\rho)\}$  is bounded. The parameter  $\epsilon$  is a smoothing parameter utilised for ensuring differentiability at 0. By [22, Lemmas 2.83, 2.87, 2.88, 2.90] and the fact that  $m_{\rho} \in C^{1}(\mathbb{R})$  with  $m_{\rho}' \in [0, 1]$ , we obtain relevant lattice properties for the spaces involved and differentiability properties for  $m_{\rho}$ . That  $m_{\rho}$  satisfies (10), (11) and (47) is clear. Let us check condition (12). Since  $\{\epsilon(\rho)\}$  is bounded, we have (for a subsequence that we relabelled)  $\epsilon(\rho) \to \bar{\epsilon}$  for some  $\bar{\epsilon} \geq 0$  and we get

$$\left\| \max_{\bar{\epsilon}}(0, z) - \max_{\epsilon(\rho)}(0, z_{\rho}) \right\|_{V^*} \leq C \left( \left\| \max_{\bar{\epsilon}}(0, z) - \max_{\bar{\epsilon}}(0, z_{\rho}) \right\|_{H} + \left\| \max_{\bar{\epsilon}}(0, z_{\rho}) - \max_{\epsilon(\rho)}(0, z_{\rho}) \right\|_{H} \right)$$

$$\leq C \left( \left\| z - z_{\rho} \right\|_{H} + \frac{3}{2} |\bar{\epsilon} - \epsilon(\rho)| \right)$$

$$\to 0$$

with the final inequality due to Lipschitz properties given in [36, Lemma 2.1 (v), (vi)] and the convergence due to the compact embedding  $V \stackrel{c}{\hookrightarrow} H$ . Hence we find  $z \leq 0$ . Finally, by the regularity of  $m_{\rho}$  (which has a bounded derivative) we have that  $m_{\rho} \colon H^1(\Omega) \to H$  is  $C^1$  (see, e.g. [21, Proposition 4]), thus we have (48). This shows that  $m_{\rho}$  is a valid choice.

**Remark 5.12.** Assumption 5.10 can be generalised as follows. Let  $\Omega \subset \mathbb{R}^n$  be a bounded Lipschitz domain, set  $H := L^2(\Omega)$  and take V to be a separable Hilbert space with  $V \stackrel{c}{\hookrightarrow} H$  and  $(V, H, V^*)$  a Gelfand triple. We assume that V is such that  $(\cdot)^+ : V \to V$  is continuous and that the superposition operator  $m_\rho$  takes V into H with  $m_\rho : V \to H$  being  $C^1$ .

The requirement for the Nemytskii operator to be Fréchet differentiable is in general a delicate issue.

*Proof of Theorem 5.11.* Elements of the proof are similar to that of [36, Theorem 3.4] but the more complicated problem structure in this paper requires additional work.

- 1. Weak C-stationarity. Observing that Assumption 5.10 implies (46), (47) and (48) (as discussed above), we have the weak C-stationarity result of Theorem 5.5 immediately at hand.
- 2. Regularity of optimal control. Owing to the characterisation of the VI relating  $u_{\rho}^*$  and  $p_{\rho}^*$  given in [42, §II.3], thanks to the strong convergence in H of  $p_{\rho}^*$  and continuity of  $(\cdot)^+$ :  $H \to H$ , we find that

$$u_{\rho}^* = \frac{1}{\nu} p_{\rho}^* + \left( u_a - \frac{p_{\rho}^*}{\nu} \right)^+ - \left( \frac{p_{\rho}^*}{\nu} - u_b \right)^+ \to \frac{1}{\nu} p^* + \left( u_a - \frac{p^*}{\nu} \right)^+ - \left( \frac{p^*}{\nu} - u_b \right)^+ = u^*.$$

It follows that  $u^* \in V$  if  $u_a$  and  $u_b$  belong to V.

3. Orthogonality condition. For the condition on  $y^* - \Phi(y^*)$  in (62f), observe that since  $m'_o$  vanishes on  $(-\infty, 0]$ ,

$$\langle \mu_{\rho}^*, (I - \Phi'(y_{\rho}^*))^{-1}(y_{\rho}^* - \Phi(y_{\rho}^*))^{-} \rangle = \frac{1}{\rho} \int_{\Omega} m_{\rho}' (y_{\rho}^* - \Phi(y_{\rho}^*))^* p_{\rho}^* (y_{\rho}^* - \Phi(y_{\rho}^*))^{-} = 0,$$

which, due to the continuity of  $(\cdot)^-: V \to V$  and the joint sequential continuity result of (56) implies that

$$\langle \mu^*, (\mathbf{I} - \Phi'(y^*))^{-1} (y^* - \Phi(y^*))^{-} \rangle = 0,$$

and since  $y^* \leq \Phi(y^*)$ , the negative part above can be dropped.

4. *E-almost statement.* Since  $y_{\rho}^* \to y^*$  in  $V, y_{\rho}^* - \Phi(y_{\rho}^*) \to y^* - \Phi(y^*)$  pointwise a.e. in  $\Omega$  for a subsequence that we do not relabel. Take  $x \in \Omega$  such that  $y^*(x) - \Phi(y^*)(x) < 0$ , then there exists a  $\hat{\rho} = \hat{\rho}(x)$  such that if  $\rho \leq \hat{\rho}$ , then

$$y_{\rho}(x) - \Phi(y_{\rho})(x) \le \frac{1}{2}(y^{*}(x) - \Phi(y^{*})(x)) < 0$$

and hence  $\rho^{-1}m_{\rho}'(y_{\rho}(x)-\Phi(y_{\rho})(x))=0$  for  $\rho\leq\hat{\rho}$ . That is,  $\rho^{-1}m_{\rho}'(y_{\rho}(x)-\Phi(y_{\rho})(x))\to 0$  pointwise a.e. on  $\{y^*<\Phi(y^*)\}$  and by Egorov's theorem, for every  $\tau>0$ , there exists  $B^{\tau}\subset\{y^*<\Phi(y^*)\}$  with  $|B^{\tau}|<\tau$  such that this convergence also holds uniformly on  $\{y^*<\Phi(y^*)\}\setminus B^{\tau}$ .

Take  $v \in V$  with v = 0 a.e. on  $\{y^* = \Phi(y^*)\} \cup B^{\tau}$ . By the uniform convergence, for any  $\gamma > 0$ , there exists  $\bar{\rho}$  such that if  $\rho \leq \bar{\rho}$ ,

$$\left| \langle \mu_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1} v \rangle \right| = \left| \int_{\{y^* < \Phi(y^*)\} \cap (B^{\tau})^c} \frac{1}{\rho} m_{\rho}' (y_{\rho} - \Phi(y_{\rho})) p_{\rho}^* v \right| \le \gamma \left\| p_{\rho}^* v \right\|_{L^1(\Omega)}.$$

The norm on the right-hand side is bounded uniformly and the left-hand side converges to  $|\langle \mu^*, (\mathbf{I} - \Phi'(y^*))^{-1}v \rangle|$  (thanks to  $\mu_\rho^* \rightharpoonup \mu^*$  in  $V^*$  from (59) and the strong convergence of  $(\mathbf{I} - \Phi'(y_\rho^*))^{-1}v$  in V given by (56)), thus giving

$$\left| \langle \mu^*, (\mathbf{I} - \Phi'(y^*))^{-1} v \rangle \right| \le C \gamma$$

for a constant C > 0. Since this holds for every  $\gamma$ , we obtain (62g) (simply set  $E^{\tau} := \mathcal{I} \setminus B^{\tau}$ ).

5. Relation between  $\xi^*$  and  $p^*$ . In order to show the remaining statement (62e), let us introduce the sets

$$M_1(\rho) := \{ 0 \le y_{\rho}^* - \Phi(y_{\rho}^*) < \epsilon \}$$
 and  $M_2(\rho) := \{ y_{\rho}^* - \Phi(y_{\rho}^*) \ge \epsilon \}.$ 

Since  $\langle \xi_{\rho}^*, y_{\rho}^* - \Phi(y_{\rho}^*) \rangle \to \langle \xi^*, y - \Phi(y) \rangle = 0$ , we find

$$(\xi_{\rho}^{*}, y_{\rho}^{*} - \Phi(y_{\rho}^{*})) = \frac{1}{\rho} \int_{\Omega} m_{\rho} (y_{\rho}^{*} - \Phi(y_{\rho}^{*})) (y_{\rho}^{*} - \Phi(y_{\rho}^{*}))$$

$$= \frac{1}{\rho} \int_{M_{1}(\rho)} \frac{(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{3}}{2\epsilon} + \frac{1}{\rho} \int_{M_{2}(\rho)} \left( y_{\rho}^{*} - \Phi(y_{\rho}^{*}) - \frac{\epsilon}{2} \right) (y_{\rho}^{*} - \Phi(y_{\rho}^{*}))$$

$$\to 0,$$
(64)

and as both integrands in (64) are non-negative, each integral must individually converge to zero too. Hence

$$\left\| \frac{\chi_{M_1(\rho)}(y_{\rho}^* - \Phi(y_{\rho}^*))^{\frac{3}{2}}}{\sqrt{\rho\epsilon}} \right\| \to 0 \quad \text{and} \quad \left\| \frac{\chi_{M_2(\rho)}(y_{\rho}^* - \Phi(y_{\rho}^*) - \frac{\epsilon}{2})}{\sqrt{\rho}} \right\| \to 0, \tag{65}$$

where for the second convergence we used the fact that  $y_{\rho}^* - \Phi(y_{\rho}^*) \ge y_{\rho}^* - \Phi(y_{\rho}^*) - \epsilon/2 \ge 0$ . We calculate

$$\langle \xi_{\rho}^{*}, p_{\rho}^{*} \rangle = \frac{1}{\rho} \int_{M_{1}(\rho)} \frac{(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{2}}{2\epsilon} p_{\rho}^{*} + \frac{1}{\rho} \int_{M_{2}(\rho)} \left( y_{\rho}^{*} - \Phi(y_{\rho}^{*}) - \frac{\epsilon}{2} \right) p_{\rho}^{*}$$

$$= \frac{1}{2} \int_{\Omega} \chi_{M_{1}(\rho)} \frac{(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{3/2}}{\sqrt{\rho \epsilon}} \frac{(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{1/2}}{\sqrt{\rho \epsilon}} \chi_{M_{1}(\rho)} p_{\rho}^{*} + \int_{\Omega} \frac{\chi_{M_{2}(\rho)} \left( y_{\rho}^{*} - \Phi(y_{\rho}^{*}) - \frac{\epsilon}{2} \right) \chi_{M_{2}(\rho)} p_{\rho}^{*}}{\sqrt{\rho}}$$

$$\leq \frac{1}{2} \left\| \chi_{M_{1}(\rho)} \frac{(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{3/2}}{\sqrt{\rho \epsilon}} \right\| \left\| \frac{(y_{\rho}^{*} - \Phi(y_{\rho}^{*}))^{1/2}}{\sqrt{\rho \epsilon}} \chi_{M_{1}(\rho)} p_{\rho}^{*} \right\| + \left\| \frac{\chi_{M_{2}(\rho)} \left( y_{\rho}^{*} - \Phi(y_{\rho}^{*}) - \frac{\epsilon}{2} \right)}{\sqrt{\rho}} \right\| \left\| \frac{\chi_{M_{2}(\rho)} p_{\rho}^{*}}{\sqrt{\rho}} \right\|. \tag{66}$$

Now, using (65), the first factor in each term above converges to zero and hence the above right-hand side will converge to zero if we are able to show that the second factor in each term remains bounded. Since  $\mu_{\rho}^*$  and  $(I - \Phi'(y_{\rho}^*))^{-1}p_{\rho}^*$  are bounded (the latter due to (50)), so is their duality product, and therefore

$$C \ge |\langle \mu_{\rho}^{*}, (\mathbf{I} - \Phi'(y_{\rho}^{*}))^{-1} p_{\rho}^{*} \rangle|$$

$$= \frac{1}{\rho} \left| \int_{\Omega} m_{\rho}' (y_{\rho}^{*} - \Phi(y_{\rho}^{*})) (p_{\rho}^{*})^{2} \right|$$

$$= \frac{1}{\rho} \left| \int_{M_{1}(\rho)} \frac{y_{\rho}^{*} - \Phi(y_{\rho}^{*})}{\epsilon} (p_{\rho}^{*})^{2} + \int_{M_{2}(\rho)} (p_{\rho}^{*})^{2} \right|$$

$$= \frac{1}{\rho} \int_{\Omega} \chi_{M_{1}(\rho)} \frac{y_{\rho}^{*} - \Phi(y_{\rho}^{*})}{\epsilon} (p_{\rho}^{*})^{2} + \frac{1}{\rho} \int_{\Omega} \chi_{M_{2}(\rho)} (p_{\rho}^{*})^{2}.$$

Both of the terms on the right-hand side are individually bounded uniformly in  $\rho$  as the integrands are non-negative. This fact then implies from (66) that

$$\langle \xi^*, p^* \rangle = 0.$$

Replacing  $p_{\rho}^*$  by  $(p_{\rho}^*)^+$  in (66) and in the above calculation, we also obtain in the same way (utilising the fact that  $v_n \rightharpoonup v$  in V implies that  $v_n^+ \rightharpoonup v^+$  in V)

$$\langle \xi^*, (p^*)^+ \rangle = 0.$$

Conclusion. Finally, setting  $\lambda^* := (I - \Phi'(y^*)^*)^{-1}\mu^*$ , we have shown the desired system (62).

We conclude this section by showing that the alternative (stronger) condition (41) occasionally used in literature for defining a C-stationarity point can be achieved under additional assumptions.

**Proposition 5.13** (Satisfaction of alternative criterion in C-stationarity). For  $q_{\rho} \rightharpoonup q$  in V, under the conditions of Theorem 5.11 and

$$\liminf_{n\to\infty} \langle A^* q_{\rho}, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi q_{\rho}) \rangle \ge \langle A^* q, (\mathbf{I} - \Phi'(y^*))^{-1}(\psi q) \rangle \quad \forall \psi \in W^{1,\infty}(\Omega) \text{ with } \psi \ge 0, \tag{67}$$

the inequality condition in (62f) can be strengthened to

$$\langle \lambda^*, \psi p^* \rangle > 0 \quad \forall \psi \in W^{1,\infty}(\Omega) \text{ with } \psi > 0.$$

*Proof.* Testing the equation for  $p_{\rho}^*$  with  $(I - \Phi'(y_{\rho}^*))^{-1}(\psi p_{\rho}^*)$ , noticing that  $\psi p_{\rho}^* \rightharpoonup \psi p^*$  in V and making use again of (57) and (58) in a similar way to the proof of Theorem 5.5,

$$\begin{split} \limsup_{\rho \to 0} \langle \mu_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi p_{\rho}^*) \rangle &= \limsup_{\rho \to 0} \langle y_d, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi p_{\rho}^*) \rangle - \liminf_{\rho \to 0} \langle y_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi p_{\rho}^*) \rangle \\ &- \liminf_{\rho \to 0} \langle A^* p_{\rho}^*, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi p_{\rho}^*) \rangle \\ &\leq \langle y_d - y^*, (\mathbf{I} - \Phi'(y^*))^{-1}(\psi p^*) \rangle - \langle A^* p, (\mathbf{I} - \Phi'(y^*))^{-1}(\psi p^*) \rangle \\ &\qquad \qquad \qquad \text{(using (67) for the last term)} \\ &= \langle \mu^*, (\mathbf{I} - \Phi'(y^*))^{-1}(\psi p^*) \rangle \\ &= \langle \lambda^*, \psi p^* \rangle. \end{split}$$

On the other hand, we have

$$\limsup_{\rho \to 0} \langle \mu_\rho^*, (\mathbf{I} - \Phi'(y_\rho^*))^{-1}(\psi p_\rho^*) \rangle = \limsup_{\rho \to 0} \langle \lambda_\rho^*, \psi p_\rho^* \rangle = \limsup_{\rho \to 0} \int_{\Omega} m_\rho' (y_\rho^* - \Phi(y_\rho^*)) (p_\rho^*)^2 \psi \ge 0$$

which implies the result.

**Remark 5.14.** Let us consider when assumption (67) of the previous proposition holds. Suppose that A is of the form

$$\langle Au, v \rangle = \sum_{i,j=1}^{n} \int_{\Omega} a_{ij} \frac{\partial u}{\partial x_i} \frac{\partial v}{\partial x_j} + \sum_{i=1}^{n} \int_{\Omega} b_i \frac{\partial u}{\partial x_i} v + \int_{\Omega} c_0 uv \qquad \forall u, v \in V,$$
 (68)

with  $a_{ij} = a_{ji} \in C^{0,1}(\bar{\Omega}), b_i \in W^{1,\infty}(\Omega), c_0 \in L^{\infty}(\Omega)$  and

$$\sum_{i,j=1}^{n} a_{ij} \xi_i \xi_j \ge C|\xi|^2 \quad a.e. \tag{69}$$

for some C>0 and  $c_0\geq \lambda>0$  a.e. with  $\lambda$  a constant such that A is coercive. Taking  $\psi$  as in the above proposition, let  $z_{\rho}=(\mathrm{I}-\Phi'(y_{\rho}^*))^{-1}(\psi q_{\rho})$ . By (54),  $z_{\rho}\rightharpoonup z:=(\mathrm{I}-\Phi'(y^*))^{-1}(\psi q^*)$  in V. We have, as done in [65, Lemma 3.6] and [70, Lemma 4.5],

$$\langle A^* q_{\rho}, (\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi q_{\rho}) \rangle = \langle A^* q_{\rho}, z_{\rho} \rangle$$

$$= \langle q_{\rho}, A z_{\rho} \rangle$$

$$= \sum_{i,j=1}^{n} \int_{\Omega} a_{ij} \frac{\partial z_{\rho}}{\partial x_{i}} \frac{\partial q_{\rho}}{\partial x_{j}} + \sum_{i=1}^{n} \int_{\Omega} b_{i} \frac{\partial z_{\rho}}{\partial x_{i}} q_{\rho} + \int_{\Omega} c_{0} z_{\rho} q_{\rho}.$$

Using the convergences  $q_{\rho} \rightharpoonup q$  and  $z_{\rho} \rightharpoonup z$  in V, the compactness of  $V \stackrel{c}{\hookrightarrow} H$  and the regularity of  $\psi$ , it is easy to pass to the limit in all but the first term. For that term, we need a weak lower semicontinuity of the form

$$\liminf_{\rho \to 0} \sum_{i,j=1}^{n} \int_{\Omega} a_{ij} \frac{\partial ((\mathbf{I} - \Phi'(y_{\rho}^*))^{-1}(\psi q_{\rho}))}{\partial x_i} \frac{\partial q_{\rho}}{\partial x_j} \ge \sum_{i,j=1}^{n} \int_{\Omega} a_{ij} \frac{\partial ((\mathbf{I} - \Phi'(y))^{-1}(\psi q))}{\partial x_i} \frac{\partial q}{\partial x_j}.$$

A condition ensuring this is the complete continuity of  $I - \Phi'(y) \colon V \to V$  (examining the proof of Lemma 5.9 shows that this condition would turn the convergence in (54) into a strong convergence so that  $z_{\rho} \to z$  in V and hence we can directly pass to the limit in that term).

## From $\mathcal{E}$ -almost to C-stationarity

In order to upgrade to C-stationarity, we need an additional condition given in the next proposition. The assumption preserves generality but is strong, however, we will explore an example below of a reasonable situation where it holds.

**Proposition 5.15** (C-stationarity). Let the assumptions of Theorem 5.11 hold and assume that

$$y_{\rho}^* - \Phi(y_{\rho}^*) \rightarrow y^* - \Phi(y^*)$$
 in  $L^{\infty}(\Omega)$ .

Then (62g) can be strengthened to

$$\langle \lambda^*, v \rangle = 0 \quad \forall v \in V : v = 0 \text{ a.e. on } \{y^* = \Phi(y^*)\}.$$

*Proof.* By assumption, the convergence of  $y_{\rho}^* - \Phi(y_{\rho}^*)$  to  $y^* - \Phi(y^*)$  is uniform and hence  $\rho^{-1}m_{\rho}'(y_{\rho}(x) - \Phi(y_{\rho})(x)) \to 0$  uniformly a.e. globally on  $\{y^* < \Phi(y^*)\}$ . This means that the argument in the proof of Theorem 5.11 can be repeated without recourse to Egorov's theorem.

Sobolev embeddings are the most obvious paths to achieve the assumption of the above proposition. We demonstrate this now with an example. Take the dimension  $n \leq 4$  and suppose that the (bounded Lipschitz) domain  $\Omega$  and operator A are such that

$$y \in H_0^1(\Omega) \cap H^2(\Omega) \implies Ay \in L^2(\Omega)$$

and 17

$$y \in H^1_0(\Omega), Ay \in L^2(\Omega) \implies \begin{cases} y \in H^2(\Omega), \\ \|y\|_{H^2(\Omega)} \le C(\|y\|_{L^2(\Omega)} + \|Ay\|_{L^2(\Omega)}). \end{cases}$$

<sup>&</sup>lt;sup>17</sup>These are elliptic regularity conditions. When  $\Omega$  is a  $C^{1,1}$  domain and A is of the form (68) with  $a_{ij} \in C^0(\bar{\Omega}) \cap W^{1,\infty}(\Omega), b_i, c_0 \in L^\infty(\Omega),$  $c_0 \ge 0$  with the strict ellipticity (69), Theorem 9.15 of [29] can be applied and it implies the first condition. The second follows from [29, Lemma 9.17].

We take  $V = H_0^1(\Omega)$  and use the fact that  $m_\rho: V \to V$  (recall that  $m_\rho$  has been chosen in (63); see [22, §2.2.3] for when this type of property could hold for other maps). Suppose that  $\Phi: V^* \to V$  is given by the solution mapping of an elliptic equation, i.e.,  $\Phi(y)$  is defined as the solution  $\phi$  of

$$B(\phi) = y$$

where B is a second-order elliptic operator with sufficient properties guaranteeing well posedness in  $H_0^1(\Omega)$  (with a continuous dependence estimate), and when  $y \in L^2(\Omega)$ , in the space  $H^2(\Omega) \cap H_0^1(\Omega)$  including a regularity estimate of the form

$$\|\phi\|_{H^2(\Omega)} \le C \|y\|_{L^2(\Omega)}.$$

Due to this, we immediately have that  $\Phi(y^*_{\rho}) \in H^2(\Omega)$  with a uniform bound:

$$\|\Phi(y_{\rho}^*)\|_{H^2(\Omega)} \le C_1 \|y_{\rho}^*\|_{L^2(\Omega)} \le C_2.$$
 (70)

Defining  $z = y_{\rho}^* - \Phi(y_{\rho}^*) \in H_0^1(\Omega)$ , we write the equation for  $y_{\rho}^*$  as

$$Az + \frac{1}{\rho}m_{\rho}(z) = u_{\rho}^* - A\Phi(y_{\rho}^*).$$

It follows from rearranging this equation that  $Az \in H$ , thus  $z \in H^2(\Omega)$  and the equation holds in a pointwise a.e. sense. Suppose for simplicity that  $A = -\Delta$  is the Dirichlet Laplacian. Test with  $-\Delta z$  and use

$$\int_{\Omega} m_{\rho}(z)(-\Delta z) = \int_{\Omega} m_{\rho}'(z)|\nabla z|^2 \ge 0$$

to obtain

$$\|-\Delta z\|_H^2 \le \|u_\rho^* - A\Phi(y_\rho^*)\|_H \|-\Delta z\|_H$$
.

Dividing through by  $\|-\Delta z\|_H$ , the resulting right-hand side is bounded due to (70), and using the regularity condition above, we obtain uniform boundedness in  $H^2(\Omega)$  of  $z=y_\rho^*-\Phi(y_\rho^*)$ . By the Sobolev embedding [1, Theorem 6.3]  $H^2(\Omega) \stackrel{c}{\hookrightarrow} C^{0,\alpha}(\bar{\Omega})$  for some  $\alpha \in (0,1)$ , we get  $y_\rho^*-\Phi(y_\rho^*) \to y^*-\Phi(y^*)$  in that Hölder space (and thus in  $L^\infty(\Omega)$ ).

#### 5.5 Strong stationarity

We now give strong stationarity conditions for (2) in the setting of  $V = H_0^1(\Omega)$ ,  $H = L^2(\Omega)$  and  $U_{ad}$  of the box constraint form (61).

Let us first of all provide some background and context. Strong stationarity for the VI obstacle problem in the absence of constraints on the control was the focus of the classical works by Mignot [50, Theorem 5.2] and Mignot and Puel [51]. The approach in the latter work is as follows. By using the results on the differentiability of the solution map associated to VIs of Mignot [50], the Bouligand stationarity condition (for example, see Proposition 5.2) reads

$$(\alpha_h, y^* - y_d)_H + \nu(u^*, h)_H \ge 0 \quad \forall h \in H$$

where  $\alpha_h$  denotes the directional derivative of the solution map with respect to the direction h. The key idea of Mignot and Puel in [51] is to use the fact that the optimal control  $u^*$  in fact belongs to V (in the unconstrained case, this follows from B-stationarity; otherwise this is a regularity result in certain situations or one may need to simply assume this) and to extend, by continuity, the above inequality to

$$(\alpha_h, y^* - y_d)_H + \nu \langle u^*, h \rangle \ge 0 \quad \forall h \in V^*$$

$$(71)$$

so that the set of feasible directions has been enlarged to  $V^*$ . Then, by writing the duality product in (71) as  $\langle AA^{-1}h, \nu u^*\rangle$  and using properties of the projection operator with respect to the bilinear form generated by A onto the critical cone, it is shown [51, Theorem 3.3] that this inequality is equivalent to a strong stationarity system.

The presence of control constraints complicates the derivation of strong stationarity conditions. In the VI setting, by using the above-mentioned technique of Mignot and Puel of enlarging the set of feasible directions onto the dual space in combination with a fine analysis of the various resulting objects and sets, strong stationarity conditions for VI optimal control problems subject to box constraints were obtained by Wachsmuth in [69]. The author also showed that certain restrictions are required on the control bounds in order to obtain a positive answer for strong stationarity, and counterexamples were given showing that violating those conditions can lead to a lack of strong stationarity. These necessary conditions (which are stated in (72)–(74) below) in the context of admissible sets as in (61) are implied [69, Lemma 5.3] by the condition

$$u_a, u_b \in H^1(\Omega)$$
 with  $u_a < 0 \le u_b$  q.e. on  $\Omega$ ,

(recall Example 3.5 for the meaning of q.e.) which in turn implies that the control space must allow for negative functions, meaning that one ultimately needs existence and directional differentiability results for QVIs with source terms and directions that may be strictly negative<sup>18</sup>.

Let  $(y^*, u^*)$  be a local optimal pair of (2). As in [51], we make the fundamental assumption that  $u^* \in V$  and we refer to Theorem 5.11 from the previous section for the satisfaction of this assumption. Let us take  $U_{ad}$  as stated in (61) where we include the possibility of taking  $u_a = -\infty$  and  $u_b = \infty$ , in which case the problem becomes one with no constraints and we can argue as in [51]. Outside of this case, we proceed as in [69]. Let the assumptions of Theorem 3.2 hold and denote by  $j: H \to V^*$  the inclusion map through the Riesz isomorphism. Then, as done in [69], the Bouligand stationarity condition (42) can be extended to

$$(\alpha_h, y^* - y_d) + \nu \langle h, u^* \rangle \ge 0 \quad \forall h \in \overline{j \mathcal{T}_{U_{ad}}(u^*)}^{V^*}.$$

This is starting point of the steps leading to the strong stationarity conditions in [69] for the VI case.

Defining the (quasi-closed) coincidence sets

$$U_a := \{x \in \Omega : u^*(x) = u_a(x)\}$$
 and  $U_b := \{x \in \Omega : u^*(x) = u_b(x)\}$ 

and arguing identically to the proof of [69, Lemma 4.3], we obtain the following sign conditions on  $u^*$ :

$$u^* = 0$$
 q.e. on  $\mathcal{A}_s(y^*) \cap (\Omega \setminus (U_a \cup U_b)),$   
 $u^* \leq 0$  q.e. on  $\mathcal{A}_s(y^*) \cap U_b,$   
 $u^* \geq 0$  q.e. on  $(\mathcal{A}_s(y^*) \cap U_a) \cup (\mathcal{B}(y^*) \cap (\Omega \setminus U_b))$ 

where  $\mathcal{B}(y^*) = \mathcal{A}(y^*) \setminus \mathcal{A}_s(y^*)$  is the *biactive set*.

Let cap(A) denote the capacity of a Borel subset A of  $\Omega$  with respect to  $H_0^1(\Omega)$  (see [20, Definition 6.47]). We have the following strong stationarity characterisation, the proof of which involves modifications of [69] and is sketched in Appendix B.

**Theorem 5.16** (Strong stationarity). Let  $(y^*, u^*)$  be a local minimiser of (2) with  $u^* \in V$ . Assume Assumption 3.1, (19), the local assumptions  $^{19}$  (25), (32) and suppose that

$$\Phi \colon V \to V \text{ is Frèchet differentiable at } y^*,$$

$$\operatorname{cap}(U_a \cap \mathcal{B}(y^*)) = 0, \tag{72}$$

$$u_b \ge 0 \text{ q.e. on } \mathcal{B}(y^*),$$
 (73)

$$u^* = 0 \text{ q.e. on } A_s(y^*).$$
 (74)

Then  $(y^*, u^*)$  is a strong stationarity point, i.e., there exist multipliers  $(p^*, \xi^*, \lambda^*) \in V \times V^* \times V^*$  such that

$$\begin{split} y^* + (\mathrm{I} - \Phi'(y^*)^*) \lambda^* + A^* p^* &= y_d, \\ Ay^* - u^* + \xi^* &= 0, \\ \xi^* &\geq 0 \text{ in } V^*, \quad y^* \leq \Phi(y^*), \quad \langle \xi^*, y^* - \Phi(y^*) \rangle &= 0, \\ u^* &\in U_{ad} : (\nu u^* - p^*, u^* - v) \leq 0 \quad \forall v \in U_{ad}, \\ p^* &\geq 0 \quad \textit{q.e. on } \mathcal{B}(y^*) \text{ and } p^* &= 0 \text{ q.e. on } \mathcal{A}_s(y^*), \\ \langle \lambda^*, v \rangle &\geq 0 \quad \forall v \in V : v \geq 0 \text{ q.e. on } \mathcal{B}(y^*) \text{ and } v = 0 \text{ q.e. on } \mathcal{A}_s(y^*). \end{split}$$

Note also that, whilst this work was under preparation, a related result has recently been obtained in [72] however only in the absence of control constraints (i.e.,  $U_{ad}$  is taken to be the whole space).

#### A **Technical proofs**

 $Proof^{20}$  of Lemma 2.3. Take an arbitrary subsequence  $\{v_{n_i}\}$ ; this remains uniformly bounded hence we can extract a

weakly convergent subsequence such that  $v_{n_{j_k}} \rightharpoonup v$  in V to some v. Select an arbitrary  $f \in V_+^*$  and set  $l_n := \langle f, v_n \rangle$  which is a monotonic sequence (since f is non-negative) and also bounded. Hence the monotone convergence theorem applies and we obtain the existence of l such that  $l_n \to l$ . Since also  $l_{n_{j_k}} \to l$ , we conclude that  $l = \langle f, v \rangle$ .

<sup>&</sup>lt;sup>18</sup>Our theory of differentiability for QVIs in the earlier paper [4] (which was for non-negative sources and directions) could not be immediately used to obtain strong stationarity by arguing in this fashion since the setting of [4] would have forced  $U_{ad}$  to be selected such that  $U_{ad} \subset H_+$ . This is why the development of the results of §2 and §3 are crucial.

<sup>&</sup>lt;sup>19</sup>These, of course, should be evaluated at  $y^*$ .

<sup>&</sup>lt;sup>20</sup>We thank Jochen Glück for the idea of the proof.

Take another subsequence of  $\{v_n\}$ , say  $\{v_{n_m}\}$ , then by the above argument, we have  $v_{n_{m_j}} \rightharpoonup \hat{v}$  for some  $\hat{v}$  and  $l = \langle f, \hat{v} \rangle$ . That is,

$$\langle f, v \rangle = \langle f, \hat{v} \rangle \quad \forall f \in V_+^*,$$

and from this, we can conclude via the weak-\* density of  $V_+^* - V_+^*$  in  $V^*$  (e.g., see [8, Lemma 2.7]) that  $\hat{v} = v$ . The subsequence principle then yields the result.

*Proof of Lemma* 5.9. Define  $T_n = (I - \Phi'(z_n))$  and  $T = (I - \Phi'(z))$ . Then

$$T_n^{-1}q_n - T^{-1}q = (T_n^{-1} - T^{-1})q_n + T^{-1}(q_n - q)$$

and we get  $T^{-1}(q_n-q) \rightharpoonup 0$  in V by continuity and linearity of  $T^{-1}$ . For the first term on the right-hand side above, we use the identity  $T_n^{-1}-T^{-1}=T_n^{-1}(T-T_n)T^{-1}$  relating the inverses of operators to see that

$$\begin{split} \left\| (T_{n}^{-1} - T^{-1}) q_{n} \right\|_{V} &= \left\| T_{n}^{-1} (T - T_{n}) T^{-1} q_{n} \right\|_{V} \\ &\leq C_{1} \left\| (T - T_{n}) T^{-1} q_{n} \right\|_{V} \\ &\leq C_{1} \left\| T - T_{n} \right\|_{\mathcal{L}(V,V)} \left\| T^{-1} q_{n} \right\|_{V} \\ &\leq C_{2} \left\| \Phi'(z_{n}) - \Phi'(z) \right\|_{\mathcal{L}(V,V)} \\ &\rightarrow 0 \end{split} \tag{because $T^{-1}$ and $q_{n}$ are bounded)}$$

with the convergence because we assumed that  $\Phi$  is continuously Fréchet differentiable and hence the derivative is continuous. Therefore,  $T_n^{-1}q_n \rightharpoonup T^{-1}q$  in V. The strong convergence follows because if  $q_n \to q$  then  $T^{-1}(q_n - q) \to 0$  in V. For the final claim, we have

$$\langle AT_n^{-1}q_n, q_n \rangle - \langle AT^{-1}q, q \rangle = \langle A(T_n^{-1}q_n - T^{-1}q_n), q_n \rangle + \langle AT^{-1}q_n, q_n \rangle - \langle AT^{-1}q, q \rangle$$

and the first term on the right-hand side tends to zero by the calculation above. Since by (50),  $AT^{-1}$  is bounded and coercive (as well as being linear), we obtain

$$\liminf_{n \to \infty} \langle AT^{-1}q_n, q_n \rangle - \langle AT^{-1}q, q \rangle \ge 0.$$

# B Sketch proof of Theorem 5.16

Recall the notation  $\alpha_h$  which stands for the directional derivative in the direction h given through Theorem 3.2.

**Lemma B.1.** Denote by  $j: H \to V^*$  the inclusion map. Then  $0 \in V^*$  is a minimiser of the problem

$$\min_{h \in \overline{jT_{U_{ad}}(u^*)}^{V^*}} (\alpha_h, y^* - y_d)_H + \nu \langle h, u^* \rangle.$$
(75)

*Proof.* Choosing the direction h=0 in the inequality of Proposition 5.2 implies  $0 \le (\alpha_0, y^* - y_d) + \nu(u^*, 0) = 0$  with the equality because  $\alpha_0 = 0$ . Hence h=0 is a minimiser of

$$\min_{h \in \mathcal{T}_{U_{a,d}}(u^*)} (\alpha_h, y^* - y_d)_H + \nu(u^*, h).$$

As in Lemma 4.1 of [69], the feasible set can be enlarged (the continuity in  $V^*$  of  $h \mapsto \alpha_h$  assured by Proposition 3.12 is needed here) to obtain the desired result.

The aim now is to rewrite (75) over the space

$$W := \{ v \in V : v = 0 \text{ q.e. in } A_s(y^*) \}.$$

Using the characterisation of the critical cone from [69, Lemma 3.1], we see that  $\mathcal{K}^{y^*} \subset W$ . Denote by  $i \colon W \to V$  the inclusion map and define the closed convex set

$$\mathcal{C}_W^{y^*} := \{v \in W : v \leq 0 \text{ q.e. in } \mathcal{B}(y^*)\},$$

which satisfies  $\mathcal{K}^{y^*}=i\mathcal{C}_W^{y^*}$ . Now, note that, using (32),  $(\mathbf{I}-\Phi'(y^*))\colon V\to V$  is invertible. Define

$$A_W: W \to W^*, \qquad A_W:=i^*A(I - \Phi'(y^*))^{-1}i$$

and observe that for any  $\tilde{d} \in W^*$  the inequality

$$\delta \in \mathcal{C}_W^{y^*}: \langle A_W \delta - \tilde{d}, \delta - w \rangle_{W^*, W} \le 0 \quad \forall w \in \mathcal{C}_W^{y^*}$$

has a unique solution by the Lions–Stampacchia theorem since  $A_W$  is bounded and coercive due to the Lipschitz condition (32) (see [72, Lemma 3.3]). Now suppose that for  $d \in V^*$ ,  $\delta$  solves

$$\delta \in \mathcal{C}_W^{y^*}: \langle A_W \delta - i^* d, \delta - w \rangle_{W^*, W} \le 0 \quad \forall w \in \mathcal{C}_W^{y^*}.$$

Consider also

$$z \in \mathcal{K}^{y^*} : \langle A(\mathbf{I} - \Phi'(y^*))^{-1}z - d, z - v \rangle_{V^*, V} \le 0 \quad \forall v \in K^{y^*}.$$

Then it is easy to see that that  $z = i\delta$ .

**Lemma B.2.** Define the operator  $\theta: W \to V$  by

$$\theta := (I - \Phi'(y^*))^{-1}i.$$

Then (0,0) is a solution of

$$\min_{(\beta_h,h)\in W\times W^*} (\theta(\beta_h), y^* - y_d)_H + \nu \langle h, u^* \rangle_{W^*,W} \text{ s.t.}$$

$$\begin{cases}
\beta_h \in \mathcal{C}_W^{y^*} \\
h = A_W \beta_h \\
h \in \overline{i^* j \mathcal{T}_{U_{ad}}(u^*)}^{W^*}.
\end{cases}$$
(76)

*Proof.* By defining  $\gamma_h := \alpha_h - \Phi'(y^*)(\alpha_h) = (I - \Phi'(y^*))\alpha_h$ , the QVI (27) satisfied by  $\alpha_h$  can be written as

$$\gamma_h \in \mathcal{K}^{y^*} : \langle A(I - \Phi'(y^*))^{-1} \gamma_h - h, \gamma_h - \varphi \rangle \le 0 \quad \forall \varphi \in \mathcal{K}^{y^*}.$$

Now if  $\beta_h$  satisfies

$$\beta_h \in \mathcal{C}_W^{y^*}: \langle A_W \beta_h - i^* h, \beta_h - \varphi \rangle \le 0 \quad \forall \varphi \in \mathcal{C}_W^{y^*},$$

we have (as discussed above)  $\gamma_h = i\beta_h$ , hence

$$i\beta_h = (I - \Phi'(y^*))\alpha_h \quad \iff \quad \alpha_h = \theta(\beta_h)$$

Therefore, (75) can be restated and we get (using the continuity of  $\Phi'(y^*)$ ) that 0 is a solution of

$$\begin{split} \min_{h \in \overline{i^* j \mathcal{T}_{U_{ad}}(u^*)}^{W^*}} (\theta(\beta_h), y^* - y_d)_H + \nu \langle h, u^* \rangle_{W^*, W} \text{ s.t.} \\ \beta_h \in \mathcal{C}_W^{y^*} : \langle A_W \beta_h - h, \beta_h - \varphi \rangle_{W_{s,*}^*, W} \leq 0 \quad \forall \varphi \in \mathcal{C}_W^{y^*}; \end{split}$$

this is well defined because  $u^* \in W$  due to (74). Hence, similarly to Proposition 3.13, (0,0,0) is a solution of

$$\min_{(\beta_h,h,\xi_h)\in W\times W^*\times W^*} (\theta(\beta_h),y^*-y_d)_H + \nu\langle h,u^*\rangle_{W^*,W} \text{ s.t.}$$

$$\begin{cases} \beta_h\in \mathcal{C}_W^{y^*} \\ \xi_h=h-A_W\beta_h \\ \xi_h\in (\mathcal{C}_W^{y^*})^\circ \\ \langle \xi_h,\beta_h\rangle=0 \\ h\in \overline{i^*i\mathcal{T}_U},(u^*)^{W^*}. \end{cases}$$

Setting  $\xi_h = 0$  leads to the result.

We need to derive stationarity conditions for this problem and then transform the resulting system back to the original spaces and operators. Let us remark that under the assumptions of the theorem, we have that  $\theta$  is linear and bounded.

#### Lemma B.3. Defining

$$D := \overline{i^* j \mathcal{T}_{U_{ad}}(u^*)}^{W^*}, \qquad \mathcal{Y} := W^* \times W \times W^*, \qquad C := (\{0\}, \mathcal{C}_W^{y^*}, D),$$

there exists  $(\tilde{p}, \tilde{\lambda}, \sigma) \in \mathcal{Y}^* \cap C^{\circ}$  such that

$$A_W^* \tilde{p} + \theta^* (j(y^* - y_d)) + \tilde{\lambda} = 0,$$
  

$$\nu u^* - \tilde{p} + \sigma = 0,$$
  

$$\tilde{\lambda} \in (\mathcal{C}_W^{y^*})^{\circ},$$
  

$$\sigma \in D^{\circ}.$$

*Proof.* In addition to the notation introduced above, let us also define the space  $\mathcal{X} := W \times W^*$ . Define the map  $g \colon \mathcal{X} \to \mathcal{Y}$  by  $g(\beta, h) := (A_W \beta - h, \beta, h)$  and observe that (76) can be compactly written as

$$\min_{q(\beta_h, h) \in C} (\theta(\beta_h), y^* - y_d) + \nu \langle h, u^* \rangle_{W^*, W}. \tag{77}$$

We now proceed with checking the Zowe–Kurcyusz constraint qualification  $g'((0,0))\mathcal{X} - \mathcal{R}_C(g(0,0)) = \mathcal{Y}$  to deduce the existence of Lagrange multipliers. First observe that D is a convex cone which in turn implies that C is a convex cone and then by [20, Example 2.62],  $\mathcal{R}_C((0,0,0)) = C$  and  $\mathcal{T}_C((0,0,0))^\circ = C^\circ$ . Now, we see that g(0,0) = (0,0,0) and  $\mathcal{R}_C(g(0,0)) = C$ . We also have

$$g'(0,0)(\gamma,d) = (A_W \gamma - d, \gamma, d) \qquad \forall (\gamma,d) \in W \times W^*.$$

Therefore, we are required to show that for every  $(w_1^*, w_2, w_3^*) \in \mathcal{Y}$ , there exist  $(\gamma, d, v, h) \in \mathcal{X} \times \mathcal{C}_W^{y^*} \times D$  such that

$$A_W \gamma - d = w_1^*,$$
  

$$\gamma - v = w_2,$$
  

$$d - h = w_3^*.$$
(78)

The first equation written in terms of v and h reads  $A_W v - (w_1^* + w_3^* - A_W w_2) = h$ . In order to force solutions to belong to the desired sets, we consider the VI

find 
$$v \in \mathcal{C}_W^{y^*}: \langle A_W v - (w_1^* + w_3^* - A_W w_2), v - \varphi \rangle \le 0 \quad \forall \varphi \in \mathcal{C}_W^{y^*}$$
 (79)

associated to the above PDE. As explained above, (79) has a solution and furthermore, the following complementarity system (which can be derived by the same arguments as before) is satisfied by any solution:

$$\begin{cases} v \in \mathcal{C}_{W}^{y^{*}} \\ \eta := (w_{1}^{*} + w_{3}^{*} - A_{W}w_{2}) - A_{W}v \\ \eta \in (\mathcal{C}_{W}^{y^{*}})^{\circ} \\ \eta \perp v. \end{cases}$$

Using this, we see that  $h:=-\eta\in -(\mathcal{C}_W^{y^*})^\circ$ . The manipulations in the paragraph after Lemma 5.1 of [69] show that  $(i^*j\mathcal{T}_{U_{ad}}(y^*))^\circ\subset -\mathcal{C}_W^{y^*}$  which implies that  $-(\mathcal{C}_W^{y^*})^\circ\subset (i^*j\mathcal{T}_{U_{ad}}(y^*))^{\circ\circ}=D$ , that is,  $g\in D$ . Then we simply define  $\gamma$  and d by (78). Thus the constraint qualification is met for (77).

Writing the objective functional in (77) as  $\hat{J}$ , we obtain the existence of a Lagrange multipler  $(\tilde{p}, \tilde{\lambda}, \sigma) \in \mathcal{Y}^* \cap C^\circ$  such that

$$\hat{J}'(0,0)(x) + \langle q'(0,0)^*(\tilde{p},\tilde{\lambda},\sigma), x \rangle = 0 \quad \forall x \in \mathcal{X}.$$

With  $x = (\gamma, d)$ , we see that since  $\theta(0) = 0$ , the first term above is

$$\hat{J}'(0,0)(x) = \langle \theta^*(j(y^* - y_d)), \gamma \rangle_{W^*,W} + \nu \langle d, u^* \rangle_{W^*,W},$$

where  $\theta^* \colon V^* \to W^*$  is the adjoint of  $\theta \colon W \to V$  (this exists due to the linearity assumption). We also have, by definition of the adjoint operator,

$$\langle g'(0,0)^*(\tilde{p},\tilde{\lambda},\sigma), x \rangle = \langle (\tilde{p},\tilde{\lambda},\sigma), (A_W\gamma - d,\gamma,d) \rangle$$
  
=  $\langle A_W^*\tilde{p}, \gamma \rangle_{W^*,W} + \langle \tilde{\lambda}, \gamma \rangle_{W^*,W} + \langle \sigma - \tilde{p}, d \rangle_{W^*,W}.$ 

This implies the result.

We now transform all quantities back to the space V.

Conclusion of sketch proof of Theorem 5.16. Observe that under the assumptions, Proposition 5.2, Lemma B.3 and Theorem 3.2 are applicable. To start with, let us define

$$p^* := i\hat{x}$$

and

$$\lambda^* := (I - \Phi'(y^*)^*)^{-1} (-A^* i \tilde{p} - j(y - y_d)),$$

and for convenience, denote  $L := \Phi'(y^*)$ .

- By definition of  $\lambda^*$  and  $p^*$ , we get the first line in the system after etching away the inclusion map j.
- We see from the definition of  $\lambda^*$  and elementary manipulations to relate it to  $\tilde{\lambda} \in (C_W)^\circ$  and the usage of the fact that  $iC_W = \mathcal{K}^{y^*}$  that  $\lambda^* \in (\mathcal{K}^{y^*})^\circ$ . This implies the final condition of the system thanks to [69, Lemma 3.1].

• Since  $\tilde{p} \in W$ , it vanishes q.e. on the strongly active set. As  $\tilde{p} = \nu u^* + \sigma$  and since  $\sigma \in D^\circ$ , Lemma 5.1 of [69] tells us that  $\sigma \geq 0$  q.e. on  $\Omega \setminus U_a$ . Thus

$$\sigma|_{\mathcal{B}(y^*)} = \sigma|_{U_a \cap \mathcal{B}(y^*)} + \sigma|_{(\Omega \setminus U_a) \cap \mathcal{B}(y^*)} \ge \sigma|_{U_a \cap \mathcal{B}(y^*)} = 0$$

with the final equality because of (72). Note also that

$$u^*|_{\mathcal{B}(y^*)} = u^*|_{\mathcal{B}(y^*) \cap U_b} + u^*|_{\mathcal{B}(y^*) \cap (\Omega \setminus U_b)} \ge u^*|_{\mathcal{B}(y^*) \cap (\Omega \setminus U_b)} \ge 0$$
 q.e.,

with the first inequality by (73) and the final inequality by the third sign condition on  $u^*$  stated in §5.5. This implies the stated condition on  $p^*$ , which is equivalent to  $-p^* \in \mathcal{K}^{y^*}$  due to the characterisation of the critical cone in [69, Lemma 3.1].

• We obtain  $\sigma \in \mathcal{N}_{U_{ad}}(u^*)$  exactly as in the proof of Theorem 5.2 in [69]<sup>21</sup> (where  $\mathcal{N}_{U_{ad}}$  denotes the normal cone to  $U_{ad}$  with respect to H), which is the polar cone of the tangent cone, see [20, §2.2.4]) and this is precisely the desired inequality constraint relating the control and the adjoint.

## References

- [1] R. A. Adams and J. J. F. Fournier. Sobolev Spaces. Academic Press, second edition, 2003.
- [2] S. Adly, M. Bergounioux, and M. Ait Mansour. Optimal control of a quasi-variational obstacle problem. *J. Global Optim.*, 47(3):421–435, 2010.
- [3] C. D. Aliprantis and O. Burkinshaw. *Positive operators*, volume 119 of *Pure and Applied Mathematics*. Academic Press, Inc., Orlando, FL, 1985.
- [4] A. Alphonse, M. Hintermüller, and C. N. Rautenberg. Directional differentiability for elliptic quasi-variational inequalities of obstacle type. *Calc. Var. Partial Differential Equations*, 58(1):Art. 39, 47, 2019.
- [5] A. Alphonse, M. Hintermüller, and C. N. Rautenberg. Recent trends and views on elliptic quasi-variational inequalities. In M. Hintermüller and J. F. Rodrigues, editors, *Topics in Applied Analysis and Optimisation*, pages 1–31, Cham, 2019. Springer International Publishing.
- [6] A. Alphonse, M. Hintermüller, and C. N. Rautenberg. Stability of the solution set of quasi-variational inequalities and optimal control. *SIAM J. Control Optim.*, 58(6):3508–3532, 2020.
- [7] A. Alphonse, C. N. Rautenberg, and J. F. Rodrigues. Analysis of a quasi-variational contact problem arising in thermoelasticity. *arXiv e-prints*. Available at https://arxiv.org/abs/2008.00890.
- [8] W. Arendt and R. Nittka. Equivalent complete norms and positivity. Arch. Math. (Basel), 92(5):414–427, 2009.
- [9] J.-P. Aubin. *Mathematical methods of game and economic theory*, volume 7 of *Studies in Mathematics and its Applications*. North-Holland Publishing Co., Amsterdam-New York, 1979.
- [10] A. Azevedo, F. Miranda, and L. Santos. Variational and quasivariational inequalities with first order constraints. *J. Math. Anal. Appl.*, 397(2):738–756, 2013.
- [11] A. Azevedo, F. Miranda, and L. Santos. *Stationary Quasivariational Inequalities with Gradient Constraint and Nonhomogeneous Boundary Conditions*, pages 95–112. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014.
- [12] C. Baiocchi and A. Capelo. Variational and Quasivariational Inequalities. Wiley-Interscience, 1984.
- [13] V. Barbu. *Optimal control of variational inequalities*, volume 100 of *Research Notes in Mathematics*. Pitman (Advanced Publishing Program), Boston, MA, 1984.
- [14] J. W. Barrett and L. Prigozhin. A quasi-variational inequality problem in superconductivity. *Math. Models Methods Appl. Sci.*, 20(5):679–706, 2010.
- [15] J. W. Barrett and L. Prigozhin. A quasi-variational inequality problem arising in the modeling of growing sandpiles. *ESAIM Math. Model. Numer. Anal.*, 47(4):1133–1165, 2013.
- [16] J. W. Barrett and L. Prigozhin. Lakes and rivers in the landscape: a quasi-variational inequality approach. *Interfaces Free Bound.*, 16(2):269–296, 2014.

<sup>&</sup>lt;sup>21</sup>In [69], the notation  $\mu$  is used instead of  $\sigma$ .

- [17] A. Bensoussan, M. Goursat, and J.-L. Lions. Contrôle impulsionnel et inéquations quasi-variationnelles stationnaires. *C. R. Acad. Sci. Paris Sér. A-B*, 276:A1279–A1284, 1973.
- [18] M. Bergounioux. Optimal control of an obstacle problem. Appl. Math. Optim., 36(2):147–172, 1997.
- [19] M. Bergounioux. Use of augmented Lagrangian methods for the optimal control of obstacle problems. *J. Optim. Theory Appl.*, 95(1):101–126, 1997.
- [20] J. F. Bonnans and A. Shapiro. *Perturbation analysis of optimization problems*. Springer Series in Operations Research. Springer-Verlag, New York, 2000.
- [21] J. T. Cal Neto and C. Tomei. Numerical analysis of semilinear elliptic equations with finite spectral interaction. *J. Math. Anal. Appl.*, 395(1):63–77, 2012.
- [22] S. Carl, V. K. Le, and D. Motreanu. *Nonsmooth variational problems and their inequalities*. Springer Monographs in Mathematics. Springer, New York, 2007. Comparison principles and applications.
- [23] H. Dietrich. Über Probleme der optimalen Steuerung mit einer Quasivariationsungleichung. In *Operations Research Proceedings 1999 (Magdeburg)*, pages 111–116. Springer, Berlin, 2000.
- [24] H. Dietrich. Optimal control problems for certain quasivariational inequalities. *Optimization*, 49(1-2):67–93, 2001. In celebration of Prof. Dr. Alfred Göpfert 65th birthday.
- [25] F. Facchinei and C. Kanzow. Generalized nash equilibrium problems. 40R, 5(3):173-210, Sep 2007.
- [26] T. Fukao and N. Kenmochi. Abstract theory of variational inequalities with Lagrange multipliers and application to nonlinear PDEs. *Math. Bohem.*, 139(2):391–399, 2014.
- [27] T. Fukao and N. Kenmochi. Quasi-variational inequality approach to heat convection problems with temperature dependent velocity constraint. *Discrete Contin. Dyn. Syst.*, 35(6):2523–2538, 2015.
- [28] M. Fukushima, Y. Oshima, and M. Takeda. *Dirichlet forms and symmetric Markov processes*, volume 19 of *De Gruyter Studies in Mathematics*. Walter de Gruyter & Co., Berlin, extended edition, 2011.
- [29] D. Gilbarg and N. S. Trudinger. Elliptic Partial Differential Equations of Second Order. Springer-Verlag, 1983.
- [30] R. Glowinski. Numerical Methods for Nonlinear Variational Problems. Springer-Verlag, 1984.
- [31] P. Hájek and M. Johanis. *Smooth analysis in Banach spaces*, volume 19 of *De Gruyter Series in Nonlinear Analysis and Applications*. De Gruyter, Berlin, 2014.
- [32] A. Haraux. How to differentiate the projection on a convex set in Hilbert space. Some applications to variational inequalities. *J. Math. Soc. Japan*, 29(4):615–631, 1977.
- [33] F. Harder and G. Wachsmuth. Comparison of optimality systems for the optimal control of the obstacle problem. *GAMM-Mitt.*, 40(4):312–338, 2018.
- [34] P. T. Harker. Generalized nash games and quasi-variational inequalities. *European Journal of Operational Research*, 54(1):81 94, 1991.
- [35] M. Hintermüller and I. Kopacka. Mathematical programs with complementarity constraints in function space: *C*-and strong stationarity and a path-following algorithm. *SIAM J. Optim.*, 20(2):868–902, 2009.
- [36] M. Hintermüller and I. Kopacka. A smooth penalty approach and a nonlinear multigrid algorithm for elliptic MPECs. *Comput. Optim. Appl.*, 50(1):111–145, 2011.
- [37] M. Hintermüller, B. S. Mordukhovich, and T. M. Surowiec. Several approaches for the derivation of stationarity conditions for elliptic MPECs with upper-level control constraints. *Math. Program.*, 146(1-2, Ser. A):555–582, 2014.
- [38] M. Hintermüller and T. Surowiec. First-order optimality conditions for elliptic mathematical programs with equilibrium constraints via variational analysis. *SIAM J. Optim.*, 21(4):1561–1593, 2011.
- [39] A. Kadoya, N. Kenmochi, and M. Niezgódka. Quasi-variational inequalities in economic growth models with technological development. *Adv. Math. Sci. Appl.*, 24(1):185–214, 2014.
- [40] R. Kano, Y. Murase, and N. Kenmochi. Nonlinear evolution equations generated by subdifferentials with nonlocal constraints. In *Nonlocal and abstract parabolic equations and their applications*, volume 86 of *Banach Center Publ.*, pages 175–194. Polish Acad. Sci. Inst. Math., Warsaw, 2009.

- [41] N. Kenmochi. Parabolic quasi-variational diffusion problems with gradient constraints. *Discrete Contin. Dyn. Syst. Ser. S*, 6(2):423–438, 2013.
- [42] D. Kinderlehrer and G. Stampacchia. An Introduction to Variational Inequalities and Their Applications. SIAM, 2000
- [43] K. Kunisch and D. Wachsmuth. Sufficient optimality conditions and semi-smooth Newton methods for optimal control of stationary variational inequalities. *ESAIM Control Optim. Calc. Var.*, 18(2):520–547, 2012.
- [44] M. Kunze and J. F. Rodrigues. An elliptic quasi-variational inequality with gradient constraints and some of its applications. *Math. Methods Appl. Sci.*, 23(10):897–908, 2000.
- [45] T. Laetsch. A uniqueness theorem for elliptic quasi-variational inequalities. *J. Functional Analysis*, 18:286–287, 1975.
- [46] H. Lewy and G. Stampacchia. On the regularity of the solution of a variational inequality. *Comm. Pure Appl. Math.*, 22:153–188, 1969.
- [47] J.-L. Lions. Quelques Méthodes de Résolutions des Problémes aux Limites non Linéaires. Dunod, Gauthier-Villars, 1969.
- [48] J.-L. Lions. Sur le côntrole optimal des systemes distribuées. *Enseigne*, 19:125–166, 1973.
- [49] P. Meyer-Nieberg. Banach lattices. Universitext. Springer-Verlag, Berlin, 1991.
- [50] F. Mignot. Contrôle dans les inéquations variationelles elliptiques. J. Functional Analysis, 22(2):130-185, 1976.
- [51] F. Mignot and J.-P. Puel. Optimal control in some variational inequalities. SIAM J. Control Optim., 22(3):466–476, 1984.
- [52] F. Miranda, J.-F. Rodrigues, and L. Santos. On a p-curl system arising in electromagnetism. *Discrete Contin. Dyn. Syst. Ser. S*, 5(3):605–629, 2012.
- [53] B. S. Mordukhovich and J. Outrata. Coderivative analysis of quasi-variational inequalities with applications to stability and optimization. *SIAM J. Optim.*, 18(2):389–412, 2007.
- [54] Y. Murase, A. Kadoya, and N. Kenmochi. Optimal control problems for quasi-variational inequalities and its numerical approximation. *Discrete Contin. Dyn. Syst.*, (Dynamical systems, differential equations and applications. 8th AIMS Conference. Suppl. Vol. II):1101–1110, 2011.
- [55] J.-S. Pang and M. Fukushima. Quasi-variational inequalities, generalized nash equilibria, and multi-leader-follower games. *Computational Management Science*, 2(1):21–56, Jan 2005.
- [56] L. Prigozhin. Sandpiles and river networks: extended systems with non-local interactions. *Phys. Rev. E*, 49:1161–1167, 1994.
- [57] L. Prigozhin. On the Bean critical-state model in superconductivity. *European Journal of Applied Mathematics*, 7:237–247, 1996.
- [58] L. Prigozhin. Sandpiles, river networks, and type-ii superconductors. Free Boundary Problems News, 10:2-4, 1996.
- [59] L. Prigozhin. Variational model of sandpile growth. European J. Appl. Math., 7(3):225-235, 1996.
- [60] S. M. Robinson. Stability theory for systems of inequalities. II. Differentiable nonlinear systems. *SIAM J. Numer. Anal.*, 13(4):497–513, 1976.
- [61] J. F. Rodrigues. *Obstacle problems in mathematical physics*, volume 134 of *North-Holland Mathematics Studies*. North-Holland Publishing Co., Amsterdam, 1987. Notas de Matemática [Mathematical Notes], 114.
- [62] J. F. Rodrigues and L. Santos. A parabolic quasi-variational inequality arising in a superconductivity model. *Ann. Scuola Norm. Sup. Pisa Cl. Sci.* (4), 29(1):153–169, 2000.
- [63] T. Roubíček. Nonlinear partial differential equations with applications, volume 153 of International Series of Numerical Mathematics. Birkhäuser Verlag, Basel, 2005.
- [64] H. Scheel and S. Scholtes. Mathematical programs with complementarity constraints: stationarity, optimality, and sensitivity. *Math. Oper. Res.*, 25(1):1–22, 2000.
- [65] A. Schiela and D. Wachsmuth. Convergence analysis of smoothing methods for optimal control of stationary variational inequalities with control constraints. *ESAIM Math. Model. Numer. Anal.*, 47(3):771–787, 2013.

- [66] R. E. Showalter. *Monotone Operators in Banach Space and Nonlinear Partial Differential Equations*. American Mathematical Society, 1997.
- [67] L. Tartar. Inéquations quasi variationnelles abstraites. CR Acad. Sci. Paris Sér. A, 278:1193–1196, 1974.
- [68] F. Tröltzsch. *Optimal control of partial differential equations*, volume 112 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2010. Theory, methods and applications, Translated from the 2005 German original by Jürgen Sprekels.
- [69] G. Wachsmuth. Strong stationarity for optimal control of the obstacle problem with control constraints. *SIAM Journal on Optimization*, 24(4):1914–1932, 2014.
- [70] G. Wachsmuth. Towards M-stationarity for optimal control of the obstacle problem with control constraints. *SIAM J. Control Optim.*, 54(2):964–986, 2016.
- [71] G. Wachsmuth. A guided tour of polyhedric sets: basic properties, new results on intersections, and applications. *J. Convex Anal.*, 26(1):153–188, 2019.
- [72] G. Wachsmuth. Elliptic quasi-variational inequalities under a smallness assumption: uniqueness, differential stability and optimal control. *Calc. Var. Partial Differential Equations*, 59(2):Paper No. 82, 15, 2020.
- [73] E. Zeidler. *Nonlinear Functional Analysis and Applications*, volume II/B: Nonlinear Monotone Operators. Springer-Verlag, 1989.
- [74] J. Zowe and S. Kurcyusz. Regularity and stability for the mathematical programming problem in Banach spaces. *Appl. Math. Optim.*, 5(1):49–62, 1979.