

**Weierstraß-Institut  
für Angewandte Analysis und Stochastik  
Leibniz-Institut im Forschungsverbund Berlin e. V.**

Preprint

ISSN 2198-5855

**State-constrained control-affine parabolic problems I: First and  
second order necessary optimality conditions**

M. Soledad Aronna, <sup>1</sup> J. Frédéric Bonnans, <sup>2</sup> Axel Kröner, <sup>3</sup>

submitted: September 23, 2020

<sup>1</sup> Escola de Matemática Aplicada EMap/FGV  
Rio de Janeiro 22250-900  
Brazil  
E-Mail: soledad.aronna@fgv.br

<sup>2</sup> Inria Saclay and CMAP  
Ecole Polytechnique  
CNRS  
Université Paris Saclay  
91128 Palaiseau  
France  
E-Mail: Frederic.Bonnans@inria.fr

<sup>3</sup> Weierstrass Institute  
Mohrenstr. 39  
10117 Berlin  
Germany  
E-Mail: axel.kroener@wias-berlin.de

No. 2762

Berlin 2020



---

2010 *Mathematics Subject Classification.* 49J20, 49K20, 35J10, 93C20.

*Key words and phrases.* Optimal control of partial differential equations, semilinear parabolic equations, state constraints, second order analysis, control-affine problems.

The first author was supported by FAPERJ, CNPq and CAPES (Brazil) and by the Alexander von Humboldt Foundation (Germany). The second author thanks the ‘Laboratoire de Finance pour les Marchés de l’Énergie’ for its support. The second and third authors were supported by a public grant as part of the Investissement d’Avenir project, reference ANR-11-LABX-0056-LMH, LabEx LMH, in a joint call with Gaspard Monge Program for optimization, operations research and their interactions with data sciences.

This is the first part of a work on optimality conditions for a control problem of a semilinear heat equation. More precisely, the full version, available at <https://arxiv.org/abs/1906.00237v1>, has been divided in two, resulting in the current manuscript (that corresponds to Part I) and <https://arxiv.org/abs/1909.05056> (which is Part II).

Edited by  
Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)  
Leibniz-Institut im Forschungsverbund Berlin e. V.  
Mohrenstraße 39  
10117 Berlin  
Germany

Fax: +49 30 20372-303  
E-Mail: [preprint@wias-berlin.de](mailto:preprint@wias-berlin.de)  
World Wide Web: <http://www.wias-berlin.de/>

# State-constrained control-affine parabolic problems I: First and second order necessary optimality conditions

M. Soledad Aronna, J. Frédéric Bonnans, Axel Kröner,

## Abstract

In this paper we consider an optimal control problem governed by a semilinear heat equation with bilinear control-state terms and subject to control and state constraints. The state constraints are of integral type, the integral being with respect to the space variable. The control is multidimensional. The cost functional is of a tracking type and contains a linear term in the control variables. We derive second order necessary conditions relying on the concept of alternative costates and quasi-radial critical directions.

## 1 Introduction

This is the first part of two papers on necessary and sufficient optimality conditions for an optimal control problem governed by a semilinear heat equation containing bilinear terms coupling the control and the state, and subject to constraints on the control and state. The control may have several components and enters in an affine way in the cost. In this first part we derive necessary optimality conditions of first and second order, in the second part [2] sufficient optimality conditions are shown.

In the context of second order conditions for problems governed by control-affine ordinary differential equations we can mention several works, starting with the early papers [18] by Goh and [19] by Kelley, later [15] by Dmitruk, and recently [1]. In this context, the case dealing with both control and state constraints was treated in e.g. Maurer [25], McDanell and Powers [28], Maurer, Kim and Vossen [27], Schättler [30], and Aronna *et al.* [3]. For a more detailed description of the contributions in this framework, we refer to [3].

In the infinite dimensional case, the issue of second order conditions for problems governed by elliptic equations and assuming state constraints was treated by several authors, see e.g. Casas, Tröltzsch and Unger [12], Bonnans [6], Casas, Mateos and Tröltzsch [11] and Casas and Tröltzsch [13].

Parabolic optimal control problems with state constraints were discussed in several articles. For a semilinear equation in the presence of pure-state constraints, Raymond and Tröltzsch [29], and Krumbiegel and Rehberg [20] obtained second order sufficient conditions. Casas, de Los Reyes, and Tröltzsch [10] and de Los Reyes, Merino, Rehberg and Tröltzsch [14] proved sufficient second order conditions for semilinear equations, both in the elliptic and parabolic cases. The articles mentioned in this paragraph did not consider bilinear terms as we do in the current work.

Further details regarding the existing results on second order analysis of control-affine state-constrained problems are given in the second part [2] of this research.

The contribution of this paper are first and second order necessary optimality conditions for an optimal control problem for a semilinear parabolic equation with cubic nonlinearity, several controls coupled with the state variable through bilinear terms, pointwise control constraints and state constraints that

are integral in space. To incorporate the state constraints we use the concept of *alternative costates* (see Bonnans and Jaisson [8]) and the concept of quasi-radial directions (see Bonnans and Shapiro [9] and Aronna, Bonnans and Goh [3]).

The paper is organized as follows. In Section 2 the problem is stated and main assumptions are formulated. In Section 3 first order analysis is done. Section 4 is devoted to second order necessary conditions. Finally, in the appendix, we give an example satisfying the hypotheses of our main results.

## Notation

Let  $\Omega$  be an open and bounded subset of  $\mathbb{R}^n$ ,  $n \leq 3$ , with  $C^\infty$  boundary  $\partial\Omega$ . Given  $p \in [1, \infty]$  and  $k \in \mathbb{N}$ , let  $W^{k,p}(\Omega)$  be the Sobolev space of functions in  $L^p(\Omega)$  with derivatives (here and after, derivatives w.r.t.  $x \in \Omega$  or w.r.t. time are taken in the sense of distributions) in  $L^p(\Omega)$  up to order  $k$ . Let  $\mathcal{D}(\Omega)$  be the set of  $C^\infty$  functions with compact support in  $\Omega$ . By  $W_0^{k,p}(\Omega)$  we denote the closure of  $\mathcal{D}(\Omega)$  with respect to the  $W^{k,p}$ -topology. Given a horizon  $T > 0$ , we write  $Q := \Omega \times (0, T)$ .  $\|\cdot\|_p$  denotes the norm in  $L^p(0, T)$ ,  $L^p(\Omega)$  and  $L^p(Q)$ , indistinctively. When a function depends on both space and time, but the norm is computed only with respect to one of these variables, we specify both the space and domain. For example, if  $y \in L^p(Q)$  and we fix  $t \in (0, T)$ , we write  $\|y(\cdot, t)\|_{L^p(\Omega)}$ . For the  $p$ -norm in  $\mathbb{R}^m$ , for  $m \in \mathbb{N}$ , we use  $|\cdot|_p$ , for the Euclidean norm we omit the index. We set  $H^k(\Omega) := W^{k,2}(\Omega)$  and  $H_0^k(\Omega) := W_0^{k,2}(\Omega)$ , with dual denoted by  $H^{-k}(\Omega)$ . By  $W^{2,1,p}(Q)$  we mean the Sobolev space of  $L^p(Q)$ -functions whose second derivative in space and first derivative in time belong to  $L^p(Q)$ . For  $p > n + 1$ , we denote by  $Y_p$  the set of elements of  $W^{2,1,p}(Q)$  with zero trace on  $\Sigma$ , and by  $Y_p^0$  its trace at time zero. We write  $H^{2,1}(Q)$  for  $W^{2,1,2}(Q)$  and, setting  $\Sigma := \partial\Omega \times (0, T)$ , we define the state space as

$$Y := \{y \in H^{2,1}(Q); y = 0 \text{ a.e. on } \Sigma\}. \quad (1.1)$$

The latter is continuously embedded in

$$W(0, T) := \{y \in L^2(0, T; H_0^1(\Omega)); \dot{y} \in L^2(0, T; H^{-1}(\Omega))\}. \quad (1.2)$$

Note that if  $y$  is a function over  $Q$ , we use  $\dot{y}$  to denote its time derivative in the sense of distributions. As usual we denote the spatial gradient and the Laplacian by  $\nabla$  and  $\Delta$ . By  $\text{dist}(t, I) := \inf\{\|t - \bar{t}\|; \bar{t} \in I\}$  for  $I \subset \mathbb{R}$ , we denote the distance of  $t$  to the set  $I$ .

## 2 Statement of the problem and main assumptions

In this section we introduce the optimal control problem we deal with and we show well-posedness of the state equation and existence of solutions of the optimal control problem.

### 2.1 Setting

Consider the *state equation*

$$\begin{cases} \dot{y}(x, t) - \Delta y(x, t) + \gamma y^3(x, t) = f(x, t) + y(x, t) \sum_{i=0}^m u_i(t) b_i(x) & \text{in } Q, \\ y = 0 & \text{on } \Sigma, \quad y(\cdot, 0) = y_0 \text{ in } \Omega, \end{cases} \quad (2.1)$$

and

$$y_0 \in H_0^1(\Omega), \quad f \in L^2(Q), \quad b \in L^\infty(\Omega)^{m+1}, \quad (2.2)$$

$\gamma \geq 0$ ,  $u_0 \equiv 1$  is a constant, and  $u := (u_1, \dots, u_m) \in L^2(0, T)^m$ . Lemma 2.3 below shows that for each control  $u \in L^2(0, T)^m$ , there is a unique associated solution  $y \in Y$  of (2.1), called the *associated state*. Let  $y[u]$  denote this solution. We consider control constraints of the form  $u \in \mathcal{U}_{\text{ad}}$ , where

$$\mathcal{U}_{\text{ad}} \text{ is a nonempty, closed convex subset of } L^2(0, T)^m. \quad (2.3)$$

In some statements, we will consider a specific form of  $\mathcal{U}_{\text{ad}}$  (see (3.26) below). In addition, we have finitely many linear running state constraints of the form

$$g_j(y(\cdot, t)) := \int_{\Omega} c_j(x)y(x, t)dx + d_j \leq 0, \quad \text{for } t \in [0, T], \quad j = 1, \dots, q, \quad (2.4)$$

where  $c_j \in H^2(\Omega) \cap H_0^1(\Omega)$  for  $j = 1, \dots, q$ , and  $d \in \mathbb{R}^q$ . The  $H_0^1(\Omega)$  regularity of  $c$  is used in Lemma 3.2 to derive regularity results for the adjoint state and the  $H^2(\Omega)$  regularity in Proposition 3.11 for results on the Lagrange multiplier associated with the state constraint.

We call any  $(u, y[u]) \in L^2(0, T)^m \times Y$  a *trajectory*, and if it additionally satisfies the control and state constraints, we say it is an *admissible trajectory*. The *cost function* is

$$\begin{aligned} J(u, y) := & \frac{1}{2} \int_Q (y(x, t) - y_d(x))^2 dx dt \\ & + \frac{1}{2} \int_{\Omega} (y(x, T) - y_{dT}(x))^2 dx + \sum_{i=1}^m \alpha_i \int_0^T u_i(t) dt, \end{aligned} \quad (2.5)$$

where

$$y_d \in L^2(Q), \quad y_{dT} \in H_0^1(\Omega), \quad (2.6)$$

and  $\alpha \in \mathbb{R}^m$ . We consider the optimal control problem

$$\text{Min}_{u \in \mathcal{U}_{\text{ad}}} J(u, y[u]); \quad \text{subject to (2.4)}. \quad (\text{P})$$

For problem (P) we consider the two types of solution given next.

**Definition 2.1.** Let  $\bar{u} \in \mathcal{U}_{\text{ad}}$ . We say that  $(\bar{u}, y[\bar{u}])$  is an  $L^2$ -local solution (resp.,  $L^\infty$ -local solution) if there exists  $\varepsilon > 0$  such that  $(\bar{u}, y[\bar{u}])$  is a minimum among the admissible trajectories  $(u, y)$  that satisfy  $\|u - \bar{u}\|_2 < \varepsilon$  (resp.,  $\|u - \bar{u}\|_\infty < \varepsilon$ ).

## 2.2 Well-posedness of the state equation

Here we study the state equation and analyze, by means of the Implicit Function Theorem, the *control-to-state mapping*, i.e. the mapping that associates to each control, the corresponding solution of the state equation. We start by the following easily checked technical result.

**Lemma 2.2.** For  $i = 0, \dots, m$ , the mapping defined on  $L^2(0, T) \times L^\infty(\Omega) \times L^\infty(0, T; L^2(\Omega))$ , given by  $(u_i, b_i, y) \mapsto u_i b_i y$ , has image in  $L^2(Q)$ , is of class  $C^\infty$ , and satisfies

$$\|u_i b_i y\|_2 \leq \|u_i\|_2 \|b_i\|_\infty \|y\|_{L^\infty(0, T; L^2(\Omega))}. \quad (2.7)$$

A uniqueness and existence result, and *a priori* estimates for the state follows.

**Lemma 2.3.** *The state equation (2.1) has a unique solution  $y = y[u, y_0, f]$  in  $Y$ . The mapping  $(u, y_0, f) \mapsto y[u, y_0, f]$  is  $C^\infty$  from  $L^2(0, T)^m \times H_0^1(\Omega) \times L^2(Q)$  to  $Y$ , and nondecreasing w.r.t.  $y_0$  and  $f$ . In addition, there exist functions  $C_i$ ,  $i = 1$  to  $2$ , not decreasing w.r.t. each component, such that*

$$\|y\|_{L^\infty(0, T; L^2(\Omega))} + \|\nabla y\|_2 \leq C_1(\|y_0\|_2, \|f\|_2, \|u\|_2 \|b\|_\infty), \quad (2.8)$$

$$\|y\|_Y \leq C_2(\|y_0\|_{H_0^1(\Omega)}, \|f\|_2, \|u\|_2 \|b\|_\infty). \quad (2.9)$$

Moreover, the state  $y$  also belongs to  $C([0, T]; H_0^1(\Omega))$ , since  $Y$  is continuously embedded in that space [24, Theorem 3.1, p.23].

In the proof that follows, we use several times the (continuous) Sobolev inclusion

$$H_0^1(\Omega) \subset L^6(\Omega), \quad \text{when } n \leq 3. \quad (2.10)$$

*Proof.* (i) Observe first that by the standard Sobolev inclusions and Lemma 2.2, any  $y \in Y$  is such that  $y^3$  and  $y \sum_{i=0}^m u_i b_i$  belong to  $L^2(Q)$ . So,  $\dot{y} - \Delta y \in L^2(Q)$  and, therefore, the notion of solution of the state equation in  $Y$  is clear. We could as well define a solution in  $W(0, T)$  but since by (2.10), for  $n \leq 3$ ,  $W(0, T) \subset L^2(0, T; L^6(\Omega))$ , and the compatibility condition (equality between the trace of the initial condition on  $\partial\Omega$  and the Dirichlet condition on  $\Sigma$ ) holds, it follows then that any solution in  $W(0, T)$  is a solution in  $Y$ .

(ii) We establish the *a priori* estimates (2.8)-(2.9). Multiplying the state equation by  $y$  and integrating over  $\Omega$ , we get

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \int_{\Omega} y(x, t)^2 dx + \int_{\Omega} |\nabla y(x, t)|^2 dx + \gamma \int_{\Omega} y(x, t)^4 dx \\ \leq \frac{1}{2} \int_{\Omega} f(x, t)^2 dx + \left(\frac{1}{2} + |u(t)|_1 \|b\|_\infty\right) \int_{\Omega} y(x, t)^2 dx. \end{aligned} \quad (2.11)$$

In particular,  $\eta(t) := \int_{\Omega} y(x, t)^2 dx$  satisfies

$$\dot{\eta}(t) \leq \int_{\Omega} f(x, t)^2 dx + (1 + 2|u(t)|_1 \|b\|_\infty) \eta(t). \quad (2.12)$$

By Gronwall's Lemma:

$$\|\eta\|_\infty \leq (\|y_0\|_2^2 + \|f\|_2^2) e^{T+2\|u\|_1 \|b\|_\infty} \quad (2.13)$$

and then (2.8) easily follows.

Now multiplying the state equation by  $\dot{y}$  we get, for all  $\varepsilon > 0$ ,

$$\begin{aligned} \int_{\Omega} \dot{y}(x, t)^2 dx + \frac{1}{2} \frac{d}{dt} \int_{\Omega} |\nabla y(x, t)|^2 dx + \frac{\gamma}{4} \frac{d}{dt} \int_{\Omega} y(x, t)^4 dx \\ \leq \frac{1}{\varepsilon} \int_{\Omega} f(x, t)^2 dx + \frac{1}{\varepsilon} |u(t)|^2 \|b\|_\infty^2 \int_{\Omega} y(x, t)^2 dx + \frac{\varepsilon}{2} \int_{\Omega} \dot{y}(x, t)^2 dx. \end{aligned} \quad (2.14)$$

Choosing  $\varepsilon = 1$  we get, after cancellation,

$$\begin{aligned} \int_{\Omega} \dot{y}(x, t)^2 dx + \frac{d}{dt} \int_{\Omega} |\nabla y(x, t)|^2 dx + \frac{\gamma}{2} \frac{d}{dt} \int_{\Omega} y(x, t)^4 dx \\ \leq 2 \int_{\Omega} f(x, t)^2 dx + 2|u(t)|^2 \|b\|_\infty^2 \int_{\Omega} y(x, t)^2 dx. \end{aligned} \quad (2.15)$$

For  $\tau \in [0, T]$ , integrating from 0 to  $\tau$ , and using (2.10), we obtain that

$$\|y\|_{H^1(0,T;L^2(\Omega))} + \|\nabla y\|_{L^\infty(0,T;L^2(\Omega))} \leq C_2(\|y_0\|_{H_0^1(\Omega)}, \|f\|_2, \|u\|_2 \|b\|_\infty). \quad (2.16)$$

We easily deduce (2.9) since we can estimate  $\|\Delta y\|_{L^2(Q)}$  and, therefore, also  $\|y\|_{L^2(0,T;H^2(\Omega))}$  with the previous relations.

(iii) We construct a sequence  $y_k$  of Galerkin approximations for which estimates analogous to (2.8) hold. Some subsequence weakly converges in  $W(0, T)$  to some  $y$  and is such that the sequence  $y_k^3$ , bounded in  $L^2(Q)$ , weakly converges in this space. By the Aubin-Lions lemma [4], the injection of  $W(0, T)$  into  $L^2(Q)$  is compact. So (extracting again a subsequence if necessary),  $y_k^3$  converges a.e. to  $y^3$ . By Lions [22, Lem. 1.3, p. 12], the weak limit of  $y_k^3$  is  $y^3$ , and  $y$  is therefore solution of the state equation.

(iv) The  $C^\infty$  regularity of  $y[u, y_0, f]$  is a consequence of the Implicit Function Theorem. In fact, let  $Y^0$  denote the trace at time 0 of elements of  $Y$ , which with the trace norm is a Banach space containing  $H_0^1(\Omega)$ . Then the mapping  $F : L^2(0, T) \times Y \times Y^0 \times L^2(Q) \rightarrow L^2(Q) \times Y^0$  defined by

$$F(u, y, y_0, f) := \left( \dot{y} - \Delta y + \gamma y^3 - y \sum_{i=1}^m u_i b_i, y(0) - y_0 \right), \quad (2.17)$$

is of class  $C^\infty$ . That the linearized mapping  $D_y F$  is bijective follows from results already shown in this proof.

(v) Uniqueness follows from the monotonicity w.r.t.  $(y_0, f)$ , that we prove as follows. Consider the difference  $z := y_2 - y_1$  of two solutions  $y_1$  and  $y_2$  of (2.1), with data  $(y_{01}, f_1) \leq (y_{02}, f_2)$ , resp. By the Mean Value Theorem,  $z$  is solution of

$$\dot{z} - \Delta z + z \sum_{i=1}^m u_i b_i + 3\gamma \hat{y}^2 z = \tilde{f}; \quad z(0) = \tilde{y}_0 \quad (2.18)$$

where  $\hat{y} \in [y_1, y_2]$  a.e.,  $\tilde{y}_0 := y_{02} - y_{01} \leq 0$  and  $\tilde{f} := f_2 - f_1 \leq 0$ . Testing the equation with  $z_+ := \max(z, 0)$  we get that  $\nu(t) := \int_\Omega z_+^2$  satisfies

$$\frac{1}{2} \dot{\nu} - |u(t)| \|b\|_\infty \nu(t) \leq \frac{1}{2} \dot{\nu} + \int_\Omega z_+^2 \sum_{i=1}^m u_i b_i \leq \int_\Omega \tilde{f} z_+ \leq 0 \quad (2.19)$$

and applying Gronwall's inequality we obtain that  $z_+ = 0$ .  $\square$

In the analysis that follows, we fix a trajectory  $(\bar{u}, \bar{y} = y[\bar{u}])$ .

For this trajectory  $(\bar{u}, \bar{y})$ , let us consider the linear continuous operator  $A$  from  $L^2(0, T; H^2(\Omega))$  to  $L^2(Q)$  such that, for each  $z \in Y$  and  $(x, t) \in Q$ ,

$$(Az)(x, t) := -\Delta z(x, t) + 3\gamma \bar{y}(x, t)^2 z(x, t) - \sum_{i=0}^m \bar{u}_i(t) b_i(x) z(x, t). \quad (2.20)$$

**Lemma 2.4.** *For any  $\bar{f} \in L^2(Q)$ , the equation*

$$\begin{cases} \dot{z} + Az = \bar{f}, & \text{in } Q, \\ z = 0 \text{ on } \Sigma, \quad z(x, 0) = 0 \text{ in } \Omega, \end{cases} \quad (2.21)$$

has a unique solution  $z \in Y$  that verifies

$$\|z\|_{L^\infty(0,T;L^2(\Omega))} \leq e^{\frac{1}{2}T + \sum_{i=0}^m \|\bar{u}_i\|_1 \|b_i\|_\infty} \|\bar{f}\|_{L^2(0,T;L^2(\Omega))}. \quad (2.22)$$

*Proof.* We follow the same method used in Lemma 2.3. Multiplying (2.21) by  $z(x, t)$  and integrating over space we obtain that for a.a.  $t \in (0, T)$

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|z(\cdot, t)\|_{L^2(\Omega)}^2 + \|\nabla z(\cdot, t)\|_{L^2(\Omega)}^2 + 3\gamma \|\bar{y}(\cdot, t)z(\cdot, t)\|_{L^2(\Omega)}^2 \\ = \int_{\Omega} z(x, t) \left( \bar{f}(x, t) + \sum_{i=0}^m \bar{u}_i(t) \cdot b_i(x)z(x, t) \right) dx. \end{aligned} \quad (2.23)$$

The r.h.s. of (2.23) can be bounded above by

$$\|\bar{f}(\cdot, t)\|_{L^2(\Omega)}^2 + \left( \frac{1}{2} + \sum_{i=0}^m |\bar{u}_i| \|b_i\|_{\infty} \right) \|z(\cdot, t)\|_{L^2(\Omega)}^2. \quad (2.24)$$

Then we deduce the estimate (2.22) with Gronwall's Lemma.  $\square$

## 2.3 Existence of solution of the optimal control problem

In order to study the existence of local solutions, we need to establish the sequential weak continuity of the control-to-state mapping. We use ' $\rightharpoonup$ ' to denote the weak convergence of a sequence, the space being indicated in each case. We need the following result (see [23, p. 14]):

$$\begin{cases} \text{For any } p \in [1, 10), \text{ the following injection is compact:} \\ Y \hookrightarrow L^p(0, T; L^{10}(\Omega)), \text{ when } n \leq 3. \end{cases} \quad (2.25)$$

**Lemma 2.5.** *The mapping  $u \mapsto y[u]$  is sequentially weakly continuous from  $L^2(0, T)^m$  into  $Y$ .*

*Proof.* Taking  $u_\ell \rightharpoonup \bar{u}$  in  $L^2(0, T)^m$ , we shall prove that  $y_\ell \rightharpoonup \bar{y}$  in  $Y$ , where  $y_\ell := y[u_\ell]$ ,  $\bar{y} := y[\bar{u}]$ . We know that it is enough to check that any subsequence of  $y_\ell$  weakly converges to  $\bar{y}$  in  $Y$ . To do this, we prove that we can pass to the limit in each term of the state equation.

(a) We know by Lemma 2.3 that  $y_\ell$  is bounded in  $Y$ , so extracting a subsequence if necessary, we may assume that it weakly converges in  $Y$  to some  $\hat{y}$ . By (2.25),  $y_\ell \rightharpoonup \hat{y}$  in  $L^6(Q)$  and, therefore, maybe for a subsequence, it converges almost everywhere in  $Q$ .

Let  $\nu \in [2, 5]$  be integer. Set  $\sigma := 6/\nu$ . By the mean value theorem,  $y_\ell^\nu - \hat{y}^\nu = \nu \tilde{y}_\ell^{\nu-1} (y_\ell - \hat{y})$ , with  $\tilde{y}_\ell(x, t) \in [y_\ell(x, t), \hat{y}(x, t)]$  a.e. Obviously  $\tilde{y}_\ell$  is measurable and bounded in  $L^6(Q)$ . By Hölder's inequality, with  $p = \nu/(\nu - 1)$  and  $q = 6/\sigma = \nu$  (note that  $1/p + 1/q = 1$ ), we get

$$\begin{aligned} \frac{1}{\nu^\sigma} \|y_\ell^\nu - \hat{y}^\nu\|_\sigma^\sigma &= \int_Q \tilde{y}_\ell^{\sigma(\nu-1)} (y_\ell - \hat{y})^\sigma dx dt \leq \|\tilde{y}_\ell^{\sigma(\nu-1)}\|_p \| (y_\ell - \hat{y})^\sigma \|_q \\ &= \|\tilde{y}_\ell\|_6^{\sigma(\nu-1)} \|y_\ell - \hat{y}\|_6^\sigma. \end{aligned} \quad (2.26)$$

Therefore,  $y_\ell^\nu \rightarrow \hat{y}^\nu$  in  $L^\sigma(Q)$ . Taking  $\nu = 3$  we get the desired result.

(b) We claim that  $u_\ell y_\ell b$  weakly converges in  $L^2(Q)$  to  $\bar{u} \hat{y} b$ . It is enough to get the result when  $m = 1$ . Fix  $\varphi$  in  $L^\infty(Q)$ . By Lemma 2.2,  $u_\ell y_\ell$  is bounded in  $L^2(Q)$  and has therefore (up to a subsequence) a weak limit  $w$  in that space. Since  $y_\ell \rightarrow \hat{y}$  in  $L^6(Q)$ ,  $\int_Q u_\ell (y_\ell - \hat{y}) b \varphi \rightarrow 0$ . On the other hand  $\int_Q u_\ell \hat{y} b \varphi \rightarrow \int_Q \bar{u} \hat{y} b \varphi$  since  $\hat{y} b \varphi \in L^2(Q)$ . Therefore  $\int_Q u_\ell y_\ell b \varphi \rightarrow \int_Q \bar{u} \hat{y} b \varphi$ . Since  $L^\infty(Q)$  is a dense subset of  $L^2(Q)$ . The claim follows.

By steps (a)-(b), we can pass to the limit in the weak formulation, and obtain (due to the uniqueness of solution) that  $\hat{y} = \bar{y}$ . The conclusion follows.  $\square$



**Theorem 2.6.** (i) The function  $u \mapsto J(u, y[u])$ , from  $L^2(0, T)^m$  to  $\mathbb{R}$ , is weakly sequentially l.s.c.  
(ii) The set of solutions of the optimal control problem (P) is weakly sequentially closed in  $L^2(0, T)^m$ .  
(iii) If (P) has a bounded minimizing sequence, the set of solutions of (P) is non empty. This is the case in particular if (P) is admissible and  $\mathcal{U}_{\text{ad}}$  is a nonempty, bounded subset of  $L^2(0, T)^m$ .

*Proof.* (i) Combine Lemma 2.5 and the fact that the cost function  $J$  is continuous and convex on  $L^2(0, T)^m \times Y$ , hence it is also weakly lower semicontinuous over this product space.

(ii) Let  $(u_\ell) \subset L^2(0, T)^m$  be a sequence of solutions weakly converging to  $\bar{u} \in L^2(0, T)^m$ , with associated states  $y_\ell$ . By Lemma 2.5,  $(y_\ell)$  weakly converge in  $Y$  to the state  $\bar{y}$  associated with  $\bar{u}$  and, by point (i),  $J(\bar{u}, \bar{y}) \leq \liminf_\ell J(u_\ell, y_\ell)$ . This lower limit being nothing but the value of problem (P), the conclusion follows.

(iii) By the previous arguments, a weak limit of a minimizing sequence is a solution of (P). This weak limit exists iff the sequence is bounded. This concludes the proof.  $\square$

### 3 First order analysis

In this section we state first order necessary optimality conditions. More precisely, we introduce the adjoint equation, and define and prove existence of associated Lagrange multipliers.

Throughout the section,  $(\bar{u}, \bar{y})$  is a trajectory of problem (P). We recall the hypotheses (2.2), (2.6) on the data, and the definition of the operator  $A$  given in (2.20).

#### 3.1 Linearized state equation and costate equation

The *linearized state equation* at  $(\bar{u}, \bar{y})$  is given by

$$\begin{cases} \dot{z} + Az = \sum_{i=1}^m v_i b_i \bar{y} & \text{in } Q; \\ z = 0 \text{ on } \Sigma, \quad z(\cdot, 0) = 0 \text{ on } \Omega, \end{cases} \quad (3.1)$$

For  $v \in L^2(0, T)^m$ , equation (3.1) above possesses a unique solution  $z[v] \in Y$  (as follows from Lemma 2.4), and the mapping  $v \mapsto z[v]$  is linear and continuous from  $L^2(0, T)^m$  to  $Y$ . Particularly, the following estimate holds.

**Proposition 3.1.** *One has*

$$\|z\|_{L^\infty(0, T; L^2(\Omega))} \leq M_1 \sum_{i=1}^m \|b_i\|_\infty \|v_i\|_1, \quad (3.2)$$

where  $M_1 := e^{\frac{T}{2} + \sum_{i=0}^m \|\bar{u}_i\|_1 \|b_i\|_\infty} \|\bar{y}\|_{L^\infty(0, T; L^2(\Omega))}$ .

*Proof.* Immediate consequence of Lemma 2.4.  $\square$

It is well-known that the dual of  $C([0, T])$  is the set of (finite) Radon measures, and that the action of a finite Radon measure coincides with the Stieltjes integral associated with a bounded variation function  $\mu \in BV(0, T)$ . We may assume w.l.g. that  $\mu(T) = 0$ , and we let  $d\mu$  denote the Radon measure

associated to  $\mu$ . Note that if  $d\mu$  belongs to the set  $\mathcal{M}_+(0, T)$  of nonnegative finite Radon measures then we may take  $\mu$  nondecreasing. Set

$$BV(0, T)_{0,+} := \{\mu \in BV(0, T) \text{ nondecreasing, right-continuous; } \mu(T) = 0\}. \quad (3.3)$$

The *generalized Lagrangian* of problem (P) is, choosing the multiplier of the state equation to be  $(p, p_0) \in L^2(Q) \times H^{-1}(\Omega)$  and taking  $\beta \in \mathbb{R}_+$ ,  $\mu \in BV(0, T)_{0,+}^q$ ,

$$\begin{aligned} \mathcal{L}[\beta, p, p_0, \mu](u, y) &:= \beta J(u, y) - \langle p_0, y(\cdot, 0) - y_0 \rangle_{H_0^1(\Omega)} \\ &+ \int_Q p \left( \Delta y(x, t) - \gamma y^3(x, t) + f(x, t) + \sum_{i=0}^m u_i(t) b_i(x) y(x, t) - \dot{y}(x, t) \right) dx dt \\ &+ \sum_{j=1}^q \int_0^T g_j(y(\cdot, t)) d\mu_j(t). \end{aligned} \quad (3.4)$$

The *costate equation* is the condition of stationarity of the Lagrangian  $\mathcal{L}$  with respect to the state that is, for any  $z \in Y$ :

$$\begin{aligned} \int_Q p(\dot{z} + Az) dx dt + \langle p_0, z(\cdot, 0) \rangle_{H_0^1(\Omega)} &= \sum_{j=1}^q \int_0^T \int_{\Omega} c_j z dx d\mu_j(t) \\ &+ \beta \int_Q (\bar{y} - y_d) z dx dt + \beta \int_{\Omega} (\bar{y}(x, T) - y_{dT}(x)) z(x, T) dx. \end{aligned} \quad (3.5)$$

To each  $(\varphi, \psi) \in L^2(Q) \times H_0^1(\Omega)$ , let us associate  $z = z[\varphi, \psi] \in Y$ , the unique solution of

$$\dot{z} + Az = \varphi; \quad z(\cdot, 0) = \psi. \quad (3.6)$$

Since this mapping is onto, the costate equation (3.5) can be rewritten, for  $z = z[\varphi, \psi]$  and arbitrary  $(\varphi, \psi) \in L^2(Q) \times H_0^1(\Omega)$ , as

$$\begin{aligned} \int_Q p \varphi dx dt + \langle p_0, \psi \rangle_{H_0^1(\Omega)} &= \sum_{j=1}^q \int_0^T \int_{\Omega} c_j z dx d\mu_j(t), \\ &+ \beta \int_Q (\bar{y} - y_d) z dx dt + \beta \int_{\Omega} (\bar{y}(x, T) - y_{dT}(x)) z(x, T) dx. \end{aligned} \quad (3.7)$$

The r.h.s. of (3.7) can be seen as a linear continuous form on the pairs  $(\varphi, \psi)$  of the space  $L^2(Q) \times H_0^1(\Omega)$ . By the Riesz Representation Theorem, there exists a unique  $(p, p_0) \in L^2(Q) \times H^{-1}(\Omega)$  satisfying (3.7), that means, there is a unique solution of the costate equation.

Next consider the *alternative costates*

$$p^1 := p + \sum_{j=1}^q c_j \mu_j; \quad p_0^1 := p_0 + \sum_{j=1}^q c_j \mu_j(0). \quad (3.8)$$

**Lemma 3.2.** *Let  $(p, p_0, \mu) \in L^2(Q) \times H^{-1}(\Omega) \times BV(0, T)_{0,+}^q$  satisfy (3.7), let  $(p^1, p_0^1)$  be given by (3.8). Then  $p^1 \in Y$ , it satisfies  $p^1(0) = p_0^1$ , and it is the unique solution of*

$$-p^1 + Ap^1 = \beta(\bar{y} - y_d) + \sum_{j=1}^q \mu_j A c_j, \quad p^1(\cdot, T) = \beta(\bar{y}(\cdot, T) - y_{dT}). \quad (3.9)$$

Moreover,  $p(x, 0)$  and  $p(x, T)$  are well-defined as elements of  $H_0^1(\Omega)$  in view of (3.8), and we have

$$p(\cdot, 0) = p_0, \quad p(\cdot, T) = \beta(\bar{y}(\cdot, T) - y_{dT}). \quad (3.10)$$

*Proof.* Let  $z \in Y$ . Note that, for  $1 \leq j \leq q$ , the function  $t \mapsto \int_{\Omega} c_j(x)z(x, t)dx$ , belongs to  $W^{1,1}(0, T)$  and is, therefore, of bounded variation. Using the integration by parts formula for the product of scalar functions with bounded variation, one of them being continuous (see e.g. [8, Lemma 3.6]), and taking into account the fact that  $\mu_j(T) = 0$ , we get that, for  $\psi = z(\cdot, 0)$ ,

$$\sum_{j=1}^q \int_Q c_j \mu_j \dot{z} dx dt + \sum_{j=1}^q \mu_j(0) \langle c_j, \psi \rangle_{L^2(\Omega)} = - \sum_{j=1}^q \int_0^T \int_{\Omega} c_j z dx d\mu_j(t). \quad (3.11)$$

By the definition (3.8) of the alternative costate, the latter equation can be rewritten as

$$\int_Q (p^1 - p) \dot{z} dx dt + \langle p_0^1 - p_0, \psi \rangle_{H_0^1(\Omega)} = - \sum_{j=1}^q \int_0^T \int_{\Omega} c_j z dx d\mu_j(t). \quad (3.12)$$

Now adding (3.7) and (3.12), as well as the identity

$$\int_Q (p^1 - p) Az = \int_Q \sum_{j=1}^q c_j \mu_j Az \quad (3.13)$$

we obtain, since  $\varphi = \dot{z} + Az$ , that (implicitly identifying, as usual,  $L^2(\Omega)$  with its dual)

$$\begin{aligned} & \int_Q p^1 \varphi dx dt + \langle p_0^1, \psi \rangle_{H_0^1(\Omega)} \\ &= \beta \int_Q (\bar{y} - y_d) z dx dt + \beta \int_{\Omega} (\bar{y}(x, T) - y_{dT}(x)) z(x, T) dx + \int_Q \sum_{j=1}^q c_j \mu_j Az. \end{aligned} \quad (3.14)$$

Since  $A$  is symmetric, using (2.6), we see that  $p^1$  is solution in  $Y$  of (3.9); the solution of the latter being clearly unique. Multiplying (3.9) by  $z \in Y$  and integrating over  $Q$ , with an integration by parts of the term with  $\dot{p}^1 z$ , we recover (using (3.8)) equation (3.14) implying that  $p^1(x, 0) = p_0^1(x)$  for a.a.  $x$  in  $\Omega$ . Conversely, it is easy to prove that any solution of (3.14) is solution of (3.9).

Since  $p^1$  and  $c_j \mu_j$  belong to  $L^\infty(0, T; H_0^1(\Omega))$ , by (3.8) also  $p$  has this regularity. Use (3.8) again, the final condition on  $p^1$  and the fact that  $\mu(T) = 0$  to get the second relation of (3.10). Furthermore, we have

$$p_0 = p^1(\cdot, 0) - \sum_{j=1}^q c_j \mu_j(0) = p(\cdot, 0). \quad (3.15)$$

□

**Corollary 3.3.** *If  $\mu \in H^1(0, T)^q$ , then  $p \in Y$  and*

$$-\dot{p} + Ap = \beta(\bar{y} - y_d) + \sum_{j=1}^q c_j \dot{\mu}_j. \quad (3.16)$$

*Proof.* This follows immediately from (3.8) and (3.9). □

### 3.2 First order optimality conditions

Let  $(\bar{u}, \bar{y})$  be an admissible trajectory of problem  $(P)$ . We say that  $\mu \in BV(0, T)_{0,+}^q$  is *complementary to the state constraint* for  $\bar{y}$  if

$$\int_0^T g_j(\bar{y}(\cdot, t)) d\mu_j(t) = \int_0^T \left( \int_{\Omega} c_j(x) \bar{y}(x, t) dx + d_j \right) d\mu_j(t) = 0, \quad j = 1, \dots, q. \quad (3.17)$$

Let  $(\beta, \mu) \in \mathbb{R}_+ \times BV(0, T)_{0,+}^q$ . We say that  $p \in L^\infty(0, T; H_0^1(\Omega))$  is the *costate associated with*  $(\bar{u}, \bar{y}, \beta, \mu)$ , or shortly to  $(\beta, \mu)$ , if it is the unique solution of (3.5) with  $p_0 = p(\cdot, 0)$ .

**Definition 3.4.** We say that the triple  $(\beta, p, \mu) \in \mathbb{R}_+ \times L^\infty(0, T; H_0^1(\Omega)) \times BV(0, T)_{0,+}^q$  is a generalized Lagrange multiplier if it satisfies the following first-order optimality conditions:  $\mu$  is complementary to the state constraint,  $p$  is the costate associated with  $(\beta, \mu)$ , the non-triviality condition

$$(\beta, d\mu) \neq 0, \quad (3.18)$$

holds and, for  $i = 1$  to  $m$ , defining the switching function by

$$\Psi_i^p(t) := \beta \alpha_i + \int_{\Omega} b_i(x) \bar{y}(x, t) p(x, t) dx, \quad \text{for } i = 1, \dots, m, \quad (3.19)$$

one has  $\Psi^p \in L^\infty(0, T)^m$  and

$$\sum_{i=1}^m \int_0^T \Psi_i^p(t) (u_i(t) - \bar{u}_i(t)) dt \geq 0, \quad \text{for every } u \in \mathcal{U}_{\text{ad}}. \quad (3.20)$$

We let  $\Lambda(\bar{u}, \bar{y})$  denote the set of generalized Lagrange multipliers  $(\beta, p, \mu)$  associated with  $(\bar{u}, \bar{y})$ . If  $\beta = 0$  we say that the corresponding multiplier is singular. Finally, we write  $\Lambda_1(\bar{u}, \bar{y})$  for the set of pairs  $(p, \mu)$  with  $(1, p, \mu) \in \Lambda(\bar{u}, \bar{y})$ . When the nominal solution is fixed and there is no place for confusion, we just write  $\Lambda$  and  $\Lambda_1$ .

Note that, in view of (3.10),  $p_0 = p(\cdot, 0)$  and hence we do not need to consider  $p_0$  as a component of the multiplier.

#### 3.2.1 The reduced abstract problem

Set  $F(u) := J(u, y[u])$ , and  $G : L^2(0, T)^m \rightarrow C([0, T])^q$ ,  $G(u) := g(y[u])$ . The *reduced problem* is

$$\text{Min}_{u \in \mathcal{U}_{\text{ad}}} F(u); \quad G(u) \in K, \quad (\text{RP})$$

where  $K := C([0, T])_+^q$  is the closed convex cone of continuous functions over  $[0, T]$ , with values in  $\mathbb{R}_+^q$ . Its interior is the set of functions in  $C([0, T])^q$  with negative values. We say that the reduced problem (RP) is *qualified* at  $\bar{u}$  if:

$$\left\{ \begin{array}{l} \text{there exists } u \in \mathcal{U}_{\text{ad}} \text{ such that } v := u - \bar{u} \text{ satisfies} \\ G(\bar{u}) + DG(\bar{u})v \in \text{int}(K). \end{array} \right. \quad (3.21)$$

Given a Banach space  $X$ , a closed convex subset  $S \subseteq X$  and a point  $\bar{s} \in S$ , the *normal cone* to  $S$  at  $\bar{s}$  is defined as

$$N_S(\bar{s}) := \{x^* \in X^*; \langle x^*, s - \bar{s} \rangle \leq 0, \text{ for all } s \in S\}. \quad (3.22)$$

We get the following first order conditions for our problem  $(P)$ :

**Lemma 3.5.** (i) If  $(\bar{u}, y[\bar{u}])$  is an  $L^2$ -local solution of (P), then the associated set  $\Lambda$  of multipliers is nonempty.

(ii) If in addition the qualification condition (3.21) holds at  $\bar{u}$ , then there is no singular multiplier, and  $\Lambda_1$  is bounded in  $L^\infty(0, T; H_0^1(\Omega)) \times BV(0, T)_{0,+}^q$ .

*Proof.* (i) Let us consider the generalized Lagrangian associated with the reduced problem (RP):

$$L[\beta, \mu](u) := \beta F(u) + \sum_{j=1}^q \int_0^T G_j(u)(t) d\mu_j(t). \quad (3.23)$$

Let  $\bar{u}$  be a local solution of (RP). By, e.g., [9, Proposition 3.18], since  $K$  has nonempty interior, there exists a generalized Lagrange multiplier associated with problem (RP), that is,  $(\beta, d\mu) \in \mathbb{R}_+ \times N_K(G(\bar{u}))$  for  $\mu \in BV(0, T)_{0,+}^q$  such that

$$(\beta, d\mu) \neq 0 \quad \text{and} \quad -D_u L[\beta, \mu](\bar{u}) \in N_{\mathcal{U}_{\text{ad}}}(\bar{u}). \quad (3.24)$$

Due to the costate equation (3.7), the latter condition is equivalent to (3.20).

(ii) That  $\Lambda_1$  is nonempty and weakly-\* compact follows from [9, Proposition 3.16].  $\square$

Observe that the qualification condition for (RP) given in (3.21) holds if and only if the following qualification condition for the original problem (P) is satisfied:

$$\begin{cases} \text{there exists } \varepsilon > 0 \text{ and } u \in \mathcal{U}_{\text{ad}} \text{ such that } v := u - \bar{u} \text{ satisfies} \\ g_j(\bar{y}(\cdot, t)) + g'_j(\bar{y}(\cdot, t))z[v](\cdot, t) < -\varepsilon, \text{ for all } t \in [0, T], \text{ and } j = 1, \dots, q. \end{cases} \quad (3.25)$$

In view of Lemma 3.5, if (3.25) is satisfied, then  $\Lambda_1$  is nonempty and weakly-\* compact.

In the sequel of this section, we consider  $(\bar{u}, \bar{y}, \beta, p, \mu)$ , with  $\bar{y}$  the state associated with the admissible control  $\bar{u}$  and  $(\beta, p, \mu) \in \Lambda$ .

### 3.3 Arcs and junction points

We assume in the remainder of the article that the admissible set of controls has the form

$$\mathcal{U}_{\text{ad}} = \{u \in L^2(0, T)^m; \check{u}_i \leq u_i(t) \leq \hat{u}_i, i = 1, \dots, m\}, \quad (3.26)$$

for some constants  $\check{u}_i < \hat{u}_i$ , for  $i = 1, \dots, m$ . Consider the *contact sets associated to the control bounds* defined, up to null measure sets, by

$$\check{I}_i := \{t \in [0, T]; \bar{u}_i(t) = \check{u}_i\}, \quad \hat{I}_i := \{t \in [0, T]; \bar{u}_i(t) = \hat{u}_i\}, \quad I_i := \check{I}_i \cup \hat{I}_i. \quad (3.27)$$

For  $j = 1, \dots, q$ , the *contact set associated with the  $j$ th state constraint* is

$$I_j^C := \{t \in [0, T]; g_j(\bar{y}(\cdot, t)) = 0\}. \quad (3.28)$$

Given  $0 \leq a < b \leq T$ , we say that  $(a, b)$  is a *maximal state constrained arc* for the  $j$ th state constraints, if  $I_j^C$  contains  $(a, b)$  but it contains no open interval strictly containing  $(a, b)$ . We define in the same way a *maximal (lower or upper) control bound constraints arc* (having in mind that the latter are defined up to a null measure set).

We will assume the following *finite arc property*:

$$\left\{ \begin{array}{l} \text{the contact sets for the state and bound constraints are,} \\ \text{up to a finite set, the union of finitely many maximal arcs.} \end{array} \right. \quad (3.29)$$

In the sequel we identify  $\bar{u}$  (defined up to a null measure set) with a function whose  $i$ th component is constant over each interval of time that is included, up to a zero-measure set, in either  $\check{I}_i$  or  $\hat{I}_i$ . For almost all  $t \in [0, T]$ , the *set of active constraints at time  $t$*  is denoted by  $(\check{B}(t), \hat{B}(t), C(t))$  where

$$\left\{ \begin{array}{l} \check{B}(t) := \{1 \leq i \leq m; \bar{u}_i(t) = \check{u}_i\}, \\ \hat{B}(t) := \{1 \leq i \leq m; \bar{u}_i(t) = \hat{u}_i\}, \\ C(t) := \{1 \leq j \leq q; g_j(\bar{y}(\cdot, t)) = 0\}. \end{array} \right. \quad (3.30)$$

These sets are well-defined over open subsets of  $(0, T)$  where the set of active constraints is constant, and by (3.29), there exist time points called *junction points*

$$0 =: \tau_0 < \dots < \tau_r := T, \quad (3.31)$$

such that the intervals  $(\tau_k, \tau_{k+1})$  are *maximal arcs with constant active constraints*, for  $k = 0, \dots, r-1$ . We may sometimes call them shortly *maximal arcs*.

**Definition 3.6.** For  $k = 0, \dots, r-1$ , let  $\check{B}_k, \hat{B}_k, C_k$  denote the set of indexes of active lower and upper bound constraints, and state constraints, on the maximal arc  $(\tau_k, \tau_{k+1})$ , and set  $B_k := \check{B}_k \cup \hat{B}_k$ .

As a consequence of above definitions and hypothesis (3.26) on the admissible set of controls, we get the following characterization of the first order condition.

**Corollary 3.7.** The first order optimality condition (3.20) is equivalent to

$$\{t \in [0, T]; \Psi_i^p(t) > 0\} \subseteq \check{I}_i, \quad \{t \in [0, T]; \Psi_i^p(t) < 0\} \subseteq \hat{I}_i, \quad (3.32)$$

for every  $(\beta, p, \mu) \in \Lambda$ .

### 3.4 About the jumps of the multiplier at junction points

Given a function  $v : [0, T] \rightarrow X$ , where  $X$  is a Banach space, we denote (if they exist) its left and right limits at  $\tau \in [0, T]$  by  $v(\tau \pm)$ , with the convention  $v(0-) := v(0)$ ,  $v(T+) := v(T)$ ; then the jump of  $v$  at time  $\tau$  is defined as  $[v(\tau)] := v(\tau+) - v(\tau-)$ .

We denote the time derivative of the state constraints by

$$\bar{g}_j^{(1)}[t] := \frac{d}{dt} g_j(\bar{y}(\cdot, t)) = \int_{\Omega} c_j(x) \dot{\bar{y}}(x, t) dx, \quad j = 1, \dots, q. \quad (3.33)$$

Note that  $\bar{g}_j^{(1)}[t]$  is an element of  $L^1(0, T)$ , for each  $j = 1, \dots, q$ .

**Lemma 3.8.** Let  $\bar{u}$  have left and right limits at  $\tau \in (0, T)$ . Then

$$[\Psi_i^p(\tau)][\bar{u}_i(\tau)] = [\bar{g}_j^{(1)}[\tau]][\mu_j(\tau)] = 0, \quad i = 1, \dots, m, \quad j = 1, \dots, q. \quad (3.34)$$

*Proof.* Since  $p = p^1 - \sum_{j=1}^q c_j \mu_j$ ,  $p^1 \in Y \subset C([0, T]; H_0^1(\Omega))$ ,  $\mu \in BV(0, T)_{0,+}^q$ , and any function with bounded variation has left and right limits, we have that  $p(\cdot, \tau)$  has left and right limits in  $H_0^1(\Omega)$  and satisfies

$$[p(\cdot, \tau)] = - \sum_{j=1}^q c_j [\mu_j(\tau)], \quad \text{for all } \tau \in [0, T]. \quad (3.35)$$

Consequently  $\Psi^p$  has left and right limits over  $[0, T]$ , and

$$[\Psi_i^p(\tau)] = - \sum_{j=1}^q [\mu_j(\tau)] \int_{\Omega} b_i(x) c_j(x) \bar{y}(x, \tau) dx, \quad \text{for all } \tau \in [0, T]. \quad (3.36)$$

Next, if  $\bar{u}$  has left and right limits at some  $\tau \in (0, T)$ , then, using the state equation and (3.33), we get

$$[\bar{g}_j^{(1)}[\tau]] = \sum_{i=1}^m [\bar{u}_i(\tau)] \int_{\Omega} b_i(x) c_j(x) \bar{y}(x, \tau) dx. \quad (3.37)$$

Thus, by (3.36) and (3.37), we have

$$\sum_{i=1}^m [\Psi_i^p(\tau)] [\bar{u}_i(\tau)] + \sum_{j=1}^q [\bar{g}_j^{(1)}[\tau]] [\mu_j(\tau)] = 0. \quad (3.38)$$

By the first order conditions (3.32) we have  $[\Psi_i^p(\tau)] [\bar{u}_i(\tau)] \leq 0$ , for  $i = 1$  to  $m$ . Also  $[\mu_j(\tau)] \geq 0$ , and if  $[\mu_j(\tau)] \neq 0$ , the corresponding state constraint has a maximum at time  $\tau$ . Then  $[\bar{g}_j^{(1)}[\tau]] \leq 0$ . So, all terms in the sums in (3.38) are nonpositive and therefore are equal to zero. The conclusion follows.  $\square$

### 3.5 Regularity of the switching function and multiplier over maximal arcs

In the discussion that follows we fix  $k$  in  $\{0, \dots, r-1\}$ , and consider a maximal arc  $(\tau_k, \tau_{k+1})$ , where the junction points are given in (3.31). Recall Definition 3.6 for  $\check{B}_k, \hat{B}_k, B_k \subset \{1, \dots, m\}$  and  $C_k \subset \{1, \dots, q\}$ . Set  $\bar{B}_k := \{1, \dots, m\} \setminus B_k$  and

$$M_{ij}(t) := \int_{\Omega} b_i(x) c_j(x) \bar{y}(x, t) dx, \quad 1 \leq i \leq m, \quad 1 \leq j \leq q. \quad (3.39)$$

Let  $\bar{M}_k(t)$  (of size  $|\bar{B}_k| \times |C_k|$ ) denote the submatrix of  $M(t)$  having rows with index in  $\bar{B}_k$  and columns with index in  $C_k$ . In the sequel we make the following assumption.

**Hypothesis 3.9.** We assume that  $|C_k| \leq |\bar{B}_k|$ , for  $k = 0, \dots, r-1$ , and that the following (uniform) local controllability condition holds:

$$\begin{cases} \text{there exists } \alpha > 0, \text{ such that } |\bar{M}_k(t)\lambda| \geq \alpha|\lambda|, \\ \text{for all } \lambda \in \mathbb{R}^{|C_k|}, \text{ a.e. on } (\tau_k, \tau_{k+1}), \text{ for } k = 0, \dots, r-1. \end{cases} \quad (3.40)$$

**Remark 3.10.** This hypothesis was already used in a different setting (i.e. higher-order state constraints in the finite dimensional case) in e.g. [7, 26]. Note that condition (3.40) implies, in particular, that the matrix  $\bar{M}_k(t)$  has rank  $|C_k|$  over  $(\tau_k, \tau_{k+1})$ .

The expression of the derivative of the  $j$ th state constraint, for  $1 \leq j \leq q$ , is

$$\bar{g}_j^{(1)}[t] = \int_{\Omega} c_j(x) (f(x, t) + \Delta \bar{y}(x, t) - \gamma \bar{y}(x, t)^3) dx + \sum_{i=1}^m M_{ij}(t) \bar{u}_i(t), \quad (3.41)$$

or, in vector form, for the active state constraints (denoting by  $\bar{g}_{C_k}^{(1)}[t]$  the vector of components  $\bar{g}_j^{(1)}[t]$  for  $j \in C_k$ ), we get

$$\bar{g}_{C_k}^{(1)}[t] = G_k(t) + \bar{M}_k(t)^\top \bar{u}_{\bar{B}_k}(t) = 0, \quad (3.42)$$

where  $\bar{u}_{\bar{B}_k}$  is the restriction of  $\bar{u}$  to the components in  $\bar{B}_k$ , and  $G_k(t)$  takes into account the contributions of the integral in (3.41) and of the components of  $\bar{u}$  in  $B_k$ , that is, for  $j \in C_k$ :

$$G_{k,j}(t) := \int_{\Omega} c_j (f(x, t) + \Delta \bar{y}(x, t) - \gamma \bar{y}(x, t)^3) dx + \sum_{i \in B_k} M_{ij}(t) \bar{u}_i(t). \quad (3.43)$$

By the controllability condition (3.40),  $\bar{M}_k(t)^\top$  is onto from  $\mathbb{R}^{|\bar{B}_k|}$  to  $\mathbb{R}^{|C_k|}$ . In view of the state equation, by an integration by parts argument,  $M(t)$  has a bounded derivative and is therefore Lipschitz continuous. So there exists a linear change of control variables of the form  $u(t) = N_k(t) \hat{u}(t)$ , for some invertible Lipschitz continuous matrix  $N_k(t)$  of size  $m \times m$ , such that, calling  $\bar{N}_k(t)$  the upper  $|\bar{B}_k| \times |\bar{B}_k|$ -diagonal block of  $N_k(t)$ , it holds that  $\bar{M}_k(t)^\top \bar{N}_k(t)$  has its first  $|C_k|$  columns being equal to the identity matrix, the other columns having null components. That is, for all  $\hat{u} \in \mathbb{R}^{|\bar{B}_k|}$ :

$$(\bar{M}_k(t)^\top \bar{N}_k(t) \hat{u})_j = \hat{u}_j, \quad \text{for } j = 1, \dots, |C_k|. \quad (3.44)$$

Over a maximal arc  $(\tau_k, \tau_{k+1})$ , we have that  $\bar{g}_j^{(1)}[t] = 0$  for  $j \in C_k$  is equivalent to

$$\hat{u}_j = -G_{k,j}(t), \quad \text{for } j = 1, \dots, |C_k|. \quad (3.45)$$

The following result on the regularity of the state constraint multiplier holds. Recall the definition of the switching function  $\Psi^p$  given in (3.19).

**Proposition 3.11.** *There exists  $a \in L^1(0, T)^m$  such that*

$$(i) \quad d\Psi^p(t) = a(t)dt - M(t)d\mu(t), \quad \text{on } [0, T]. \quad (3.46)$$

(ii) *We have that  $\dot{\mu}_{C_k}$  is locally integrable over  $(\tau_k, \tau_{k+1})$ , hence  $\mu_{C_k}$  is locally absolutely continuous, and the following expression holds*

$$0 = \dot{\Psi}_{\bar{B}_k}^p(t) = a_{\bar{B}_k}(t)dt - \bar{M}_k(t) \dot{\mu}_{C_k}(t), \quad \text{on } (\tau_k, \tau_{k+1}). \quad (3.47)$$

*Proof.* By (3.8) and (3.19), one has, for  $i \in \{1, \dots, m\}$ :

$$\Psi_i^p(t) = \alpha_i + \int_{\Omega} b_i(x) \bar{y}(x, t) p^1(x, t) dx - \sum_{j=1}^q M_{ij}(t) \mu_j(t), \quad i = 1, \dots, m. \quad (3.48)$$

Let  $a: (0, T) \rightarrow \mathbb{R}^m$  be given by

$$a_i(t) := \frac{d}{dt} \int_{\Omega} b_i(x) \bar{y}(x, t) p^1(x, t) dx - \sum_{j=1}^q \dot{M}_{ij}(t) \mu_j(t), \quad \text{for } i = 1, \dots, m. \quad (3.49)$$



Note that  $\dot{M}_{ij}(t) = \int_{\Omega} b_i(x)c_j(x)\dot{\bar{y}}(x, t)dx$  is integrable (this follows integrating by parts the contribution of  $\Delta\bar{y}$  and since  $Y \subset C([0, T]; H_0^1(\Omega))$ ), and that

$$\frac{d}{dt}(\bar{y}p^1) = p^1 \Delta\bar{y} - \bar{y} \Delta p^1 + fp^1 + 2\gamma\bar{y}^3 p^1 - \beta\bar{y}(\bar{y} - y_d) - \sum_{j=1}^q \mu_j \bar{y} A c_j. \quad (3.50)$$

Integrating by parts the terms in (3.50) containing Laplacians, we get, for the integral term in (3.49),

$$\begin{aligned} \int_{\Omega} b_i(x) \frac{d}{dt}(\bar{y}p^1) dx &= \int_{\Omega} b_i \left( fp^1 + 2\gamma\bar{y}^3 p^1 - \beta\bar{y}(\bar{y} - y_d) - \sum_{j=1}^q \mu_j \bar{y} A c_j \right) dx \\ &\quad - \int_{\Omega} \nabla b_i (p^1 \nabla \bar{y} - \bar{y} \nabla p^1) dx. \end{aligned} \quad (3.51)$$

It follows that  $a \in L^1(0, T)^m$  and (3.46) holds. Consequently  $\Psi^p$  has bounded variation.

Over  $(\tau_k, \tau_{k+1})$ , we have  $d\mu_j(t) = 0$  whenever  $j \notin C_k$ , and so

$$0 = d\Psi_{\bar{B}_k}^p(t) = a_{\bar{B}_k}(t)dt - \bar{M}_k(t)d\mu_{C_k}(t). \quad (3.52)$$

Since  $\bar{M}_k(t)$  is continuous and injective, and  $a$  is integrable, this implies the existence of  $\dot{\mu}_j(t) \in L^1(0, T)$ , for  $j \in C_k$ . This yields (3.47).

And so,  $\mu_{C_k}(t)$  is locally absolutely continuous.  $\square$

**Corollary 3.12.** *Let the finite maximal arc property (3.29) and the uniform controllability condition (3.40) hold.*

- (i) *If  $f, y_d \in L^\infty(0, T; L^2(\Omega))$ , then  $a \in L^\infty(0, T)^m$ .*
- (ii) *If additionally  $f, y_d \in C([0, T]; L^2(\Omega))$ , then  $\mu$  is  $C^1$  over each maximal arc  $(\tau_k, \tau_{k+1})$ .*

*Proof.* Indeed, a careful inspection of the previous proof shows that  $a$  is a sum of essentially bounded terms, so (i) follows. If the additional regularity hypotheses of item (ii) hold, then  $a$  is continuous. The regularity of  $\mu$  follows from (3.52) and the local controllability assumption (3.40). This concludes the proof.  $\square$

## 4 Second order necessary conditions

In this section we derive second order necessary optimality conditions, based on the concept of *radiality* of critical directions.

Let us consider an admissible trajectory  $(\bar{u}, \bar{y})$ .

### 4.1 Assumptions and additional regularity

For the remainder of the article we make the following set of assumptions.

**Hypothesis 4.1.** *The following conditions hold:*

1. the control set has the form (3.26),
2. the finite maximal arc property (3.29),
3. the qualification hypothesis (3.25),
4. the local (uniform) controllability condition (3.40) over each maximal arc  $(\tau_k, \tau_{k+1})$ ,
5. the discontinuity of the derivative of the state constraints at corresponding junction points, i.e.,

$$\text{for some } c > 0: g_j(\bar{y}(\cdot, t)) \leq -c \operatorname{dist}(t, I_j^C), \text{ for all } t \in [0, T], j = 1, \dots, q, \quad (4.1)$$

6. the uniform distance to control bounds whenever they are not active, i.e. there exists  $\delta > 0$  such that,

$$\operatorname{dist}(\bar{u}_i(t), \{\tilde{u}_i, \hat{u}_i\}) \geq \delta, \quad \text{for a.a. } t \notin I_i, \text{ for all } i = 1, \dots, m, \quad (4.2)$$

7. the following regularity for the data (we do not try to take the weakest hypotheses) for some  $r > n + 1$ :

$$y_0, y_{dT} \in W_0^{1,r}(\Omega) \cap W^{2,r}(\Omega), \quad y_d, f \in L^\infty(Q), \quad b \in L^\infty(\Omega)^{m+1}, \quad (4.3)$$

8. the control  $\bar{u}$  has left and right limits at the junction points  $\tau_k \in (0, T)$ , (this will allow to apply Lemma 3.8).

In view of point 3 above, we consider from now on  $\beta = 1$  and thus we omit the component  $\beta$  of the multipliers.

**Theorem 4.2.** *The following assertions hold.*

(i) *For any  $u \in L^\infty(0, T)^m$ , the associated state  $y[u]$  belongs to  $C(\bar{Q})$ . If  $u$  remains in a bounded subset of  $L^\infty(0, T)^m$  then the corresponding states form a bounded set in  $C(\bar{Q})$ . In addition, if the sequence  $(u_\ell)$  of admissible controls converges to  $\bar{u}$  a.e. on  $(0, T)$ , then the associated sequence of states  $(y_\ell := y[u_\ell])$  converges uniformly to  $\bar{y}$  in  $\bar{Q}$ .*

(ii) *For every  $(p, \mu) \in \Lambda_1$ , one has that  $\mu \in W^{1,\infty}(0, T)^q$  and  $p$  is essentially bounded in  $Q$ .*

*Proof.* (i) Let  $r \in [2, \infty)$ . That  $y \in W^{2,1,r}(Q)$  follows from Theorem A.3 in the Appendix. Taking  $r > n + 1$ , it follows from the Sobolev Embedding Theorem (see e.g. [17, Theorem 5, p. 269]) that  $y$  is continuous (and even Hölder-continuous) on the closure of  $Q$ , with uniform bound over the set of admissible controls. If the sequence  $(u_\ell)$  of admissible controls converges a.e. to  $\bar{u}$ , by the Dominated Convergence Theorem,  $u_\ell \rightarrow \bar{u}$  in  $L^q(0, T)$  for all  $q \in [1, \infty)$ . So, by similar arguments it can be proved that the associated sequence of states converges uniformly to  $\bar{y}$ .

(ii) By Hypothesis 4.1,  $y_{dT}$  is the trace at time  $T$  of an element of  $W^{2,1,r}(Q)$  vanishing on  $\Sigma$  and this obviously holds also for  $y(T)$  in view of Theorem A.3 in the Appendix. It follows then from corollary A.2 that  $p^1 \in W^{2,1,r}(Q)$ . The continuity of  $\mu$  at junction points follows from (4.1) in Hypothesis 4.1 and Lemma 3.8. The boundedness on each arc of the derivative of  $\mu$  follows from (3.47) for  $\dot{\mu}$ , since by Corollary 3.12,  $a \in L^\infty(0, T)^m$  and by (3.40),  $\bar{M}(t)$  is ‘uniformly injective’ over each arc. The conclusion follows.  $\square$

## 4.2 Second variation

For  $(p, \mu) \in \Lambda_1$ , set

$$\kappa(x, t) := 1 - 6\gamma\bar{y}(x, t)p(x, t), \quad (4.4)$$

and consider the quadratic form

$$\mathcal{Q}[p](z, v) := \int_Q \left( \kappa z^2 + 2p \sum_{i=1}^m v_i b_i z \right) dxdt + \int_{\Omega} z(x, T)^2 dx. \quad (4.5)$$

Let  $(u, y)$  be a trajectory, and set

$$(\delta y, v) := (y - \bar{y}, u - \bar{u}). \quad (4.6)$$

Recall the definition of the operator  $A$  given in (2.20). Subtracting the state equation at  $(\bar{u}, \bar{y})$  from the one at  $(u, y)$ , we get that

$$\begin{cases} \frac{d}{dt}\delta y + A\delta y = \sum_{i=1}^m v_i b_i y - 3\gamma\bar{y}(\delta y)^2 - \gamma(\delta y)^3 & \text{in } Q, \\ \delta y = 0 & \text{on } \Sigma, \quad \delta y(\cdot, 0) = 0 & \text{in } \Omega. \end{cases} \quad (4.7)$$

Combining with the linearized state equation (3.1), we deduce that  $\eta$  given by

$$\eta := \delta y - z, \quad (4.8)$$

satisfies the equation

$$\begin{cases} \dot{\eta} - \Delta\eta = r\eta + \tilde{r} & \text{in } Q, \\ \eta = 0 & \text{on } \Sigma, \quad \eta(\cdot, 0) = 0 & \text{in } \Omega \end{cases} \quad (4.9)$$

where  $r$  and  $\tilde{r}$  are defined as

$$r := -3\gamma\bar{y}^2 + \sum_{i=0}^m \bar{u}_i b_i, \quad \tilde{r} := \sum_{i=1}^m v_i b_i \delta y - 3\gamma\bar{y}(\delta y)^2 - \gamma(\delta y)^3. \quad (4.10)$$

**Proposition 4.3.** *Let  $(p, \mu) \in \Lambda_1$ , and let  $(u, y)$  be a trajectory. Then*

$$\begin{aligned} & \mathcal{L}[p, \mu](u, y, p) - \mathcal{L}[p, \mu](\bar{u}, \bar{y}, p) \\ &= \int_0^T \Psi^p(t) \cdot v(t) dt + \frac{1}{2} \mathcal{Q}[p](\delta y, v) - \gamma \int_Q p(\delta y)^3 dxdt. \end{aligned} \quad (4.11)$$

Here, we omit the dependence of the Lagrangian on  $(\beta, p_0)$  being equal to  $(1, p(\cdot, 0))$ .

*Proof.* Use  $\Delta\mathcal{L}$  to denote the l.h.s. of (4.11). We have

$$\begin{aligned} \Delta\mathcal{L} &= J(u, y) - J(\bar{u}, \bar{y}) + \int_Q p \left( -\frac{d}{dt}\delta y + \Delta\delta y - \gamma(y^3 - \bar{y}^3) \right) dxdt \\ &+ \int_Q p \left( \sum_{i=1}^m v_i b_i y + \sum_{i=0}^m \bar{u}_i b_i \delta y \right) dxdt + \sum_{j=1}^q \int_0^T \int_{\Omega} c_j \delta y dx d\mu_j(t) \\ &= \int_Q \delta y \left( \frac{1}{2}\delta y + \bar{y} - y_d \right) dxdt + \int_{\Omega} \delta y(x, T) \left( \frac{1}{2}\delta y(x, T) + \bar{y}(x, T) - y_{dT}(x) \right) dx \\ &+ \sum_{i=1}^m \alpha_i \int_0^T v_i dt + \int_Q p \left( -\frac{d}{dt}\delta y + \Delta\delta y - \gamma(\delta y^3 + 3\bar{y}\delta y^2 + 3\bar{y}^2\delta y) \right) dxdt \\ &+ \int_Q p \left( \sum_{i=1}^m v_i b_i y + \sum_{i=0}^m \bar{u}_i b_i \delta y \right) + \sum_{j=1}^q \int_0^T \int_{\Omega} c_j \delta y dx d\mu_j(t). \end{aligned} \quad (4.12)$$

By (3.5) we obtain

$$\begin{aligned} \int_Q p \frac{d}{dt} \delta y \, dx dt &= - \int_Q p A \delta y \, dx dt + \sum_{j=1}^q \int_0^T \int_{\Omega} c_j \delta y \, dx d\mu_j(t) \\ &+ \int_Q \delta y (\bar{y} - y_d) \, dx dt + \int_{\Omega} \delta y(x, T) (\bar{y}(x, T) - y_{dT}(x)) \, dx. \end{aligned} \quad (4.13)$$

Thus, from (4.12) and (4.13) we get

$$\begin{aligned} \Delta \mathcal{L} &= \frac{1}{2} \int_Q \delta y^2 \, dx dt + \frac{1}{2} \int_{\Omega} \delta y(\cdot, T)^2 \, dx + \sum_{i=1}^m \alpha_i \int_0^T v_i \, dt \\ &+ \int_Q p \left( -\gamma [\delta y^3 + 3\bar{y} \delta y^2] + \sum_{i=1}^m v_i b_i y \right) \, dx dt, \end{aligned} \quad (4.14)$$

which leads to (4.11) in view of the definition of  $\Psi_i^p$  given in (3.19). This concludes the proof.  $\square$

### 4.3 Critical directions

Recall the definitions of  $\check{I}_i, \hat{I}_i$  and  $I_j^C$  given in (3.27) and (3.28), and remember that we use  $z[v]$  to denote the solution of the linearized state equation (3.1) associated to  $v$ .

Let us define the *cone of critical directions* at  $\bar{u}$  in  $L^2$ , or in short *critical cone*, by

$$C := \left\{ \begin{array}{l} (z[v], v) \in Y \times L^2(0, T)^m; \\ v_i(t) \Psi_i^p(t) = 0 \text{ a.e. on } [0, T], \text{ for all } (p, \mu) \in \Lambda_1 \\ v_i(t) \geq 0 \text{ a.e. on } \check{I}_i, v_i(t) \leq 0 \text{ a.e. on } \hat{I}_i, \text{ for } i = 1, \dots, m, \\ \int_{\Omega} c_j(x) z[v](x, t) \, dx \leq 0 \text{ on } I_j^C, \text{ for } j = 1, \dots, q \end{array} \right\}. \quad (4.15)$$

The *strict critical cone* is defined below, and it is obtained by imposing that the linearization of active constraints is zero,

$$C_s := \left\{ \begin{array}{l} (z[v], v) \in Y \times L^2(0, T)^m; v_i(t) = 0 \text{ a.e. on } I_i, \text{ for } i = 1, \dots, m, \\ \int_{\Omega} c_j(x) z[v](x, t) \, dx = 0 \text{ on } I_j^C, \text{ for } j = 1, \dots, q \end{array} \right\}. \quad (4.16)$$

Hence, clearly  $C_s \subseteq C$ , and  $C_s$  is a closed subspace of  $Y \times L^2(0, T)^m$ . Now, note that in the interior of each  $I_j^C$  one has, for every  $(z[v], v) \in C_s$ ,

$$\begin{aligned} 0 &= \frac{d}{dt} (g'_j(\bar{y}(\cdot, t)) z[v](\cdot, t)) = \frac{d}{dt} \int_{\Omega} c_j(x) z[v](x, t) \, dx \\ &= \int_{\Omega} c_j(x) \dot{z}[v](x, t) \, dx = \int_{\Omega} c_j(x) \left( -(Az[v])(x, t) + (v(t) \cdot b(x)) \bar{y}(x, t) \right) \, dx, \end{aligned} \quad (4.17)$$

which can be rewritten as

$$\sum_{i=1}^m v_i(t) M_{ij}(t) = \int_{\Omega} c_j(x) (Az[v])(x, t) \, dx, \quad (4.18)$$

in view of the definition of  $M_{ij}$  given in (3.39). Therefore, over any arc  $(a, b)$  we have  $g'_j(\bar{y}(\cdot, t))z[v](\cdot, t) = 0$  for  $t \in (a, b)$  if and only if  $g'_j(\bar{y}(\cdot, a))z[v](\cdot, a) = 0$  and (4.18) holds over  $(a, b)$ . We define the *entry (resp. exit) point* of a time interval  $(t', t'')$  as  $t'$  (resp.  $t''$ ). This induces the consideration of the following sets

$$C_e := \left\{ (z[v], v) \in Y \times L^2(0, T)^m; \right. \\ \left. g'_j(\bar{y}(\cdot, \tau_k))z[v](\tau_k) = 0, \text{ if } j \in C_k, \text{ for } k = 0, \dots, r-1 \right\}, \quad (4.19)$$

$$C_n := \left\{ (z[v], v) \in Y \times L^2(0, T)^m; v_i(t) = 0 \text{ a.e. on } I_i, \text{ for } i = 1, \dots, m, \right. \\ \left. \sum_{i=1}^m v_i(t)M_{ij}(t) = \int_{\Omega} c_j(x)(Az[v])(x, t)dx \text{ a.e. on } I_j^C, \text{ for } j = 1, \dots, q \right\}. \quad (4.20)$$

With these definitions, we can write the strict critical cone as

$$C_s = C_e \cap C_n, \quad (4.21)$$

and prove the following result.

**Lemma 4.4.**  $C_s \cap (Y \times L^\infty(0, T)^m)$  is dense in  $C_s$ , with respect to the  $Y \times L^2(0, T)^m$ -topology.

*Proof.* In view of Dmitruk's density lemma (see [16, Lemma 1]), it is enough to prove that  $C_n \cap (Y \times L^\infty(0, T)^m)$  is a dense subset of  $C_n$ .

Let us then take  $(z, v) \in C_n$ . Recall the definition of the junction times  $\tau_k$  given after equation (3.39). Fix  $k \in \{0, \dots, r-1\}$ . Note that we can take a partition of  $[0, T]$ , say  $0 = t_0 \leq \dots \leq t_\ell \leq \dots \leq t_N = T$ , such that  $(t_\ell, t_{\ell+1})$  is contained in some  $(\tau_k, \tau_{k+1})$ , and on  $(t_\ell, t_{\ell+1})$  a fixed set of the rows of  $M(t)$  is linearly independent with rank equal to the one of  $M(t)$ . Now consider the matrix  $\bar{M}_k$  given after (3.39). Using the same notation as in (3.42), let us write  $v_{\bar{B}_k}$  to refer to the restriction of  $v$  to the components in  $\bar{B}_k$ . For each  $t \in (t_\ell, t_{\ell+1})$ , we can write

$$v_{\bar{B}_k}(t) = v_{\bar{B}_k,0}(t) + v_{\bar{B}_k,1}(t), \quad (4.22)$$

where  $v_{\bar{B}_k,0}(t) \in \text{Ker } \bar{M}_k(t)^\top$  and  $v_{\bar{B}_k,1}(t) \in \text{Im } \bar{M}_k(t)$  for almost all  $t$ , hence  $v_{\bar{B}_k,1}(t) = \bar{M}_k(t)\lambda_k(t)$  for some  $\lambda_k(t) \in \mathbb{R}^{|C_k|}$ . Let  $E_{C_k}(t)$  be the  $|C_k|$ -dimensional vector with components

$$E_{C_k,j}(t) := \int_{\Omega} c_j(x)(Az)(x, t)dx, \quad j \in C_k. \quad (4.23)$$

Then (4.18) can be rewritten as

$$E_{C_k}(t) = \bar{M}_k(t)^\top v_{\bar{B}_k}(t) = \bar{M}_k(t)^\top v_{\bar{B}_k,1}(t) = \bar{M}_k(t)^\top \bar{M}_k(t)\lambda_k(t), \quad (4.24)$$

and, therefore,  $\lambda_k(t) = (\bar{M}_k(t)^\top \bar{M}_k(t))^{-1} E_{C_k}(t)$ , so that

$$v_{\bar{B}_k,1}(t) = \bar{M}_k(t)\lambda_k(t) = \bar{M}_k(t)(\bar{M}_k(t)^\top \bar{M}_k(t))^{-1} E_{C_k}(t). \quad (4.25)$$

By an integration by parts (in space) argument, it follows that  $E_{C_k}(t)$  is a continuous function, and so is  $\bar{M}_k(t)$ . Therefore,  $v_{\bar{B}_k,1}$  is continuous on each maximal arc. We may also view the application  $z \mapsto v_{\bar{B}_k,1}$  as a linear and continuous mapping say

$$L_1 : Y \rightarrow \prod_{k=0}^{r-1} \text{Lip}(\tau_k, \tau_{k+1})^{|C_k|} \quad (4.26)$$

where  $C_k$  is the set of active state constraints on  $(\tau_k, \tau_{k+1})$  and, for  $t' < t''$ ,  $\text{Lip}(t', t'')$  is the Banach space of continuous real functions with domain  $(t', t'')$ , endowed with the norm

$$\|f\|_{\text{Lip}(t', t'')} := \sup_{t \in (t', t'')} |f(t)| + \sup_{t, \tau \in (t', t'')} \frac{|f(t) - f(\tau)|}{|t - \tau|}, \quad (4.27)$$

with the convention “ $0/0 = 0$ ”.

For any  $\varepsilon > 0$ , there exists  $v_{\bar{B}_k, 0}^\varepsilon$  in  $L^\infty(0, T)^{|B_k|}$  such that  $\|v_{\bar{B}_k, 0}^\varepsilon - v_{\bar{B}_k, 0}\|_2 < \varepsilon$ , it has zero components for indexes corresponding to active control bound constraints, and  $v_{\bar{B}_k, 0}^\varepsilon(t) \in \text{Ker } \bar{M}_k(t)^\top$  for a.a.  $t$ . In fact, to construct this  $v_{\bar{B}_k, 0}^\varepsilon$  it suffices to project an approximation of  $v_{\bar{B}_k, 0}$  obtained by a truncation argument on the kernel  $\text{Ker } \bar{M}_k(t)^\top$ . In what follows we shall abuse notation and use the same symbol to denote a vector and its canonical immersion in  $\mathbb{R}^m$ . Let  $z_\varepsilon$  be the unique solution in  $Y$  of the linearized equation

$$\dot{z}_\varepsilon + Az_\varepsilon = \sum_{i=1}^m (L_1(z_\varepsilon) + v_{\bar{B}, 0}^\varepsilon + v_B)_i b_i \bar{y}, \quad (4.28)$$

with the usual initial and boundary conditions, and where  $v_B$  is the restriction of  $v$  to the set  $B$ . Set  $v_{\bar{B}, 1}^\varepsilon := L_1(z_\varepsilon)$ ,  $v_{\bar{B}_k}^\varepsilon := v_{\bar{B}_k, 1}^\varepsilon + v_{\bar{B}_k, 0}^\varepsilon$ , and define  $v_\varepsilon$  to have the restriction to  $\bar{B}_k$  equal to  $v_{\bar{B}_k}^\varepsilon$  and the restriction to  $B_k$  equal to  $v$ . Then  $v_\varepsilon$  is in  $C_n \cap (Y \times L^\infty(0, T)^m)$  and  $\|v_\varepsilon - v\|_2 = O(\varepsilon)$ . Hence,  $C_n \cap (Y \times L^\infty(0, T)^m)$  is a dense subset of  $C_n$ . The conclusion follows.  $\square$

### 4.3.1 Radiality of critical directions

According to Aronna *et al.* [3, Definition 6], a critical direction  $(z, v)$  is *quasi radial* if there exists  $\tau_0 > 0$  such that, for  $\tau \in [0, \tau_0]$ , the following conditions are satisfied:

$$\max_{t \in [0, T]} \{g_j(\bar{y}(\cdot, t)) + \tau g'_j(\bar{y}(\cdot, t))z(t)\} = o(\tau^2), \quad \text{for } j = 1, \dots, q, \quad (4.29)$$

$$\check{u}_i \leq \bar{u}_i(t) + \tau v_i(t) \leq \hat{u}_i, \quad \text{a.e. on } [0, T], \quad \text{for } i = 1, \dots, m. \quad (4.30)$$

**Lemma 4.5.** *Every direction in  $C_s \cap (Y \times L^\infty(0, T)^m)$  is quasi radial.*

*Proof.* Let  $(z, v) \in C_s \cap (Y \times L^\infty(0, T)^m)$ . Then (4.30) follows from (4.2). Let us next prove (4.29). The function  $h(t) := g'_j(\bar{y}(t))z(t)$  has the derivative  $\dot{h}(t) = \int_\Omega c_j(x) \dot{z}(x, t) dx$ , so that  $|\dot{h}(t)| \leq \|c_j\|_{L^2(\Omega)} \|\dot{z}(\cdot, t)\|_{L^2(\Omega)}$  and hence,  $\dot{h} \in L^2(0, T)$ . Let  $0 \leq t' < t'' \leq T$ . By the Cauchy-Schwarz inequality, for any  $\varepsilon > 0$ :

$$|h(t'') - h(t')| \leq \int_{t'}^{t''} |\dot{h}(t)| dt \leq \sqrt{t'' - t'} \|\dot{h}\|_{L^2(t', t'')}. \quad (4.31)$$

Let  $(a, b)$  be a maximal constrained arc with say  $a > 0$ . Take  $t' < a$ , and  $t'' = a$ . When  $t' \uparrow a$ , by the Dominated Convergence Theorem,  $\|\dot{h}\|_{L^2(t', t'')} \rightarrow 0$ . Given  $\varepsilon > 0$ , we deduce with (4.1) that for  $\tau > 0$  and  $t' < a$  close enough to  $a$ :

$$g_j(\bar{y}(\cdot, t)) + \tau g'_j(\bar{y}(\cdot, t))z(t) \leq -c(a - t) + \tau \varepsilon \sqrt{a - t}, \quad \text{for all } t \in (t', a). \quad (4.32)$$

The maximum of the r.h.s. of (4.32) over  $t \in [a - \varepsilon, a]$  is attained when

$$c\sqrt{a-t} = \frac{1}{2}\tau\varepsilon, \quad a-t = \frac{\tau^2\varepsilon^2}{4c^2}. \quad (4.33)$$

So the r.h.s. of (4.32) is less or equal than  $\tau^2\varepsilon^2/(4c)$ . Since we can take  $\varepsilon$  arbitrarily small, it is of order  $o(\tau^2)$ . For  $t > b$  close to  $b$ , we have a similar result. For  $t$  far from the boundary, (4.29) is a consequence of hypothesis (4.1). The conclusion follows.  $\square$

Combining the previous result with Lemma 4.4, we deduce that:

**Corollary 4.6.** *The set of quasi radial critical directions of  $C_s$  is dense in  $C_s$ .*

#### 4.4 Second order necessary condition

We obtain the following result applying Corollary 4.6 above and the second order condition in an abstract setting proved in [3, Theorem 8].

**Theorem 4.7** (Second order necessary condition). *Let the admissible trajectory  $(\bar{u}, \bar{y})$  be an  $L^\infty$ -local solution of  $(P)$ . Then*

$$\max_{(p,\mu) \in \Lambda_1} \mathcal{Q}[p](z, v) \geq 0, \quad \text{for all } (z, v) \in C_s. \quad (4.34)$$

*Proof.* Let  $(z, v) \in C_s$ . By Corollary 4.6, there exists a sequence  $(z^\ell, v^\ell)$  of quasi radial directions converging to  $(z, v)$  in  $Y \times L^2(0, T)^m$ . Doing as in [3, Theorem 8], we get the existence of a multiplier  $(p^\ell, \mu^\ell) \in \Lambda_1$  (with  $\Lambda_1$  defined in Section 3.2.1), such that

$$\mathcal{Q}[p^\ell](z^\ell, v^\ell) \geq 0. \quad (4.35)$$

By Lemma 3.5,  $\Lambda_1$  is bounded so that  $d\mu^\ell$  is also bounded. Extracting if necessary a subsequence, we may assume that  $d\mu^\ell$  weakly-\* converges to some  $d\mu$  with  $\mu \in BV(0, T)_{0,+}^q$ , and since  $L^\infty(0, T, H_0^1(\Omega))$  is included in  $L^2(Q)$ ,  $p^\ell$  weakly converges in  $L^2(Q)$  to some  $p \in L^2(Q)$ , such that  $(p, \mu) \in \Lambda_1$ . Since  $(z^\ell, v^\ell) \rightarrow (z[v], v)$  in  $Y \times L^2(0, T)^m$ , by lemma 2.2,  $\sum_i v_i^\ell b_i z^\ell$  strongly converges to  $\sum_i v_i b_i z$ , and so we easily deduce that  $\mathcal{Q}[p^\ell](z^\ell, v^\ell) \rightarrow \mathcal{Q}[p](z[v], v)$ . The conclusion follows.  $\square$

## A Strong solutions of the heat equation

We consider the heat equation with Dirichlet boundary condition:

$$\dot{y} - \Delta y = f \text{ in } Q, \quad y(x, 0) = y_0(x); \quad y = h \text{ on } \Sigma. \quad (A.1)$$

We have the following result, see Lieberman [21, Thm 7.32, p. 182]:

**Theorem A.1.** *Let  $r \geq 2$ ,  $w \in W^{2,1,r}(Q)$  and  $f \in L^r(Q)$ . Setting  $y_0 := w(\cdot, 0)$  and  $h := \tau_\Sigma w$  (trace of  $w$  over  $\Sigma$ ), equation (A.1) has a unique solution  $y \in W^{2,1,r}(Q)$ . In addition there exists  $C > 0$  such that*

$$\|y\|_{W^{2,1,r}(Q)} \leq C (\|f\|_{L^r(Q)} + \|w\|_{W^{2,1,r}(Q)}). \quad (A.2)$$

**Corollary A.2.** *Given  $r \geq 2$ ,  $y_0 \in W_0^{1,r}(\Omega) \cap W^{2,r}(\Omega)$  and  $f \in L^r(Q)$ , equation (A.1) has, for  $h = 0$ , a unique solution  $y \in W^{2,1,r}(Q)$  that satisfies*

$$\|y\|_{W^{2,1,r}(Q)} \leq C (\|f\|_{L^r(Q)} + \|y_0\|_{W^{2,r}(\Omega)}). \quad (\text{A.3})$$

*Proof.* Apply Theorem A.1 with  $w(x, t) := y_0(x)$ . It is clear that  $w \in W^{2,1,r}(Q)$  and that  $w$  has trace  $y_0$  at time 0 and zero trace over  $\Sigma$ . The conclusion follows.  $\square$

By the standard Sobolev embeddings, we have the continuous inclusion

$$W^{2,1,r}(Q) \subset W^{1,r}(Q) \subset L^\infty(Q), \quad \text{if } r > n + 1. \quad (\text{A.4})$$

This allows to prove the following.

**Theorem A.3.** *Assume that  $u \in L^\infty(0, T)$ ,  $y_0 \in W_0^{1,r}(\Omega) \cap W^{2,r}(\Omega)$  and  $f \in L^r(Q)$ , with  $r > n + 1$ . Then the state equation (2.1) has a unique solution  $y[u, y_0, f]$  in  $W^{2,1,r}(Q)$ , and the mapping  $y[u, y_0, f]$  is of class  $C^\infty$  from  $L^\infty(0, T) \times W_0^{1,r}(\Omega) \cap W^{2,r}(\Omega) \times L^r(Q)$  into  $W^{2,1,r}(Q)$ .*

*Proof.* We have that  $g := -\Delta y_0$  belongs to  $L^r(\Omega)$ . Let  $y_0^\pm$  be the unique solution of  $-\Delta y_0^\pm = g^\pm$  in  $\Omega$ , where  $g^+ := \max(g, 0)$  and  $g^- := -\min(g, 0)$ , with homogeneous Dirichlet condition on the boundary. Set  $f^+ := \max(f, 0)$  and  $f^- := -\min(f, 0)$ . Denote by  $y^+$  (resp.,  $y^-$ ) the solution of the state equation (2.1) when  $(y_0, f)$  is  $(y_0^+, f^+)$  (resp.  $(y_0^-, f^-)$ ). By the monotonicity results in Lemma 2.3, we have that  $-y^- \leq y \leq y^+$ . Now let  $y^{++}$ ,  $y^{--}$  denote the solutions of the state equation (2.1) when  $(y_0, f)$  is  $(y_0^+, f^+)$ ,  $(y_0^-, f^-)$ , respectively and, in addition,  $\gamma = 0$ . We claim that  $-y^{--} \leq -y^- \leq y \leq y^+ \leq y^{++}$ . Indeed, for  $z \in Y$ , set  $H_u z := \dot{z} - \Delta z - z \sum_i u_i b_i$ . Then

$$H_u y^+ = f^+ - \gamma(y^+)^3 \leq f^+ = H_u y^{++}. \quad (\text{A.5})$$

Since  $y^+$  and  $y^{++}$  have the same initial conditions, it follows that  $y^+ \leq y^{++}$ . In an analogous way, it can be proved that  $-y^{--} \leq -y^-$ .

Since  $y_0^\pm \in W_0^{1,r}(\Omega) \cap W^{2,r}(\Omega)$  and  $f^\pm \in L^r(Q)$ , by Corollary A.2,  $y^{++}$  and  $y^{--}$  belong to  $W^{2,1,r}(Q)$  and, therefore, since  $r > n + 1$ , they are also elements of  $L^\infty(Q)$ . So,  $y \in L^\infty(Q)$ . Consequently,  $H_u y = f - \gamma y^3 \in L^r(Q)$  and, by Theorem A.1 again,  $y \in W^{2,1,r}(Q)$ .

We recall that, for  $r > n + 1$ ,  $Y_r$  denotes the set of elements of  $W^{2,1,r}(Q)$  with zero trace on  $\Sigma$ , and  $Y_r^0$  denotes the trace of  $Y_r$  at time zero. Endowed with the "trace norm",  $Y_r^0$  is a Banach space that contains  $W_0^{1,r}(\Omega) \cap W^{2,r}(\Omega)$  in view of the proof of the above Corollary A.2 (by Lions [23, p. 20],  $Y_r^0$  is a subset of  $W^{2-2/r,r}(\Omega)$ ). That  $(u, y_0, f) \mapsto y[u, y_0, f]$  is of class  $C^\infty$  is a consequence of the Implicit Function Theorem applied to the mapping  $F$  from  $Y_r \times L^\infty(0, T) \times Y_r^0 \times L^r(Q)$  into  $L^r(Q) \times Y_r^0$ , defined by

$$F(y, u, y_0, f) := (H_u y + \gamma y^3, y(0) - y_0). \quad (\text{A.6})$$

The key step is to prove that the partial derivative  $D_y F$  is bijective; this can be done easily, taking advantage of the fact that  $W^{2,1,r}(Q) \subset L^\infty(Q)$  when  $r > n + 1$ .  $\square$

## B An example

Since we made a number of hypotheses about the optimal trajectory, especially at junction points, it is useful to give an example where these hypotheses are satisfied. For that purpose we discuss a



particular case in which the original optimal control problem can be reduced to the optimal control of a scalar ODE.

Let  $\Omega = (0, 1)$ , and denote by  $c_1(x) := \sqrt{2} \sin \pi x$  the first (normalized) eigenvector of the Laplace operator.

We assume that  $\gamma = 0$ , the control is scalar ( $m = 1$ ),  $b_0 \equiv 0$  and  $b_1 \equiv 1$  in  $\Omega$ , and that  $f \equiv 0$  in  $Q$ . Then the state equation with initial condition  $c_1$  reads

$$\dot{y}(x, t) - \Delta y(x, t) = u(t)y(x, t); \quad (x, t) \in (0, 1) \times (0, T), \quad y(x, 0) = c_1(x), \quad x \in \Omega. \quad (\text{B.1})$$

It is easily seen that the state satisfies  $y(x, t) = y_1(t)c_1(x)$ , where  $y_1$  is solution of

$$\dot{y}_1(t) + \pi^2 y_1(t) = u(t)y_1(t); \quad t \in (0, T), \quad y_1(0) = y_{10} = 1. \quad (\text{B.2})$$

We set  $T = 3$  and consider the state constraint (3.17) with  $q = 1$  and  $d_1 := -2$ , and the cost function (2.5) with  $\alpha_1 = 0$ . The state constraint reduces to

$$y_1(t) \leq 2, \quad t \in [0, 3]. \quad (\text{B.3})$$

As target functions take  $y_{dT} := c_1$  and  $y_d(x, t) := \hat{y}_d(t)c_1(x)$  with

$$\hat{y}_d(t) := \begin{cases} 1.5e^t & \text{for } t \in (0, \log 2), \\ 3 & \text{for } t \in (\log 2, 1), \\ 4 - t & \text{for } t \in (1, 3). \end{cases} \quad (\text{B.4})$$

We assume that the lower and upper bounds for the control are  $\check{u} := -1$  and  $\hat{u} := \pi^2 + 1$ . We will check that the optimal control is

$$\bar{u}(t) := \begin{cases} \hat{u} & \text{for } t \in (0, \log 2), \\ \pi^2 & \text{for } t \in (\log 2, 2), \\ \pi^2 - 1/\hat{y}_d & \text{for } t \in (2, 3). \end{cases} \quad (\text{B.5})$$

Thus, for the optimal state we have

$$\bar{y}_1(t) := \begin{cases} e^t & \text{for } t \in (0, \log 2), \\ 2 & \text{for } t \in (\log 2, 2), \\ 4 - t & \text{for } t \in (2, 3). \end{cases} \quad (\text{B.6})$$

The above control is feasible. The trajectory  $(\bar{u}, \bar{y})$  is optimal since for any  $t \in (0, T)$ , the state  $\bar{y}_1(t)$  has the best possible value (in order to approach  $\hat{y}_d$  and minimize the cost function) that respects the state constraint.

Let us check Hypothesis 4.1 for this example. Conditions 1 and 2 are obviously satisfied. For the constraint qualification in Condition 3 consider the linearized state equation with unique  $z_1[v]$ :

$$\dot{z}_1 = (\bar{u} - \pi^2)z_1 + v\bar{y}_1; \quad z_1(0) = 0, \quad (\text{B.7})$$

with  $v(t) := \check{u} - \bar{u}(t) < 0$ . One easily checks that  $z_1[v](t) < 0$  for all  $t > 0$ . Hence, we can find  $\varepsilon > 0$  such that

$$g_1(\bar{y}(\cdot, t)) + g_1'(\bar{y}(\cdot, t))z_1[v](\cdot, t) = \bar{y}_1(t) - 2 + z_1(t) < -\varepsilon, \quad \text{for all } t \in (0, T). \quad (\text{B.8})$$

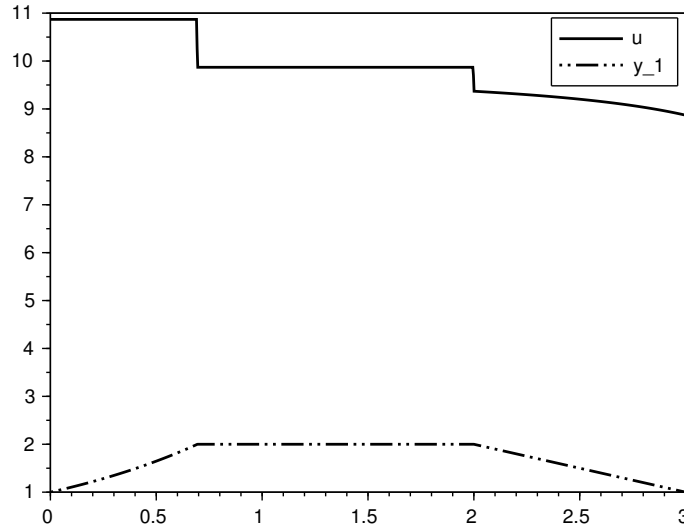


Figure 1: Optimal control and state for the example

Conditions 4 holds, since

$$M(t) = \bar{M}_1(t) = \int_{\Omega} c_1(x) \bar{y}(x, t) dx = \bar{y}_1(t) > 0 \quad \text{for } t \in (0, T). \quad (\text{B.9})$$

For Condition 5 we have

$$\text{dist}(t, I_1^C) = \begin{cases} \log 2 - t & \text{for } t \in (0, \log 2), \\ 0 & \text{for } t \in (\log 2, 2), \\ t - 2 & \text{for } t \in (2, 3), \end{cases} \quad (\text{B.10})$$

and hence,

$$g_1(\bar{y}(\cdot, t)) = \bar{y}_1(t) - 2 \leq -\text{dist}(t, I_1^C). \quad (\text{B.11})$$

Conditions 6 and 8 hold by the choice of the control in (B.5). Condition 7 holds by definition.

We solve this problem numerically using BOCOP [5] and get the optimal control and state given in Figure 1.

We now discuss the second order optimality condition for this example. The costate equation is

$$-\dot{p} + Ap = c_1(\bar{y}_1 - \hat{y}_d) + c_1 \dot{\mu}_1, \quad p(\cdot, T) = \bar{y}(T) - y_{dT} = 0 \quad (\text{B.12})$$

with  $A$  as defined in (2.20). Since  $\bar{y}$  and  $y_d$  are colinear to  $c_1$ , it follows that  $p(x, t) = p_1(t)c_1(x)$ , and

$$-\dot{p}_1 + \pi^2 p_1 = \bar{u} p_1 + \bar{y}_1 - \hat{y}_d + \dot{\mu}_1; \quad p_1(3) = 0. \quad (\text{B.13})$$

Over  $(2, 3)$ ,  $\dot{\mu}_1 = 0$  (state constraint not active) and  $\bar{y}_1 = \hat{y}_d$ , therefore  $p_1$  and  $p$  identically vanish. Over  $(\log 2, 2)$ ,  $\bar{u}$  is out of bounds and therefore

$$0 = \int_{\Omega} p(x, t) \bar{y}(x, t) = p_1(t) \bar{y}_1(t) \int_{\Omega} c_1(x)^2 = 2p_1(t). \quad (\text{B.14})$$

It follows that  $p_1$  and  $p$  also vanish on  $(\log 2, 2)$  and that

$$\dot{\mu}_1 = -(\bar{y}_1 - \hat{y}_d) > 0, \quad \text{a.a. } t \in (\log 2, 2). \quad (\text{B.15})$$

Over  $(0, \log 2)$ , the control attains its upper bound, then

$$-\dot{p}_1 = p_1 - \frac{1}{2}e^t \quad (\text{B.16})$$

with final condition  $p_1(\log 2) = 0$ , so that

$$p_1(t) = \frac{e^t}{4} - e^{-t}. \quad (\text{B.17})$$

As expected,  $p_1$  is negative.

Next, the linearized state equation at  $(\bar{u}, \bar{y})$  reads

$$\dot{z} - \Delta z = \bar{u}z + v\bar{y}; \quad z(\cdot, 0) = 0. \quad (\text{B.18})$$

Since  $\bar{y} = \bar{y}_1(t)c_1(x)$ , we deduce that  $z = z_1(t)c_1(x)$ , with  $z_1$  solution of

$$\dot{z}_1 + \pi^2 z_1 = \bar{u}z_1 + v\bar{y}_1; \quad z_1(0) = 0. \quad (\text{B.19})$$

Therefore if  $(v, z)$  satisfy the linearized state equation

$$\mathcal{Q}[p](z, v) = \int_Q (z^2 + pvz) dx dt + \int_\Omega z(x, T)^2 dx = \int_0^3 (z_1(t)^2 + p_1(t)v(t)z_1(t)) dt + z_1(3)^2. \quad (\text{B.20})$$

If in addition  $v$  is a critical direction, since  $v = 0$  and  $z_1 = 0$  a.e. on  $(0, 2)$ , and  $p_1(t) = 0$  on  $(2, 3)$ , we get

$$\mathcal{Q}[p](z, v) = \int_2^3 z_1(t)^2 dt + z_1(3)^2. \quad (\text{B.21})$$

Thus,  $\mathcal{Q}$  is non-negative for any critical directions  $(z[v], v)$ , in accordance with the second-order necessary condition of Theorem 4.7.

## References

- [1] M. S. Aronna, J. F. Bonnans, A. V. Dmitruk, and P. A. Lotito. Quadratic order conditions for bang-singular extremals. *Numerical Algebra, Control and Optimization, AIMS Journal*, 2(3):511–546, 2012.
- [2] M. S. Aronna, J. F. Bonnans, and A. Kröner. State-constrained control-affine parabolic problems II: Second-order sufficient optimality conditions. 2019.
- [3] M. S. Aronna, J.F. Bonnans, and B. S. Goh. Second order analysis of control-affine problems with scalar state constraint. *Math. Program.*, 160(1-2, Ser. A):115–147, 2016.
- [4] J.-P. Aubin. Un théorème de compacité. *C. R. Acad. Sci. Paris*, 256:5042–5044, 1963.
- [5] J. Bonnans, J.F., D. Giorgi, V. Grélard, B. Heymann, S. Maindrault, P. Martinon, O. Tissot, and J. Liu. Bocop – A collection of examples. Technical report, INRIA, 2017.

- [6] J.F. Bonnans. Second-order analysis for control constrained optimal control problems of semilinear elliptic systems. *Appl. Math. Optim.*, 38(3):303–325, 1998.
- [7] J.F. Bonnans and A. Hermant. Second-order analysis for optimal control problems with pure state constraints and mixed control-state constraints. *Ann. Inst. H. Poincaré Anal. Non Linéaire*, 26(2):561–598, 2009.
- [8] J.F. Bonnans and P. Jaisson. Optimal control of a parabolic equation with time-dependent state constraints. *SIAM J. Control Optim.*, 48(7):4550–4571, 2010.
- [9] J.F. Bonnans and A. Shapiro. *Perturbation analysis of optimization problems*. Springer Series in Operations Research. Springer-Verlag, New York, 2000.
- [10] E. Casas, J.C. de Los Reyes, and F. Tröltzsch. Sufficient second-order optimality conditions for semilinear control problems with pointwise state constraints. *SIAM J. Optim.*, 19(2):616–643, 2008.
- [11] E. Casas, Mariano Mateos, and Fredi Tröltzsch. Necessary and sufficient optimality conditions for optimization problems in function spaces and applications to control theory. In *Proceedings of 2003 MODE-SMAI Conference*, volume 13 of *ESAIM Proceedings*, pages 18–30. EDP Sciences, 2003.
- [12] E. Casas, F. Tröltzsch, and A. Unger. Second order sufficient optimality conditions for a nonlinear elliptic control problem. *J. for Analysis and its Applications (ZAA)*, 15:687–707, 1996.
- [13] E. Casas and F. Tröltzsch. Recent advances in the analysis of pointwise state-constrained elliptic optimal control problems. *ESAIM Control Optim. Calc. Var.*, 16(3):581–600, 2010.
- [14] J.C. de Los Reyes, P. Merino, J. Rehberg, and F. Tröltzsch. Optimality conditions for state-constrained PDE control problems with time-dependent controls. *Control and Cybernetics*, 37(1):5–38, 2008.
- [15] A.V. Dmitruk. Quadratic conditions for a weak minimum for singular regimes in optimal control problems. *Soviet Math. Doklady*, 18(2):418–422, 1977.
- [16] A.V. Dmitruk. Jacobi type conditions for singular extremals. *Control & Cybernetics*, 37(2):285–306, 2008.
- [17] L.C. Evans. *Partial differential equations*. Amer. Math Soc., Providence, RI, 1998. Graduate Studies in Mathematics 19.
- [18] B.S. Goh. Necessary conditions for singular extremals involving multiple control variables. *SIAM J. Control*, 4:716–731, 1966.
- [19] H.J. Kelley. A second variation test for singular extremals. *AIAA Journal*, 2:1380–1382, 1964.
- [20] K. Krumbiegel and J. Rehberg. Second order sufficient optimality conditions for parabolic optimal control problems with pointwise state constraints. *SIAM J. Control Optim.*, 51(1):304–331, 2013.
- [21] Gary M. Lieberman. *Second order parabolic differential equations*. World Scientific Publishing Co., Inc., River Edge, NJ, 1996.
- [22] J.-L. Lions. *Quelques méthodes de résolution des problèmes aux limites non linéaires*. Dunod, Paris, 1969.

- [23] J.-L. Lions. *Contrôle des systèmes distribués singuliers*, volume 13 of *Méthodes Mathématiques de l'Informatique*. Gauthier-Villars, Montrouge, 1983.
- [24] J.-L. Lions and E. Magenes. *Problèmes aux limites non homogènes et applications. Vol. 1*. Dunod, Paris, 1968.
- [25] H. Maurer. On optimal control problems with bounded state variables and control appearing linearly. *SIAM J. Control Optimization*, 15(3):345–362, 1977.
- [26] H. Maurer. On the minimum principle for optimal control problems with state constraints. *Schriftenreihe des Rechenzentrum 41*, Universität Münster, 1979.
- [27] H. Maurer, J.-H. R. Kim, and G. Vossen. *On A State-Constrained Control Problem in Optimal Production and Maintenance*, pages 289–308. Springer US, Boston, MA, 2005.
- [28] J.P. McDanell and W.F. Powers. Necessary conditions for joining optimal singular and nonsingular subarcs. *SIAM J. Control*, 9:161–173, 1971.
- [29] J.-P. Raymond and F. Tröltzsch. Second order sufficient optimality conditions for nonlinear parabolic control problems with state constraints. *Discrete Contin. Dynam. Systems*, 6(2):431–450, 2000.
- [30] H. Schättler. Local fields of extremals for optimal control problems with state constraints of relative degree 1. *J. Dyn. Control Syst.*, 12(4):563–599, 2006.