

# Some results on the optimal control with unilateral state constraints

Bernard Brogliato

► **To cite this version:**

| Bernard Brogliato. Some results on the optimal control with unilateral state constraints. [Research Report] RR-5992, INRIA. 2006. <inria-00103775v2>

**HAL Id: inria-00103775**

**<https://hal.inria.fr/inria-00103775v2>**

Submitted on 12 Oct 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

*Some results on the optimal control with unilateral  
state constraints*

Bernard Brogliato

**N° 5992**

Octobre 2006

Thème NUM



*Rapport  
de recherche*



## Some results on the optimal control with unilateral state constraints

Bernard Brogliato\*

Thème NUM — Systèmes numériques  
Projets Bipop

Rapport de recherche n° 5992 — Octobre 2006 — 44 pages

**Abstract:** In this paper we study the problem of quadratic optimal control with state variables unilateral constraints, for linear time-invariant systems. The necessary conditions are formulated as a linear invariant system with complementary slackness conditions. Some structural properties of this system are examined. Then it is shown that the problem can benefit from the higher order Moreau's sweeping process, that is a specific distributional differential inclusion, and from ten Dam's geometric theory for the partitioning of the admissible domain boundary. In fact the first step may be also seen as follows: does the higher order Moreau's sweeping process (developed in [1, 2]) correspond to the necessary conditions of some optimal control problem with an adapted integral action? The knowledge of the qualitative behaviour of optimal trajectories is improved with the approach, which also paves the way towards efficient time-stepping numerical algorithms to solve the optimal control boundary value problem.

**Key-words:** Complementarity systems, Convex analysis, Measure differential inclusions, Zero dynamics, Time-stepping algorithm, Optimal control, State constraints, Boundary Value Problem, Moreau's sweeping process.

\* INRIA Rhône-Alpes, 655 avenue de l'Europe, 38334 Saint Ismier, France, Bernard.Brogliato@inrialpes.fr

## Quelques résultats sur la commande optimale avec contraintes d'inégalités sur l'état

**Résumé :** Dans cet article nous étudions la commande optimale quadratique avec contraintes d'inégalités sur l'état. Les conditions nécessaires sont formulées comme un système linéaire de complémentarité. Quelques propriétés structurelles de ce système sont examinées. Nous montrons ensuite comment ce problème peut bénéficier d'une étude via le processus de raffinement d'ordre élevé, qui est une inclusion différentielle à distribution. La théorie géométrique de Ten Dam est ensuite utilisée pour étudier les propriétés qualitatives des trajectoires optimales.

**Mots-clés :** Systèmes de complémentarité, analyse convexe, inclusions différentielles à mesure, dynamique zéro, algorithme à pas de temps, contraintes sur l'état, problème à valeurs frontières, processus de raffinement de Moreau.

## 1 Introduction

Optimal control with state constraints is a topic of major importance, and which has attracted the attention of researchers since a long time [57], see [34, 28, 13, 26] to cite a few. Most of these works study the first order necessary conditions stemming from Boltyanskii-Pontryagin's maximum principle. Dynamic programming under inequality state constraints is examined in [5, 58], and second order optimality conditions are studied in [31]. The qualitative properties of optimal trajectories are studied in [28, 34, 25], while their regularity properties are examined in [58]. Numerical studies have received attention in [42, 56] and applications are presented in [9, 8, 7, 5, 55]. The optimal control of systems with inputs which are measures has also received attention, without state constraints [51, 43, 53], and with state constraints [45, 35]. The theory of systems with distributional inputs is progressing significantly [17, 43, 21], so that it becomes possible to consider control problems involving inputs which are higher order distributions basing on a solid ground. A further motivation for considering higher degree distributions, will be explained next. However as we shall see later, even problems with a continuous optimal controller, may involve higher degree distributions in the costate differential equation.

The following Bolza problem is of interest

$$\text{minimize}_{u(\cdot) \in \mathbf{U}} I(u) = \frac{1}{2} \int_0^{T_1} [x(t)^T Q x(t) + u(t)^T R u(t)] dt + \frac{1}{2} x(T_1)^T F x(T_1) \quad (1)$$

subject to

$$\begin{cases} \dot{x}(t) = Ax(t) + Bu(t), & x(0) = \bar{x}_0, & x(T_1) = \bar{x}_1 \\ w(t) = Cx(t) + D \geq 0 \end{cases} \quad (2)$$

where  $A, B, C$  and  $D$  are constant matrices,  $(A, B, C)$  is a minimal state space representation,  $\bar{x}_0, \bar{x}_1 \in \Phi = \{x \mid Cx + D \geq 0\}$ ,  $\mathbf{U}$  is the set of admissible inputs,  $w(t) \in \mathbb{R}^m$ ,  $u(t) \in \mathbb{R}^{n_u}$ ,  $x(t) \in \mathbb{R}^n$  and  $Q = Q^T \geq 0$ ,  $R = R^T > 0$  (the Legendre-Clebsch condition). This is a particular example of optimal control of a dynamical system with state constraints. The least requirement for problem (1)-(2) to possess a solution in  $\mathbf{U}$  is that  $\bar{x}_1$  belongs to the reachable set from  $\bar{x}_0$ ,  $0 \leq T_1 < +\infty$ . A fundamental ingredient is therefore  $\mathbf{U}$ , as illustrated by example 6. This example shows that distributions can easily appear in unilaterally constrained controlled systems. Clearly then the formulation of the optimal control problem in (1)-(2) will have to be modified so that the integral action makes sense. This will be the object of the first part of this paper and detailed in the next sections.

It is noteworthy the strong similarities between this optimal control problem, and the so-called continuous linear programming problem (CTLP) [3, 46], introduced by Bellman [6] to model some economic processes

$$\begin{cases} \text{maximize} & \int_{[0, T_1]} c^T(t) x(t) dt \\ \text{subject to:} & H(t)x(t) + \int_0^t G(s, t)x(s) ds \leq a(t) \\ & x(t) \geq 0, t \in [0, T_1] \end{cases} \quad (3)$$

When the integrand of  $I(u)$  in (1) is linear in  $x(\cdot)$  and  $u(\cdot)$ , then the linear optimal control problem with both state and input constraints, is a particular case of the CTLP in (3) [46]. The optimal control problem (1) (2) and the dual problem of the CTLP in (3) share important features concerning the presence of higher degree distributions in their solution [46] [3, §2.2], and the need for time-stepping numerical algorithms to numerically solve them. Some macroeconomics models directly yield optimal control problems with unilateral state constraints and optimal controllers with singular measures [21].

This is why we will embed the optimization problem (1) (2) into a suitable framework which allows us to rigorously take into account the possible presence of higher degree distributions.

The first order necessary conditions for the optimal control problem in (1) (2) can be formulated as [26, 47, 50, 58]

$$\begin{cases} \begin{pmatrix} \dot{x}(t) \\ \dot{\eta}(t) \end{pmatrix} = \begin{pmatrix} A & BR^{-1}B^T \\ Q & -A^T \end{pmatrix} \begin{pmatrix} x(t) \\ \eta(t) \end{pmatrix} + \begin{pmatrix} 0 \\ -C^T \end{pmatrix} \lambda(t) & \text{(a)} \\ x(0) = \bar{x}_0, \eta(T_1) = Fx(T_1) + C^T\gamma + \beta = F\bar{x}_1 + C^T\gamma + \beta = \eta_1 & \text{(b)} \\ 0 \leq w(t) = Cx(t) + D \perp \lambda(t) \geq 0 & \text{(c)} \end{cases} \quad (4)$$

where  $\eta(t) \in \mathbb{R}^n$ ,  $0 \leq \gamma \in \mathbb{R}^m$ ,  $\gamma^T(Cx(T_1) + D) = 0$ ,  $\beta \in \mathbb{R}^n$ , and the optimal control is given on intervals where  $\eta(\cdot)$  is a function by

$$u(t) = \operatorname{argmax}_{u \in \mathbf{U}} \left[ -\frac{1}{2}u^T Ru + \eta(t)^T Bu \right] = R^{-1}B^T \eta(t) \quad (5)$$

An implicit assumption which allows one to write both (1) and (4) is that the multiplier  $\lambda$  is a measure. If not then both (1) and (4) have to be generalized to be meaningful. This is the first objective of this work. In this paper the optimal trajectory and optimal controller will not be given a special notation in order to lighten the presentation. The values of the multipliers  $\gamma$  and  $\beta$  can be calculated from some boundary conditions [13, §2.8]. If  $x(T_1)$  is not specified and  $F = 0$  then  $\eta(T_1) = 0$ . We assume that the Slater constraint qualification ( $\exists x$  such that  $Cx + D > 0$ ) is satisfied and that the problem (1) (2) is normal, a property that is related to the controllability of the constrained system (see e.g. [47, Proposition 4.1] in case  $\lambda$  is a measure). We conjecture that reachability and normality hold for  $r^{wu} \geq 2$  by taking  $u(\cdot)$  in an enlarged set of distributions that will be described later (see section 3.5). From a general point of view however, the controllability problem for (2) remains open and we will not tackle it here (see Appendix C for some preliminary hints). In the sequel we shall denote  $\tilde{x}^T = (x^T, \eta^T)$ ,  $(\tilde{A}, \tilde{B}, \tilde{C})$  the triple associated to the system in (4) (a) (c), which we shall often refer to in the sequel as the *necessary conditions system*. It is clear that if  $\lambda(t)$  is a measure, then from (4) (c)

$$-C^T \lambda(t) \in \partial \psi_{\Phi}(x(t)) \quad (6)$$

(using [23, Proposition 1.3.11] on differentiation of convex functions) and

$$\operatorname{supp}(\lambda) \subset \{t \mid Cx(t) + D = 0\} \quad (7)$$

It is noteworthy that when  $\lambda$  is not a function (e.g. a Dirac measure), then the meaning of the inclusion in (6) has to be stated rigorously, see [1, Section 2] or [32, p.76]. Then the dynamical system in (4) (a) (c) may be viewed as a measure differential inclusion

$$\dot{\tilde{x}}(t) - \tilde{A}\tilde{x}(t) \in \partial \psi_{\tilde{\Phi}}(\tilde{x}(t)) \quad (8)$$

with  $\tilde{\Phi} = \{\tilde{x} \mid \tilde{C}\tilde{x} + D \geq 0\}$ , and where a state reinitialisation mapping is missing in (8) (let us notice from (7) that the control input  $u(\cdot)$  must be a continuous function of time (even analytic) everywhere outside  $\operatorname{supp}(\lambda)$ ). However as we shall see in section 3, in general one has to resort to a more complex formalism to correctly handle (4). Except if additional constraints are added to (4) to assure that  $\lambda$  is a measure, then the formalism in (4) and especially the *complementary slackness* condition (4) (c) are meaningless. When  $m = 1$ , the relative degree  $r$  of the triple  $(\tilde{A}, \tilde{B}, \tilde{C})$  in (4) is twice the relative degree  $r^{wu}$  of the triple  $(A, B, C)$  in (2) [52] (what we call the relative degree in this paper is often referred to as the constraints order in the optimal control literature [28]). We will say that the triple  $(A, B, C)$  has a uniform vector relative degree  $r^{wu}$  if the matrix  $CA^{r^{wu}-1}B \in \mathbb{R}^{m \times m}$  whose entries are the scalars  $C_i A^{r^{wu}-1} B_j$ , is full rank, and  $C_i A^{k-1} B_j = 0$  for all  $1 \leq i \leq m$ ,  $1 \leq j \leq m$ , and  $k \leq r^{wu} - 2$ .

**Remark 1** It is known [50, 58] that (4) (a) (c) is a Hamiltonian system of the form  $\dot{\tilde{x}}(t) = J \begin{pmatrix} \frac{\partial H(x,\eta)}{\partial x} \\ \frac{\partial H(x,\eta)}{\partial \eta} \end{pmatrix} + \begin{pmatrix} 0_{n \times n} \\ -I_n \end{pmatrix} C^T \lambda$ ,  $J = \begin{pmatrix} 0_{n \times n} & I_n \\ -I_n & 0_{n \times n} \end{pmatrix}$ , whose unconstrained part has the Hamiltonian function

$$H(x, \eta) = \frac{1}{2} \tilde{x}^T \begin{pmatrix} -Q & A^T \\ A & BR^{-1}B^T \end{pmatrix} \tilde{x} \quad (9)$$

Notice that in (5)  $u(t) = \operatorname{argmax}_{u \in \mathbf{U}} H(x(t), u, \eta(t))$  with

$$H(x, u, \eta) = -\frac{1}{2} x^T Q x - \frac{1}{2} u^T R u + \eta^T (Ax + Bu). \quad (10)$$

Let us assume that  $\lambda$  is a measure. Defining the nonsmooth Hamiltonian function  $H_{ns}(x, \eta) = H(x, \eta) + \psi_{\Phi}(x)$  (equivalently  $H_{ns}(x, \eta) = H(x, \eta) - \lambda^T (Cx + D)$  and  $\lambda$  in (4) c)), one sees that (8) can also be rewritten as  $(-\dot{\eta}(t), \dot{x}(t)) \in \partial H_{ns}(x(t), \eta(t))$ , where  $\partial$  is the subdifferential of convex analysis.

**Remark 2** The support condition (7) is fundamental. It will also hold when  $\lambda$  is a distribution of higher degree. It implies that in the problem we are examining, the control  $u$  and costate  $\eta$  are allowed to have a degree larger than 1, only at junction times.

In this paper we provide a framework that allows us to give a meaning to the dynamical system in (4) (a) (c) when  $\lambda$  is not a measure; in particular new multipliers are introduced and it is shown that this framework may be useful even if the optimal controller  $u(\cdot)$  is a function (but the costate  $\eta$  may be a distribution). This allows us to introduce a generalized action  $I(u)$  that handles distributional cases. Then we provide a detailed study of what happens at the entry times, using the geometrical approach of ten Dam. This allows us to generalize some results on the existence (or non-existence) of entry states for odd relative degree systems, and to tackle the multivariable case. Finally a study of boundary arcs is proposed using complementarity theory. All these tools are introduced for the first time in the context of optimal control with inequality state constraints and are proved to significantly improve the qualitative knowledge of the optimal trajectories. The approach also paves the way towards the design of time-stepping numerical algorithms to solve the BVP in (4), thereby providing nice alternative to multiple shooting schemes.

The paper is organised as follows: in section 2 some properties of the system in (4) (a) are proved. In section 3 a distributional differential inclusion is presented in which the system in (4) (a) (c) is embedded; the properties of its solutions are recalled. In section 4 the geometric approach of ten Dam [20, 19] is used to study the behaviour of the trajectories of (2) at junction times. In section 5 the behaviour of the system when trajectories of (4) (a) evolve on the boundary of  $\Phi$ , is analysed. Some mathematical definitions are recalled in Appendix A. An example of a function that belongs to the considered space of solutions is given in Appendix B. Preliminary results on controllability of (2) are in Appendix C.

**Notation and nomenclature:** entry time and state (state on the boundary followed by a boundary arc), exit time and state (state on the boundary followed by an interior arc), contact time and state (state on the boundary coming from an interior arc), touch time and state (state coming from an interior arc and followed by an interior arc), junction time and state (entry, contact, exit, touch time and state), CP (complementarity problem), LCP (linear CP),  $\psi_K$  (the indicator function of a set  $K$  [27, p.82]),  $\partial\psi_K$  (the subdifferential of the indicator of a convex set  $K$  [23, p.67]),  $\operatorname{supp}(T)$  (the support of a distribution  $T$  [22, p.142]),  $x \succ 0$  ( $x$  is not zero and the first nonzero element of the vector  $x$  is positive),  $\succcurlyeq$  (the first nonzero element of the vector  $x$  is positive), similarly for  $\prec$  and  $\preccurlyeq$ ,  $\sigma_f(t) = f(t^+) - f(t^-)$  the jump of the function  $f(\cdot)$  at time  $t$ , BV (bounded variation), RCLBV (right continuous of local bounded variation), RCSLBV (right continuous of special local bounded variation), BVP (boundary value problem), IVP (initial value problem),  $\delta_t$  the Dirac measure at  $t$ ,  $0^m = (0, \dots, 0)^T \in \mathbb{R}^m$ ,  $0_m = (0^m)^T$ ,  $I_n$  is the  $n$ -identity matrix. Contact times are denoted as  $\tau$  (when only one contact is analysed) or  $t_k$  when a sequence of contact times is examined,  $E_0(h)$  denotes the set of jumps of the function  $h(\cdot)$ .  $e_i$  is the  $i$ -th unit vector of  $\mathbb{R}^n$ ; a function in  $C_0^\infty(I)$  is infinitely differentiable and with compact support.



## 2 Some properties of the necessary conditions system

The next Lemma is important for the well-posedness of (4) (a) seen as an IVP ( $x(0^-) = x_0$ ,  $\eta(0^-) = \eta_0$ ), and will also be useful for the qualitative analysis of the BVP solutions, see sections 4 and 5. We consider  $R = I_n$  and  $m = 1$  in Lemma 1.

**Lemma 1** *If  $r^{wu} = 1$ , then the leading Markov parameter of the triple  $(\tilde{A}, \tilde{B}, \tilde{C})$  is  $M^{(2)} = -CBB^T C^T < 0$ . More generally if the transfer function of the triple  $(A, B, C)$  in (2) has relative degree  $r^{wu} \geq 2$ , the leading Markov parameter of the triple  $(A, \tilde{B}, \tilde{C})$  is  $M^{(2r^{wu})} = (-1)^{r^{wu}} CA^{r^{wu}-1} B (CA^{r^{wu}-1} B)^T (= \tilde{C} \tilde{A}^{r-1} \tilde{B})$ .*

**Proof:** Let us denote  $\tilde{A} = \begin{pmatrix} A & BB^T \\ Q & -A^T \end{pmatrix} = \begin{pmatrix} \tilde{A}_1 & \tilde{A}_2 \\ \tilde{A}_3 & \tilde{A}_4 \end{pmatrix}$  and  $\tilde{A}^r = \begin{pmatrix} \tilde{A}_1^{(r)} & \tilde{A}_2^{(r)} \\ \tilde{A}_3^{(r)} & \tilde{A}_4^{(r)} \end{pmatrix}$  for some  $r \geq 1$ .

The first assertion of the Lemma is a simple calculation with  $\tilde{A}_2^{(1)} = BB^T$  and  $M^{(1)} = -C\tilde{A}_2^{(1)}C^T$ . The leading Markov parameter we want to compute is  $M^{(2r+2)} = (C \ 0)\tilde{A}^{2r+1}(0 \ -C)^T = -C\tilde{A}_2^{(2r+1)}C^T$ . If  $r^{wu} = \alpha + 1$ , then  $CA^i B = 0$  for all  $0 \leq i \leq \alpha - 1$ . Our objective is to show that  $M^{(2\alpha+2)} = (-1)^{\alpha+1} CA^\alpha BB^T (A^T)^\alpha C^T$ ,  $\alpha \geq 1$ . By direct calculation we have for any integer  $i \geq 1$

$$\begin{aligned} \tilde{A}_1^{(i)} &= \tilde{A}_1^{(i-1)} A + \tilde{A}_2^{(i-1)} Q \\ \tilde{A}_2^{(i+1)} &= \tilde{A}_1^{(i)} BB^T - \tilde{A}_2^{(i)} A^T \end{aligned} \quad (11)$$

It follows that

$$\tilde{A}_2^{(2\alpha)} = \sum_{i=1}^{2\alpha-1} (-1)^{i+1} \tilde{A}_1^{(i)} BB^T (A^T)^{2\alpha-i-1} - \tilde{A}_2^{(1)} (A^T)^{2\alpha-1} \quad (a)$$

(12)

$$\tilde{A}_2^{(2\alpha+1)} = \sum_{i=1}^{2\alpha} (-1)^i \tilde{A}_1^{(i)} BB^T (A^T)^{2\alpha-i} + \tilde{A}_2^{(1)} (A^T)^{2\alpha} \quad (b)$$

This can be shown using (11) and then by induction. From (11) it follows by induction that  $\tilde{A}_1^{(\alpha)} = A^\alpha + L_\alpha(Q)$  where  $L_\alpha(Q)$  is some matrix depending on the matrix  $Q$ . Indeed let us assume that  $\tilde{A}_1^{(i-1)} = A^{i-1} + L_{i-1}(Q)$ . Then  $\tilde{A}_1^{(i)} = (A^{i-1} + L_{i-1}(Q))A + \tilde{A}_2^{(i-1)}Q = A^i + L_i(Q)$  with  $L_i(Q) = L_{i-1}(Q)A + \tilde{A}_2^{(i-1)}Q$ . Using (12) and since  $\tilde{A}_1 = A$  and  $\tilde{A}_1^{(2)} = A^2 + BB^T Q$  the proof is complete. From (12) (b) one deduces that  $M^{(2\alpha+2)} = \sum_{i=1}^{2\alpha} (-1)^{i+1} C\tilde{A}_1^{(i)} BB^T (A^T)^{2\alpha-i} C^T + C\tilde{A}_2^{(1)} (A^T)^{2\alpha} C^T$ . Assume for the time being that  $CL_i(Q)BB^T (A^T)^{2\alpha-i} C^T = 0$  for all  $1 \leq i \leq \alpha$  (the terms with  $\alpha+1 \leq i \leq 2\alpha$  need not be considered as  $CA^i B = 0$  for all  $0 \leq i \leq \alpha - 1$ ). Then we obtain  $M^{(2\alpha+2)} = \sum_{i=1}^{2\alpha} (-1)^{i+1} CA^i BB^T (A^T)^{2\alpha-i} C^T$  and the only nonzero term in this sum is for  $i = \alpha$ . Therefore  $M^{(2\alpha+2)} = (-1)^{\alpha+1} CA^\alpha BB^T (A^T)^\alpha C^T$ .

To end the proof, it can be shown again by induction using (12) that  $L_i(Q)$  is composed of terms of the general form  $A^j B \star$  or  $\star B^T (A^T)^j$ , with  $0 \leq j \leq i-2$  and  $\star$  is some matrix. Indeed one has  $L_i(Q) = L_{i-1}(Q)A + \tilde{A}_2^{(i-1)}Q$ , as shown above. Assume that  $L_{i-1}(Q) = \sum_{j=0}^{i-3} A^j B \star + \star B^T (A^T)^j$ . Then  $L_i(Q) = \sum_{j=0}^{i-3} A^j B \star A + \star B^T (A^T)^{j+1} + \tilde{A}_2^{(i-1)}Q$ . One concludes from (12). Thus  $CL_i(Q)BB^T (A^T)^{2\alpha-i} C^T = 0$  for all  $1 \leq i \leq \alpha$  as required. ■

Lemma 1 continues to hold for systems with a uniform relative degree and  $m \geq 1$  (recall that  $C \in \mathbb{R}^{m \times n}$ ) since  $CA^iB = 0$  for all  $0 \leq i \leq r^{wu} - 2$  in this multivariable case also. In such a case, having  $M^{(2r^{wu})} \in \mathbb{R}^{m \times m}$  full-rank implies that the vectors  $C_i^T$  are independent, where  $C_i$  is the  $i$ -th row of  $C$ , and  $B$  has rank  $m$ .

**Example 1** Consider the system

$$\begin{cases} \dot{x}_1(t) = x_1(t) + x_2(t) \\ \dot{x}_2(t) = x_3(t) \\ \dot{x}_3(t) = x_4(t) + x_3(t) \\ \dot{x}_4(t) = x_4(t) + x_6(t) \\ \dot{x}_5(t) = x_5(t) + x_1(t) \\ \dot{x}_6(t) = u(t) \\ w(t) = x_1(t) \geq 0 \end{cases} \quad (13)$$

Then  $r^{wu} = 5$ ,  $CA^4B = 1$  and  $M^{(10)} = -1$ .

Let us introduce a canonical state space representation for the system  $(\tilde{A}, \tilde{B}, \tilde{C}, D)$ . This is a canonical representation which makes the so-called zero-dynamics explicitly appear. It will be useful for some subsequent developments (see also remark 5 iv)). For a triple  $(\tilde{A}, \tilde{B}, \tilde{C})$ , with input  $\lambda \in \mathbb{R}^m$ , output  $w \in \mathbb{R}^m$ , and uniform relative degree  $r \geq 1$ , it reads [48]:

$$\begin{cases} \dot{z}_1(t) = z_2(t) \\ \dot{z}_2(t) = z_3(t) \\ \dot{z}_3(t) = z_4(t) \\ \vdots \\ \dot{z}_{r-1}(t) = z_r(t) \\ \dot{z}_r(t) = \tilde{C}\tilde{A}^r\tilde{W}^{-1}\tilde{z}(t) + \tilde{C}\tilde{A}^{r-1}\tilde{B}\lambda(t) \\ \dot{\xi}(t) = \tilde{A}_\xi\xi(t) + \tilde{B}_\xi z_1(t) \\ \\ w(t) = z_1(t) + D, \quad z_1(t) = Cx(t) \\ z(0^-) = z_0. \end{cases} \quad (14)$$

We call it the  $\tilde{z}$ -dynamics, and the full-rank state transformation matrix is  $\tilde{W}$ :  $\tilde{z} = \tilde{W}\tilde{x}$  <sup>(1)</sup>. We denote  $\tilde{z}^T = (z_1, \dots, z_r)$ . In the sequel, the components of  $\tilde{z}$  and of  $z$  will both be denoted as  $z_i$ , except if not clear from the context. Since the zero dynamics plays an important role in the systems we are dealing with, the following is of interest.

**Lemma 2** Let  $m = 1$ . Consider the  $z$ - and  $\tilde{z}$ -canonical forms that correspond to the systems  $(A, B, C)$  in (2) and  $(\tilde{A}, \tilde{B}, \tilde{C})$  in (4), and let us denote the transitions matrices of their zero dynamics as  $A_\xi$  and  $\tilde{A}_\xi$  respectively. Then if  $\sigma$  is an eigenvalue of  $A_\xi$  with multiplicity  $m_\sigma$ ,  $-\sigma$  and  $\sigma$  are both eigenvalues of  $\tilde{A}_\xi$ , both with multiplicities  $m_\sigma$ .

**Proof:** Let  $Q \geq 0$  and let  $L \in \mathbb{R}^{n \times n}$  be such that  $L^T L = Q$ . The transfer function that corresponds to the operator  $\lambda \mapsto w$  in (4) is equal to  $H_{w\lambda}(s) = \tilde{C}(sI_{2n} - \tilde{A})^{-1}\tilde{B} = G(s)G(-s)(1 - H^*(s)H(s))^{-1}$  [52], where  $G(s)$  is the transfer of the triple  $(A, B, C)$  in (2).  $H(s)$  is the transfer matrix of the triple  $(A, B, L)$ ,  $H^*(s)$  is the transfer matrix of the triple  $(-A^T, L^T, B^T)$ , i.e. the adjoint system to  $(A, B, L)$  [54, p.280]. Consider now the two ZD canonical forms associated to the systems in (2) and (4), respectively. From Lemma 1 one has  $\xi \in \mathbb{R}^{n-r^{wu}}$  and  $\tilde{\xi} \in \mathbb{R}^{2(n-r^{wu})}$ . From the fact that  $H_{w\lambda}(s) = G(s)G(-s)(1 - H^*(s)H(s))^{-1}$  and since the  $\xi_1$ -dynamics corresponds to the numerator  $B(s)$  of the transfer function  $G(s)$ , we deduce that the roots of the polynomial  $B(\alpha)B(-\alpha) = 0$  are modes of the matrix  $\tilde{A}_\xi$ . Since the order of  $B(\alpha)B(-\alpha)$  is precisely  $2(n - r^{wu})$ , the result follows. ■

<sup>1</sup>Clearly a similar transformation can be applied to  $(A, B, C)$ , and we then call the obtained representation the  $z$ -dynamics.

We therefore have a complete description of the triple  $(\tilde{A}, \tilde{B}, \tilde{C})$  in (4) in terms of its relative degree and zero dynamics.

### 3 The higher order Moreau's sweeping process

#### 3.1 Presentation of the differential inclusion

Let us consider the system in (4) (a) (c) and let us forget for the time being that it may represent the necessary conditions of the maximum principle. An important point is first to understand the dynamics of such a dynamical system involving complementary-slackness conditions. For instance, what is the meaning of  $\lambda \geq 0$  if  $\lambda$  is not a measure (distributional multipliers may easily be needed to integrate (4) (a) (c))? To this end let us recall some facts about the higher order sweeping process (HOSP) as it is introduced in [1, 2], which provides a rigorous framework to study the dynamics in (4) (a) (c), and extends the well-known first and second order sweeping process (see references in [1] and [11, §5.3] for a non-mathematical introduction). Roughly speaking, the HOSP is a specific differential inclusion, whose solutions are distributions, and which permits to give a meaning to the system (4) (especially the positiveness of  $\lambda$  when it is not a measure). We denote as  $\mathcal{T}_n(I)$  the set of distributions of degree  $n + 1$  [22] which are generated by RCSLBV functions on  $I$ , whose successive derivatives possess an absolutely continuous part (denoted as  $[\cdot]$ ) that is also RCSLBV on  $I$ . The right derivative of  $[h]$  is denoted as  $\hat{h}^{(1)} = \frac{d^+ [h]}{dt}(t) = \lim_{\sigma \rightarrow 0^+} \frac{[h](t+\sigma) - [h](t)}{\sigma}$ . The set of such functions is denoted as  $\mathcal{F}_\infty(I; \mathbb{R}) = \bigcap_{k \in \mathbb{N}} \mathcal{F}_k(I; \mathbb{R})$ , with  $\mathcal{F}_k(I; \mathbb{R}) = \{h \in \mathcal{F}_{k-1}(I; \mathbb{R}) : \hat{h}^{(k)} := \frac{d^+}{dt} [\hat{h}^{(k-1)}] \in RCSLBV(I; \mathbb{R})\}$ . In particular  $\mathcal{F}_0(I; \mathbb{R}) = RCSLBV(I; \mathbb{R})$ . If  $T \in \mathcal{T}_n(I)$  and is generated by a function  $F \in \mathcal{F}_\infty(I; \mathbb{R})$ , it has a “function” part denoted as  $\{T\}(\cdot) = [\hat{F}^{(n)}](\cdot)$ , and a “measure” part denoted as  $\ll T \gg$  such that  $\langle \ll T \gg, \varphi \rangle = \int_{-\infty}^{+\infty} \varphi d[\hat{F}^{(n-1)}]$ ,  $\forall \varphi \in C_0^\infty(I)$ .  $D$  denotes the distributional derivative, and  $dz$  denotes the Stieltjes or differential measure generated by a function  $z$  of local bounded variation [32]. Thus  $\mathcal{T}_n(I)$  denotes the set of all Schwartz' distributions such that there exists a function  $F \in \mathcal{F}_\infty(I; \mathbb{R})$  such that  $T = D^n F$ . Let  $n$  be the smallest integer such that  $T \in \mathcal{T}_n(I)$ , we set

$$\deg(T) = \begin{cases} n + 1 & \text{if } n \geq 1 \\ 1 & \text{if } n = 0 \text{ and } E_0(\{T\}) \neq \emptyset \\ 0 & \text{if } n = 0 \text{ and } E_0(\{T\}) = \emptyset \end{cases} \quad (15)$$

The Dirac measure  $\delta_0$  has degree 2, its derivative  $\dot{\delta}_0$  has degree 3, etc. Continuous functions have degree 0, and discontinuous functions have degree 1. The fact that the distributions we work with originate from LBV functions, is crucial for the characterization of their support. Since the notion of a solution for the formalism in which we shall embed (4) is crucial, an example is provided in detail in appendix B. We insist here on the fact that the solutions that will be considered next, are of special bounded variation. In another words, their derivatives do not contain any singular Lebesgue integrable part. See Definition 1 and equation (32).

The core of the HOSP lies in the definition of the following set of tangent cones. Let  $\Phi$  be a nonempty closed convex subset of  $\mathbb{R}$ . We denote by  $T_\Phi(x)$  the tangent cone of  $\Phi$  at  $x \in \mathbb{R}$  defined by

$$T_\Phi(x) = \overline{\text{cone}(\Phi - \{x\})} \quad (16)$$

where  $\text{cone}(\Phi - \{x\})$  denotes the cone generated by  $\Phi - \{x\}$ , defined as in [37] to take into account constraint violations, and  $\overline{\text{cone}(\Phi - \{x\})}$  is its closure. Given a closed nonempty convex set  $\Phi$ , we set

$$T_\Phi^0(z_1) = \Phi, \quad T_\Phi^1(z_1) = T_\Phi(z_1), \quad T_\Phi^2(z_1, z_2) = T_{T_\Phi^1(z_1)}(z_2),$$

$$T_\Phi^i(z_1, \dots, z_i) = T_{T_\Phi^{i-1}(z_1, \dots, z_{i-1})}(z_i), \quad i \geq 0.$$

**Remark 3** If  $m \geq 2$  and  $\Phi = (\mathbb{R}^+)^m$  then

$$T_{\Phi}^i(z_1, \dots, z_i) = \times_{l=1}^m T_{\Phi}^i(z_1^l, \dots, z_i^l).$$

Note that

$$T_{\mathbb{R}^+}(x) = \begin{cases} \mathbb{R} & \text{if } x > 0 \\ \mathbb{R}^+ & \text{if } x \leq 0 \end{cases}$$

and

$$T_{\mathbb{R}}(x) = \mathbb{R}.$$

**Remark 4** i) The subdifferential of the indicatrix function of the cone  $T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})$  is given by [23, §1.3.1]

$$\partial\psi_{T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})}(z_i) = \{w \in \mathbb{R}^n : \langle w, v - z_i \rangle \geq 0, \forall v \in T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})\}$$

and is the outward normal cone to  $T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})$  at  $z_i$ .

ii) It can be shown that

$$T_{\Phi}^{i-1}(z_1, \dots, z_{i-1}) = \mathbb{R} \Rightarrow \partial\psi_{T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})}(z_i) = \{0\},$$

$$T_{\Phi}^{i-1}(z_1, \dots, z_{i-1}) = \mathbb{R}^+ \text{ and } z_i > 0 \Rightarrow \partial\psi_{T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})}(z_i) = \{0\},$$

$$T_{\Phi}^{i-1}(z_1, \dots, z_{i-1}) = \mathbb{R}^+ \text{ and } z_i \leq 0 \Rightarrow \partial\psi_{T_{\Phi}^{i-1}(z_1, \dots, z_{i-1})}(z_i) = \mathbb{R}^-.$$

Let us introduce now two mathematical formalisms of the HOSP. For this we rely upon the special state space representation in (14), the reason for this particular choice being explained in [1].

**Distributional Formalism.** Find  $z_1, \dots, z_r \in \mathcal{T}_{r-1}(\mathbb{R}^+)$  and  $\xi_i \in \mathcal{T}_{r-1}(\mathbb{R}^+)$  ( $1 \leq i \leq n-r$ ) satisfying the distributional equations

$$\begin{cases} Dz_1 - z_2 = 0 \\ Dz_2 - z_3 = 0 \\ Dz_3 - z_4 = 0 \\ \vdots \\ Dz_{r-1} - z_r = 0 \\ Dz_r - \tilde{C}\tilde{A}^r\tilde{W}^{-1}\tilde{z} = \tilde{C}\tilde{A}^{r-1}\tilde{B}\lambda \\ D\xi = \tilde{A}_{\xi}\xi + \tilde{B}_{\xi}z_1 \end{cases} \quad (17)$$

$$\lambda = (\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}[D^{(r-1)} \ll Dz_1 - \{z_2\} \gg + \dots +$$

$$+ D \ll Dz_{r-1} - \{z_r\} \gg] + \ll Dz_r - \tilde{C}\tilde{A}^r\tilde{W}^{-1}\{z\} \gg \quad (18)$$

and

$$\left\{ \begin{array}{l} d\{z_1\} - \{z_2\}(t)dt \in -\partial\psi_{T_\Phi^0}(\{z_1\}(t^+)) \\ \vdots \\ d\{z_i\} - \{z_{i+1}\}(t)dt \in -\partial\psi_{T_\Phi^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)) \\ \vdots \\ d\{z_{r-1}\} - \{z_r\}(t)dt \in -\partial\psi_{T_\Phi^{r-2}(\{z_1\}(t^-), \dots, \{z_{r-2}\}(t^-))}(\{z_{r-1}\}(t^+)) \\ (\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}[d\{z_r\} - \tilde{C}\tilde{A}^r\tilde{W}^{-1}\{z\}(t)dt] \in -\partial\psi_{T_\Phi^{r-1}(\{z_1\}(t^-), \dots, \{z_{r-1}\}(t^-))}(\{z_r\}(t^+)) \end{array} \right. \quad (19)$$

More compactly (17) is rewritten as  $D\tilde{z} - \tilde{W}\tilde{A}\tilde{W}^{-1}\tilde{z} - \tilde{W}\tilde{B}\lambda = 0$  or  $D\tilde{x} - \tilde{A}\tilde{x} - \tilde{B}\lambda = 0$ . The relations given in (19) have to be interpreted in the following sense: Find nonnegative real-valued Radon measures  $d\mu_i$  ( $1 \leq i \leq r$ ) relative to which the Lebesgue measure  $dt$  and the Stieltjes measure  $d\{z_i\}$  possess densities  $\frac{dt}{d\mu_i}$  and  $\frac{d\{z_i\}}{d\mu_i}$  respectively such that:

$$\begin{aligned} \frac{d\{z_i\}}{d\mu_i}(t) - \{z_{i+1}\}(t) \frac{dt}{d\mu_i}(t) &\in -\partial\psi_{T_\Phi^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)), \\ d\mu_i - \text{a.e. } t \in \mathbb{R} \quad (1 \leq i \leq r-1) & \end{aligned} \quad (20)$$

and

$$\begin{aligned} (\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}[\frac{d\{z_r\}}{d\mu_r}(t) - \tilde{C}\tilde{A}^r\tilde{W}^{-1}\{z\}(t) \frac{dt}{d\mu_r}(t)] \\ \in -\partial\psi_{T_\Phi^{r-1}(\{z_1\}(t^-), \dots, \{z_{r-1}\}(t^-))}(\{z_r\}(t^+)), \quad d\mu_r - \text{a.e. } t \in \mathbb{R}. \end{aligned} \quad (21)$$

The solutions of the distributional formalism will be shown to be distributions of degree possibly larger than 1. Let us now introduce a second formalism whose solutions are functions, independently of the relative degree  $r^{wu}$ . This will be quite useful when the optimal control problem is embedded into the HOSP, see sections 3.5, 3.6.

**Measure Differential Formalism.** Find  $z_i \in \mathcal{F}_\infty(\mathbb{R}^+; \mathbb{R})$  ( $1 \leq i \leq r$ ) and  $\xi_i \in \mathcal{F}_\infty(\mathbb{R}^+; \mathbb{R})$  ( $1 \leq i \leq 2n - r$ ) such that <sup>(2)</sup>

$$dz_i - z_{i+1}(t)dt \in -\partial\psi_{T_\Phi^{i-1}(z_1(t^-), \dots, z_{i-1}(t^-))}(z_i(t^+)) \quad (1 \leq i \leq r-1) \quad (22)$$

$$(\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}[dz_r - \tilde{C}\tilde{A}^r\tilde{W}^{-1}z(t)dt] \in -\partial\psi_{T_\Phi^{r-1}(z_1(t^-), \dots, z_{r-1}(t^-))}(z_r(t^+)) \quad (23)$$

and

$$\dot{\xi}(t) - \tilde{A}_\xi\xi(t) - \tilde{B}_\xi z_1(t) = 0, \quad dt - \text{a.e. } t \in \mathbb{R} \quad (24)$$

The system in (22) and (23) has to be interpreted in the following sense: Find nonnegative real-valued Radon measure  $d\mu_i$  relative to which the Lebesgue measure  $dt$  and the Stieltjes measure  $dz_i$  possess densities  $\frac{dt}{d\mu_i}$  and  $\frac{dz_i}{d\mu_i}$  respectively such that

$$\frac{dz_i}{d\mu_i}(t) - z_{i+1}(t) \frac{dt}{d\mu_i}(t) \in -\partial\psi_{T_\Phi^{i-1}(z_1(t^-), \dots, z_{i-1}(t^-))}(z_i(t^+)), \quad d\mu_i - \text{a.e. } t \in \mathbb{R} \quad (1 \leq i \leq r-1) \quad (25)$$

and

$$(\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}[\frac{dz_r}{d\mu_r}(t) - \tilde{C}\tilde{A}^r\tilde{W}^{-1}z(t) \frac{dt}{d\mu_r}(t)]$$

<sup>2</sup>In this formalism,  $\{z_i\} = z_i$  because  $z_i$  is a function and  $dz_i$  is a Stieltjes measure.

$$\in -\partial\psi_{T_{\Phi}^{r-1}(z_1(t^-), \dots, z_{r-1}(t^-))}(z_r(t^+)), \quad d\mu_r - \text{a.e. } t \in \mathbb{R}. \quad (26)$$

Another, more intuitive way to see the measure differential formalism, is to consider the following: Find  $z_1, \dots, z_r, \xi_1, \dots, \xi_{n-r} \in \mathcal{F}_{\infty}(\mathbb{R}^+; \mathbb{R})$  and measures  $d\nu_1, \dots, d\nu_r$  such that

$$\left\{ \begin{array}{l} dz_1 = z_2(t)dt + d\nu_1 \\ dz_2 = z_3(t)dt + d\nu_2 \\ dz_3 = z_4(t)dt + d\nu_3 \\ \vdots \\ dz_i = z_{i+1}(t)dt + d\nu_i \\ \vdots \\ dz_{r-1} = z_r(t)dt + d\nu_{r-1} \\ dz_r = \tilde{C}\tilde{A}^r\tilde{W}^{-1}z(t)dt + \tilde{C}\tilde{A}^{r-1}\tilde{B}d\nu_r \\ \dot{\xi}(t) = \tilde{A}_{\xi}\xi(t) + \tilde{B}_{\xi}z_1(t) \end{array} \right. \quad (27)$$

with the inclusions

$$d\nu_i \in -\partial\psi_{T_{\Phi}^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)) \quad \text{for all } 1 \leq i \leq r \quad (28)$$

and  $\lambda$  is in (18):

$$\lambda = (\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}[D^{(r-1)}\nu_1 + \dots + D\nu_{r-1}] + \nu_r \quad (29)$$

In other words  $\lambda$  is a distribution whose measure part is  $d\nu_r$  and its positivity is understood as the positivity of  $d\nu_r$ . The measures  $d\nu_i$  are multipliers whose meaning in the context of optimality will be made clear later (see Corollaries 8 and 9, and remark 8). From (27) and (28) it follows that

$$d\nu_i(\{t\}) = \{z_i\}(t^+) - \{z_i\}(t^-) \in -\partial\psi_{T_{\Phi}^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)), \quad (1 \leq i \leq r-1) \quad (30)$$

and

$$d\nu_r(\{t\}) = \{z_r\}(t^+) - \{z_r\}(t^-) \in -\tilde{C}\tilde{A}^{r-1}\tilde{B} \partial\psi_{T_{\Phi}^{r-1}(\{z_1\}(t^-), \dots, \{z_{r-1}\}(t^-))}(\{z_r\}(t^+)). \quad (31)$$

We recall that given a closed convex nonempty set  $K$  and  $M = M^T > 0$ ,  $\text{prox}_M[K; \tilde{z}]$  is the closest point to  $\tilde{z}$  in  $K$ , in the metric defined by  $M$ . To better understand (30) and (31) it is useful to recall here the equivalences for two vectors of appropriate dimension,  $M = M^T > 0$  and  $K$  a closed convex set:  $x - y \in -M^{-1}\partial\psi_K(x) \iff x = \text{prox}_M[K, y] \iff x = \text{proj}_M(K; y) \iff x = \text{argmin}_{z \in K} \frac{1}{2}(z - y)^T M (z - y) \iff \langle x - y, v - x \rangle \geq 0$  for all  $v \in K$ . As we shall see, the measure differential formalism is useful in the context of optimal control, see section 3.5. Indeed it is a formalism which retains only the ‘‘measure part’’ of the distributional formalism, i.e. the states discontinuities. Most importantly its solutions are functions of time independently of the degree of  $\lambda$ . Both formalisms are related through the following.

**Proposition 1** [1] *i) Let  $(z_1, \dots, z_r, \xi) \in (\mathcal{T}_{r-1}(\mathbb{R}^+))^n$  be a solution of Problem (17) (18) (19), with  $\{z\}(0^-) = z_0 \in \mathbb{R}^n$ . Then  $z_1 = \{z_1\} \in \mathcal{F}_{\infty}(\mathbb{R}^+; \mathbb{R})$ ,  $z_i \in \mathcal{T}_{i-1}(\mathbb{R}^+)$  ( $2 \leq i \leq r$ ),  $\xi = \{\xi\} \in (\mathcal{F}_{\infty}(\mathbb{R}^+; \mathbb{R}))^{n-r}$  and  $(\{z_1\}, \dots, \{z_r\}, \xi)$  is a solution of Problem (22) (23) (24), with  $\{z\}(0^-) = z_0 \in \mathbb{R}^n$ .*

*ii) Let  $(w_1, \dots, w_r, \xi) \in (\mathcal{F}_{\infty}(\mathbb{R}^+; \mathbb{R}))^n$  be a solution of Problem (22) (23) (24), with  $\{z\}(0^-) = z_0 \in \mathbb{R}^n$ . Then  $(z_1, \dots, z_r, \xi) \in (\mathcal{T}_{r-1}(\mathbb{R}^+))^n$ , where  $z_1 := w_1$  and*

$$z_i = w_i + \sum_{j=1}^{i-1} \left( \sum_{t_k \in E_0(w_j)} (w_j(t_k^+) - w_i(t_k^-)) \delta_{t_k}^{(i-j-1)} \right) \quad (2 \leq i \leq r),$$

is a solution of Problem (17) (18) (19), with  $\{z\}(0^-) = z_0 \in \mathbb{R}^n$ . ■

An example is treated in detail in section 4.4.2 which helps understanding how (17)-(19) is integrated in time and how the state jump mappings (31) (30) work. Notice that Proposition 1 applies to the solutions of the BVP as well. It is noteworthy that if the distribution  $\tilde{z}$  is a solution of the distributional problem in (17) (18) (19) and if the function  $\tilde{\zeta}(\cdot)$  is a solution of the measure problem (22) (23) (24), then  $\{\tilde{z}\}(\cdot) = \tilde{\zeta}(\cdot)$  almost everywhere on  $[0, T_1]$ . It is also clear that if  $\lambda$  is a measure, then both formalisms are equivalent one to each other. From a Control engineer point of view, the “real” solution is the solution of the distributional formalism. However the measure differential formalism will prove to be quite useful to formulate an extended optimal control problem.

**Definition 1** Let  $0 \leq a < b \in \mathbb{R} \cup \{+\infty\}$  be given. We say that a solution  $z \in (\mathcal{T}_{r-1}(\mathbb{R}^+))^n$  of (17)(18)(19), with  $\{z\}(0^-) = z_0 \in \mathbb{R}^n$ , is regular on  $[a, b)$  if for each  $t \in [a, b)$ , there exists a right neighborhood  $[t, \sigma)$  ( $\sigma > 0$ ) such that the restriction of  $\{z\}$  to  $[t, \sigma)$  is analytic.

Regular solutions may possess accumulations of jumps, but only on the left of any  $t \in [a, b)$ , as right accumulations cannot exist by definition. They encompass right-analytic solutions [46]. The following is proved in [1] and concerns the IVP in (17) (18) (19).

**Theorem 1** Suppose that  $\tilde{C}\tilde{A}^{r-1}\tilde{B} > 0$  and  $m = 1$ . For each  $z_0 \in \mathbb{R}^n$ , the system in (17) (18) (19), with  $\{z\}(0^-) = z_0 \in \mathbb{R}^n$  has at least one regular solution.

Moreover:

i)  $z_1(\cdot) \equiv \{z_1\}(\cdot) \geq 0$  on  $\mathbb{R}^+$

ii)  $\{\tilde{z}\}(0^+) = \tilde{z}'_0$

iii)  $\|\{\tilde{z}\}(t)\| \leq \sqrt{e^{ct}}\|\tilde{z}_0\|$ ,  $\forall t \in \mathbb{R}^+$ , for some  $c > 0$ ,

iv) If  $\tilde{z}^1$  and  $\tilde{z}^2$  are two regular solutions of (17) (18) (19), with  $\{\tilde{z}\}(0^-) = \tilde{z}_0 \in \mathbb{R}^n$  then  $\langle \tilde{z}^1, \varphi \rangle = \langle \tilde{z}^2, \varphi \rangle$ ,  $\forall \varphi \in C_0^\infty(\mathbb{R}^+; \mathbb{R}^n)$ ; where  $\tilde{z}'_0$  is uniquely defined by

$$z'_{0,i} = \text{prox} [T_\Phi^{i-1}(z_{0,1}, \dots, z_{0,i-1}); z_{0,i}], \forall 1 \leq i \leq r-1,$$

$$z'_{0,r} = \text{prox}_{(\tilde{C}\tilde{A}^{r-1}\tilde{B})^{-1}} [T_\Phi^{r-1}(z_{0,1}, \dots, z_{0,r-1}); z_{0,r}]$$

and

$$z'_{0,l} = z_{0,l}, \quad (r+1 \leq l \leq 2n).$$

It is proved that uniqueness holds in the class of regular solutions only. A crucial step is to understand the relationship between  $\mathbf{U}$  and the class of solutions of the HOSP. Indeed it is possible that reachability of  $\bar{x}_1$  from  $\bar{x}_0$  in (2), holds for a  $u$  which belongs to a set  $\mathbf{U}$  of distributions. However is there a correspondance between this  $\mathbf{U}$  and  $\mathcal{T}_{r-1}([0, T_1])$ ? An element of answer will be given in section 3.5. An important property which helps understanding how the HOSP works is the following. Notice first that we may write  $d\nu_i$  in (27) (28) as

$$d\nu_i = g_i(t)dt + d\mathcal{J}_i, \tag{32}$$

where  $g_i \in \mathcal{F}_\infty(\mathbb{R}^+; \mathbb{R})$  and  $d\mathcal{J}_i$  is an atomic measure with countable and orderable set of atoms generated by the right continuous jump function  $\mathcal{J}_i$ . The following holds.

**Proposition 2** [1] Let  $z(\cdot)$  be a solution of (20)-(21). We have

$$g_i(t) = 0, \quad dt - \text{a.e. } t \in \mathbb{R}^+, \quad (1 \leq i \leq r-1), \tag{33}$$

$$g_r(t) \in -\partial\psi_{T_{\Phi}^{r-1}(\{z_1\}(t^-), \dots, \{z_{r-1}\}(t^-))}(\{z_r\}(t^+)), \quad dt - \text{a.e. } t \in \mathbb{R}^+, \quad (34)$$

and

$$0 \leq z_1(t^+) \perp d\nu_r(t^+) \geq 0, \quad \text{for each } t \in \mathbb{R}^+ \quad (35)$$

$$0 \leq z_1(t^+) \perp g_r(t^+) \geq 0, \quad dt - \text{a.e. } t \in \mathbb{R}^+. \quad (36)$$

$$T_{\Phi}^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-)) \ni \{z_i\}(t^+) \perp -d\mathcal{J}_i(\{t\}) \in \partial\psi_{T_{\Phi}^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)) \quad (37)$$

for all  $1 \leq i \leq r$  and all  $t \in \mathbb{R}^+$ .

Thus, the complementarity in (4) (c) is given a meaning with (36) and (37). The cone CP in (37) shows what CP the measures  $d\nu_i$  satisfy (the rigorous meaning of the inclusion in the right-hand-side of (37) is explained in [1, Proposition 1] or [32, p.76]). It is important to notice that given  $\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-)$ , the measure value depends on  $\{z_i\}(t^+)$ . We shall come back on (37) in section 4.2. From (33) it follows that the non-atomic parts of the Lagrange multipliers are zero almost everywhere except for  $g_r(\cdot)$ . Consequently the Lagrange multipliers associated to additional constraints to be imposed on the problem, are purely atomic. This may explain why they have been ignored in the literature since the condition that  $\lambda$  be a measure has always been imposed.

**Remark 5** *i) Theorem 1 remains true if  $m \geq 2$  and  $\tilde{C}\tilde{A}^{r-1}\tilde{B}$  is a nonsingular symmetric  $\mathbf{M}$ -matrix [18].*

*ii) The HOSP is implementable with the knowledge of  $A, B, C, D$ , and  $x(\cdot)$ . The canonical  $\tilde{z}$ -dynamics in (14) is introduced only for the sake of clarity of the presentation of the differential inclusion.*

*iii) The sets  $\text{supp}(\lambda)$  and  $\text{supp}(d\nu_i)$  are countable and orderable sets.*

*iv) If  $C\bar{x}_0 + D \geq 0$ , then along solutions of the HOSP IVP one has  $\text{deg}(w(\cdot)) = 0$  on  $[0, T_1]$ .*

Another important tool is the following LCP which holds on time intervals  $[\tau - \epsilon, \tau)$ ,  $\epsilon > 0$ , on which  $z_1(t)$  is identically zero:

$$0 \leq \lambda(t) \perp \tilde{C}\tilde{A}^r\tilde{W}^{-1}\tilde{z}(t^+) + \tilde{C}\tilde{A}^{r-1}\tilde{B}\lambda(t) \geq 0 \quad (38)$$

In (38) one could write equivalently  $g_r(t)$  instead of  $\lambda(t)$ . This LCP monitors the evolution of the multiplier  $\lambda(\cdot)$  for  $t \in [\tau - \epsilon, \tau)$ . The matrix of this LCP is the leading Markov parameter  $M^{(2r^{w_u})}$ . We therefore have at our disposal a complete dynamical system which allows us to give a meaning to the system in (4), for any initial data and for all  $t \geq 0$ . This is a first crucial step for the understanding of the necessary conditions of optimality.

**Corollary 1** *Along trajectories of (17) (18) (19) one has  $\frac{d}{dt}H(\{x\}(t), \{\eta\}(t)) = 0$  almost everywhere.*

**Proof:** follows from easy calculations and the fact that the support of the measures  $d\nu_i$  is a zero Lebesgue measure on  $[0, T_1]$ . ■

However along the optimal trajectory the Hamiltonian undergoes jumps at atoms of  $\lambda$  (see Corollary 3).



### 3.2 Solutions of the necessary conditions IVP

From now on we will choose to embed the system in (4) (a) (c) in the HOSP. In particular, the degree of distributions refers to the degree of distributions in spaces  $\mathcal{T}_n([0, T_1])$  for some  $n$ . The following is going to be useful subsequently and follows from the state variable change  $\tilde{z} = \tilde{W}\tilde{x}$  and (17)–(19) and (22)–(24).

*Costate equation distributional formalism*

$$D\eta - Qx + A^T\eta = [0_{n \times n} \quad I_n]\tilde{W}^{-1} \begin{pmatrix} 0^{r-1} \\ M^{(r)}\lambda \\ 0^{2n-r-2} \end{pmatrix} \quad (39)$$

*Costate equation measure formalism*

$$d\eta - Qx(t)dt + A^T\eta(t)dt = [0_{n \times n} \quad I_n]\tilde{W}^{-1} \begin{pmatrix} d\nu_1 \\ \vdots \\ d\nu_{r-1} \\ M^{(r)}d\nu_r \end{pmatrix} \quad (40)$$

Both  $x(\cdot)$  and  $\eta(\cdot)$  are  $\mathcal{F}_\infty([0, T_1], \mathbb{R})$ -functions in (40), while in (39)  $x$  and  $\eta$  are distributions in  $\mathcal{T}_{r-1}([0, T_1])$ . We also denote  $d\nu_\eta = [0_{n \times n} \quad I_n]\tilde{W}^{-1}(d\nu_1 \dots d\nu_{r-1} M^{(r)}d\nu_r)^T$ . In a first step, it is necessary to give a meaning to (4) (a) (c), as a well-posed dynamical system.

**Lemma 3** *Let us embed the system in (4) (a) (c) in the framework of the HOSP. Assume that  $r^{wu} = 2k$ ,  $k \geq 1$  and  $m = 1$ . Then there exists a global solution starting from any initial data  $(x(0^-), \eta(0^-))$  which is regular in the sense of Definition 1, and uniqueness holds in the set of regular solutions.*

**Proof :** From Lemma 1, the leading Markov parameter  $M^{(2r^{wu})} = \tilde{C}\tilde{A}^{2r^{wu}-1}\tilde{B} > 0$ . Apply Theorem 1 and the fact that the solution  $\tilde{x}$  of (4) (a) (c) is regular on  $\mathbb{R}^+$  for any initial data. ■

The condition  $\tilde{z}(\tau^-) = 0$  at an entry time  $\tau$  is sufficient to obtain a bounded  $u(\tau)$  that keeps the state inside  $\Phi$ . Thus  $\lambda$  is an analytic function of time and from (4) (a) and (5) we deduce that  $u(\cdot)$  is also an analytic function of time. If  $\tilde{z}(\tau^-) < 0$  then from (30) (31) a state jump occurs and  $\deg(\lambda) \geq 2$ . However the condition  $\tilde{z}(\tau^-) = 0$  is not necessary to get a bounded  $u(\cdot)$ : if  $z_1(\tau^-) = \dots = z_{r-1}(\tau^-) = 0$  and  $z_r(\tau^-) < 0$ ,  $\deg(\lambda) = 2$  and  $u(\tau)$  may be a (discontinuous) function of time, provided  $\eta(\cdot)$  jumps at  $\tau$  while  $x(\cdot)$  remains continuous at  $\tau$ , see Proposition 6 below.

It is crucial to realise that Lemma 3 neither states uniqueness of solutions of the BVP (4) (a) (b) (c) when  $r^{wu}$  is even, nor says that when  $r^{wu}$  is odd, the BVP (4) (a) (b) (c) has no solution. The well-posedness of the IVP (4) (a) (c) and of the BVP (4) (a) (b) (c), are two very distinct notions (as said above, it could even be that the BVP is well-posed when elastic jumps are introduced). However studying a good formalism for the IVP, will allow us to better understand the BVP. In particular the study of local properties of the optimal solution can take advantage of the IVP study, as well as for numerical integration. As pointed out in the introduction, the existence of solutions to the minimization problem (1) (2), is a reachability problem. Concerning the BVP in (4) (a) (b) (c), we make the choice to embed it into the HOSP, i.e. we allow for its solutions to be regular in the sense of Definition 1. The consequences of such a choice are important, as it may imply to extend the integral action  $I(u)$  in (1) since  $x$ ,  $\eta$  and  $u$  may no longer be functions at contact times. Actually the central question is: let us allow the solutions of the BVP in (4) to be regular. Then how should one extend the minimization problem in (1) (2) so that the set of equations/conditions in (17) (18) (19) represent the necessary optimality conditions?

### 3.3 Motivation example

The first question that comes to one's mind is: apart from getting a better understanding of the dynamical behaviour of the system in (4) (a) (c), what do we gain by embedding the necessary conditions system in

(4) into the HOSP? In order to get a preliminary answer, let us examine the simple case where the system (2) is given in the  $z$ -dynamics representation,  $Q = I_n$  and  $R = I_{n_u}$ . Let us rewrite  $I(u)$  as the sum of a regular part  $\frac{1}{2} \int_0^{T_1} [\{z\}(t)^T \{z\}(t) + \frac{1}{2} \{u\}(t)^T R \{u\}(t)] dt$  and an algebraic part  $\frac{1}{2} \sum_k \sigma_{\{z\}}^T(t_k) \sigma_{\{z\}}(t_k)$  which takes into account the norm of the state vector jump. We notice that in addition to the minimisation of the algebraic part has to take into account some constraints on  $\{z\}(t^+)$ , due to the unilateral constraints. In other words, we disregard what happens at instants  $t_k$  and retain only the right and left-limits of the function part of the state: we choose to work with the measure differential formalism (22) (24) or (25) (26). Assume now that  $\{z_1\}(t) = 0$  and  $\{z_i\}(t) < 0$ , for some  $t \in [0, T_1]$  and all  $2 \leq i \leq r$ . Thus from (30) (31) a jump in  $\{z_i\}(\cdot)$ ,  $2 \leq i \leq r$ , occurs at  $t$  while the other variables are continuous. Then we have  $\{z_i\}(t^+) - \{z_i\}(t^-) = \operatorname{argmin}_{\sigma_{\{z_i\}}(t^-) \geq 0} \frac{1}{2} \sigma_i^2$ . The HOSP jump rule therefore minimizes  $\frac{1}{2} \sum_k \sum_i (\sigma_{\{z_i\}}(t_k))^2$ , over the admissible values of the post-impact state.

### 3.4 Meaning of the costate $\eta(\cdot)$ jump condition

In this subsection we rely on the embedding of the system (4) (a)-(c) in the HOSP formalism to characterize both  $x$  and  $\eta$  as distributions. The condition  $\eta(\tau^+) = \eta(\tau^-) - C^T \lambda_1$  with  $\lambda_1 \geq 0$  is usually given in the set of the necessary conditions. It is also sometimes indicated that if the optimal controller  $u$  is discontinuous, then a jump in  $\eta(\cdot)$  occurs. In the light of the HOSP distributional and measure differential formalisms in (17)-(21) and (22)-(26) respectively, one may wonder what this jump condition really means <sup>(3)</sup>: is  $\lambda$  a measure at  $\tau$  (i.e.,  $\lambda_1$  is the magnitude of the atom of  $d\mathcal{J}_\tau$  at  $\tau$ , see (32)), or is it  $\{\eta\}(\tau^+) = \{\eta\}(\tau^-) - C^T \lambda_1$ ? This is not at all equivalent.

**Proposition 3** *The degree of  $\lambda$  at a time  $t$  is  $\leq 2$  if and only if all the  $z_i(\cdot)$ ,  $1 \leq i \leq r-1$ , are continuous at  $t$ , while  $z_r(\cdot)$  is a possibly discontinuous function at  $t$ . Moreover in such a case  $x(\cdot) = \{x\}(\cdot)$  is a continuous function at  $t$  while  $\eta(\cdot) = \{\eta\}(\cdot)$  jumps at  $t$ .*

**Proof:** The first assertion comes from the  $\tilde{z}$ -dynamics. Now if  $x(\cdot)$  is a function and has a jump, necessarily from (4)  $\deg(\eta) \geq 2$ . Thus necessarily  $\deg(\lambda) \geq 3$  which is a contradiction. Finally  $\tilde{z} = \tilde{W}\tilde{x}$  and the transformation matrix  $\tilde{W}$  is square full rank. Thus if  $\eta(\cdot)$  is continuous at  $t$ , so is  $\tilde{z}(\cdot)$ . So necessarily  $\eta(\cdot)$  has to be discontinuous if  $z_r(\cdot)$  is. The second part of the Proposition is proved. ■

Proposition 3 says nothing on the definition of the jump: clearly if  $M^{(r)} \leq 0$  and if we admit that the system in (4) is embedded in the HOSP, then the post-jump state is not well defined. Finally the jump in  $\eta(\cdot) = \{\eta\}(\cdot)$  has to satisfy stringent conditions to assure that only  $z_r(\cdot)$  jumps while the lower order variables  $z_i(\cdot)$  remain continuous.

**Proposition 4** *Let  $m = n_u = 1$ . All the  $z_i(\cdot)$ ,  $1 \leq i \leq r-1$ , are continuous at  $t$ , while  $z_r(\cdot)$  is a discontinuous function at  $t$ , if and only if  $(A^T)^j \sigma_\eta(t) \in \operatorname{Ker}(B^T)$  for  $0 \leq j \leq r^{wu} - 1$ .*

**Proof:** Let us consider (4) and the state transformation in (14). It is easy to see that  $\dot{z}_i(t) = z_{i+1}(t) = CA^i x(t)$  for all  $1 \leq i \leq r^{wu} - 1$  when  $r^{wu} \geq 2$ . Let us assume that  $\dot{z}_i(t) = M_i x(t) +$

$\sum_{j=0}^{i-r^{wu}} \star B^T (A^T)^j \eta(t)$ , for some  $r^{wu} \leq i \leq r-1$ , where  $\star$  generically denotes some constant scalar and

$M_i \in \mathbb{R}^{1 \times n}$  is a row vector. Then  $\dot{z}_{i+1}(t) = M_i (Ax(t) + BB^T \eta(t)) + \sum_{j=0}^{i-r^{wu}} \star B^T (A^T)^j (Qx(t) - A^T \eta(t)) =$

$M_{i+1} x(t) + \sum_{j=0}^{i-r^{wu}+1} \star B^T (A^T)^j \eta(t)$ , where  $M_{i+1} = M_i A + \sum_{j=0}^{i-r^{wu}} \star B^T (A^T)^j Q$  (notice that  $M_i B \in \mathbb{R}$ ).

Since  $\dot{z}_{r^{wu}}(t) = CA^{r^{wu}} x(t) + CA^{r^{wu}-1} BB^T \eta(t)$  the proof is complete by induction with  $M^{r^{wu}} = CA^{r^{wu}}$

<sup>3</sup>We say that a vector jumps if at least one of its components jumps.

and  $\star = CA^{r^{wu}-1}B$ . From Proposition 3  $x(\cdot)$  is continuous. The “if” part of the proof thus follows by letting  $j$  vary from 0 to  $i - r^{wu}$  and  $r^{wu} \leq i \leq r - 1$ . The “only if” part follows from the fact that the pair  $(A, B)$  is controllable. Therefore the rows  $B^T(A^T)^j$  are independent for all  $0 \leq j \leq n - 1$ . Since  $M_i$  is an  $n$ -row and  $\star$  is a scalar, and since  $r^{wu} \leq n$ , the proof follows. ■

**Corollary 2** *Let  $\eta(\cdot)$  have a jump  $\sigma_\eta(t)$  at  $t$  and  $(A^T)^j \sigma_\eta(t) \in \text{Ker}(B^T)$  for  $0 \leq j \leq r^{wu} - 1$ . Then  $u^{(r^{wu}-1)}(\cdot)$  is the lowest order derivative of  $u(\cdot)$  which is discontinuous at  $t$ .*

**Proof:** From Proposition 4 it follows that  $\deg(\lambda) \leq 2$  so that (4) does represent the optimality conditions for (1) (2) and (5) holds [47]. Since  $u(\cdot) = B^T \eta(\cdot)$ , we deduce that  $\sigma_u(t) = B^T \sigma_\eta(t) = 0$ . Since  $\dot{u}(\cdot) = B^T \dot{\eta}(\cdot) = B^T(Qx(t) - A^T \eta(t))$  (4), we deduce that  $\sigma_{\dot{u}}(t) = -B^T A^T \sigma_\eta(t) = 0$ . The reasoning can be continued until one attains  $u^{(r^{wu}-1)}$ , recalling that  $z_1^{(r-1)}(\cdot) = z_r(\cdot)$  depends on  $u^{(r^{wu}-1)}(\cdot)$ . ■

One notices from Corollary 2 that the sum of the lowest derivative of  $u(\cdot)$  that is discontinuous and of the relative degree are always equal to  $2r^{wu} - 1$  whenever there is a jump in  $\eta(\cdot)$  at  $t$  (i.e.  $\lambda_1 > 0$ ) and  $\lambda$  is a measure at  $t$ . This is in agreement with [28, Theorem 6] which states that  $\lambda_1 \neq 0$  if and only if  $u^{(r^{wu}-1)}(\cdot)$  is discontinuous at  $t$ , at an entry time. However our proof completely differs from that in [28].

**Corollary 3** *Let  $\lambda$  be a measure and consider the Hamiltonian function in (9) evaluated along the solutions of the HOSP in (17)-(19) (equivalently along the solutions of (22)-(24)). Then  $H(t^+) = H(t^-)$  at  $t \in [0, T_1]$  if and only if  $(\eta(t^+) - \eta(t^-))^T Ax(t) = 0$ .*

**Proof:** One finds from Proposition 3 that  $H(t^+) - H(t^-) = -\frac{1}{2}(\eta(t^+) - \eta(t^-))^T BB^T(\eta(t^+) - \eta(t^-)) + (\eta(t^+) - \eta(t^-))^T Ax(t)$ . Thus from Proposition 4  $H(t^+) - H(t^-) = (\eta(t^+) - \eta(t^-))^T Ax(t)$ . ■

Therefore requiring the continuity of the Hamiltonian at contact states usually implies that  $\deg(\lambda) \leq 1$ . It is of interest to investigate whether or not conditions exist such that the Hamiltonian function, the optimal trajectory and the optimal control are functions of time (so that  $I(u)$  in (1) and (5) have a meaning), while  $\lambda$  is a distribution of degree  $\geq 3$ . In a sense, the costate  $\eta$  should “incorporate” all the higher degree distributions. Let us denote the transformation matrix  $\tilde{z} = \tilde{W}\tilde{x}$  as  $\tilde{W} = \begin{pmatrix} \tilde{W}_1 & 0_{n \times n} \\ \tilde{W}_2 & \tilde{W}_3 \end{pmatrix}$ , where  $\tilde{W}_j \in \mathbb{R}^{n \times n}$ ,  $1 \leq j \leq 3$ .

**Proposition 5** *Let  $m = n_u = 1$  and  $r^{wu} = n$ . Then (a)  $\deg(\lambda) \leq 2 \Leftrightarrow$  (b)  $\tilde{z}$  and  $\tilde{x}$  are functions  $\Leftrightarrow$  (c)  $B^T \eta$  and  $A^T \eta$  are functions. Also  $\deg(\lambda) \geq 3 \Leftrightarrow \deg(B^T \eta) \geq 2$  or  $\deg(A^T \eta) \geq 2$ .*

**Proof:** We first notice that under the stated conditions  $\tilde{W}$  indeed possesses the above structure, and since it is full rank then necessarily  $\tilde{W}_3$  and  $\tilde{W}_1$  are full rank [29, p.38]. Let us denote  $\zeta_i^T = (z_i, \dots, z_r)$ ,  $i \geq r^{wu} + 1$ . We first prove that (c)  $\Rightarrow$  (a). If  $z_j(\cdot)$ ,  $n + 1 \leq j \leq i - 1$  are functions and  $z_{i-1}(\cdot)$  is discontinuous, then  $\deg(\lambda) = r + 2 - i$ , and reciprocally. We notice that  $\deg(\lambda) = r + 2 - i$ ,  $i \leq r - 1$ , implies that  $\deg(\lambda) \geq 3$ . The vector  $\zeta_i$  may be considered as a  $(n - 1)$ -vector of distributions of degrees

$\geq 2$ . Since  $\eta = -(\tilde{W}_3^{-1} \tilde{W}_2 \tilde{W}_1^{-1}) \tilde{z} + \tilde{W}_3^{-1} \begin{pmatrix} z_{n+1} \\ \vdots \\ z_{i-1} \\ \zeta_i \end{pmatrix}$ , then  $B^T \eta$  and  $A^T \eta$  are time functions if and only

if  $\tilde{W}_3^{-1} \begin{pmatrix} 0^{i-n+1} \\ \zeta_i \end{pmatrix}$  belongs to  $\text{Ker}(B^T)$  and to  $\text{Ker}(A^T)$ . However since the pair  $(A, B)$  is controllable, the Kalman controllability matrix  $\mathcal{K}$  and its transpose are square full rank  $(n \times n)$  matrices. Since

$\mathcal{K}^T = \begin{pmatrix} B^T \\ B^T A^T \\ B^T (A^T)^2 \\ \vdots \\ B^T (A^T)^{n-1} \end{pmatrix}$ , one sees that  $\mathbb{R}^n \ni v \in \text{Ker}(B^T) \cap \text{Ker}(A^T)$  implies that  $\mathcal{K}^T v = 0$ . Hence

<sup>4</sup>We can safely differentiate  $B^T \eta(t)$  as  $B^T \dot{\eta}(t)$ , as we know that the function  $B^T \dot{\eta}(\cdot)$  is continuous at  $t$ .

$v = 0$  and we deduce that  $\zeta_i = 0$  which is a contradiction. Therefore  $\deg(B^T \eta) \leq 1$  and  $\deg(A^T \eta) \leq 1$  implies that  $\deg(\lambda) \leq 2$ . The other implications/equivalences follow from the  $\tilde{z}$ -dynamics. The last equivalence is just a rewriting of **(a)**  $\Leftrightarrow$  **(c)**. ■

Consequently one sees that the only case where  $B^T \eta$  and  $A^T \eta$  are functions, is when all  $z_i(\cdot)$ ,  $r^{wu} + 1 \leq i \leq r$ , are functions. In this latter case  $z_r(\cdot)$  may be time discontinuous, so that  $\lambda$  is a measure. However there may exist cases when  $\deg(\lambda) \geq 3$ ,  $\deg(B^T \eta) \leq 1$  and  $\deg(A^T \eta) \geq 2$ .

**Proposition 6** *Let  $m = n_u = 1$  and  $r^{wu} = n$ . Let us embed the necessary condition system in (4) in the HOSP in (17)-(21). Let a trajectory of (4) make contact with  $\partial\Phi$  at  $t = \tau$ . Assume that  $\bar{z}(\tau^-) = 0$ . Then  $\deg(\lambda) \leq n$  and  $x(\cdot)$  is a continuous function at  $\tau$  while  $\deg(B^T \eta) \leq 1$ . Finally  $\deg(A^T \eta) \leq n - 1$  at  $t$ .*

**Proof:** We notice first that  $\bar{z}(\tau^-) = 0$  ( $= \bar{z}(\tau^+)$  from (28)-(31)) is equivalent to having  $x(\cdot)$  is continuous at  $\tau$ , because of the minimality of  $(A, B, C)$  and consequently the structure of the transformation matrix  $\tilde{W}$  (since  $r^{wu} = n$ , the matrix  $\tilde{W}_1$  is the Kalman observability matrix, see Proposition 5). Using the  $\tilde{z}$ -dynamics it is easily deduced that  $\deg(\lambda) \leq r^{wu}$  at  $\tau$  and that  $\deg(B^T \eta) \leq 1$ . The last statement is a consequence of the structure of  $\tilde{A}$  in (4). ■

It is worth recalling that all the distributions which appear in the problem, have a atomic part whose support is of zero Lebesgue measure, see remark 5 v). As a consequence of Proposition 6, when  $\bar{z}(\tau^-) = 0$  at all entry times  $\tau$ , the state jumps are jumps of  $\{\eta\}(\cdot)$ , i.e.  $\{\eta\}(t_k^+) \neq \{\eta\}(t_k^-)$ . But imposing the constraint  $\bar{z}(\tau^-) = 0$  at all entry times  $\tau$  does not imply that  $\lambda$  is a measure.

**Example 2** *Let us consider a triple integrator  $x^{(3)}(t) = u(t)$  with  $x_1(t) \geq 0$ . We have  $\tilde{A} =$*

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 & -1 & 0 \end{pmatrix},$$

$$\tilde{W} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 1 & 0 & 1 \end{pmatrix}, \tilde{B}^T = (0 \ 0 \ 0 \ -1 \ 0 \ 0), B^T \eta = \eta_3, A^T \eta = (0 \ \eta_1 \ \eta_2)^T.$$

*Assume that  $\lambda = \delta_t$ , which happens if  $z_5(\cdot)$  jumps at  $t$  while  $z_6(\cdot)$  is continuous (for instance if  $z_1(t^-) = z_2(t^-) = z_3(t^-) = z_4(t^-) = 0$ ,  $z_5(t^-) < 0$ ,  $z_6(t^-) > 0$ ). From the HOSP dynamics, this implies that  $z_1(\cdot)$ ,  $z_2(\cdot)$ ,  $z_3(\cdot)$ ,  $z_4(\cdot)$  and  $z_6(\cdot)$  are continuous at  $t$ , while  $z_5(t^+) = 0$ . Thus  $\deg(\eta_1) = 0$ ,  $\deg(\eta_2) = 2$  and  $\deg(\eta_3) = 0$ . This is a case where  $\deg(B^T \eta) = 0$  while  $\deg(A^T \eta) = 2$ , and  $\deg(\lambda) = 3$ . This is in agreement with Proposition 5: if  $\deg(\lambda) = 3$ , necessarily either  $\deg(B^T \eta) \geq 2$  or  $\deg(A^T \eta) \geq 2$ . Since  $\deg(B^T \eta) \leq 1$  necessarily  $\deg(A^T \eta) \geq 2$ .*

The preceding analysis suggests that there are two basic situations concerning the system in (4) (a) (c):

- (1)  $\lambda$  is a measure, i.e.  $\lambda = g_r(t)dt + d\mathcal{J}_r$ ,
- (2)  $\deg(\lambda) \geq 3$  and
  - (2.1)  $B^T \eta(\cdot)$  and  $x(\cdot)$  are time functions on  $[0, T_1]$ ,  $\deg(A^T \eta) \geq 2$ ,
  - (2.2)  $\deg(\lambda) \geq 3$ ,  $\deg(B^T \eta) \geq 2$ .

Quantum electronics and laser control, portfolio optimisation, optimisation of the loan policy of a company, are problems which do involve optimal controls with  $\deg(u(\cdot)) = 2$  [21]. Considering optimal

inputs which are not functions therefore makes sense from the physical point of view. There are several basic ingredients in the optimality conditions, among which: i) the value of the optimal control in (5) which comes from  $\frac{\partial H}{\partial u}(x, u, \eta) = 0$ , ii) the costate equation  $\dot{\eta}(t) = Qx(t) - A^T \eta(t) - C^T \lambda$  which comes from calculus of variations and Kuhn-Tucker conditions stating that  $\dot{\hat{x}}(t) = J \frac{\partial H_{n_s}(\hat{x})}{\partial \hat{x}}(t)$  (see (9) and paragraph after), iii) the expression of the integral action (or cost) in (1).

In case (1) the three ingredients have a meaning and the BVP in (4) represents the optimality necessary conditions with  $I(u)$  in (1) [4, Theorem 1] [8][47, (4.1)-(4.5)]. From [47, Proposition 4.1] normal extremals are minimizers of (1) (2) in such a case. Also we note that  $x = \{x\}$ ,  $u = \{u\}$ , and  $\eta = \{\eta\}$  on  $[0, T_1]$ .

**Corollary 4** *Assume the conditions of Proposition 6 are satisfied and that  $\mathbf{U}$  is a set of functions of time. Let  $\tau \in [0, T_1]$  be an entry time. Let us consider the solution  $\hat{x} = \tilde{W}^{-1} \bar{z}$  of the distributional formalism in (17) (18) (19) and the function  $H(x, u, \eta)$  in (10). Then the control which satisfies  $\frac{\partial H}{\partial u}(x, u, \eta) = 0$  at  $\tau$  is given in (5) if and only if  $(\sigma_{z_{n+1}}(\tau), \dots, \sigma_{z_r}(\tau)) \tilde{W}_3^{-T} \in \text{Ker}(A)$ . In this case  $\deg(\lambda) \leq 2$ .*

**Proof:** From Proposition 6, one sees that  $\deg(A^T \eta) \leq r^{wu} - 1$  at  $\tau$ . One may rewrite the Hamiltonian as  $H(x, u, \eta) = -\frac{1}{2} x^T Q x - \frac{1}{2} u^T R u + \eta^T B u + \langle A^T \eta, x \rangle$  since  $A^T \eta \in \mathcal{T}_{n-2}$  is a distribution. At  $\tau$  we may

write  $\eta$  as  $\eta = \sum_{i=0}^{n-3} \eta_i \delta_\tau^{(i)}$ . We know that  $\eta = -\tilde{W}_3^{-1} \tilde{W}_2 \tilde{W}_1^{-1} \bar{z} + \tilde{W}_3^{-1} \begin{pmatrix} z_{n+1} \\ \vdots \\ z_r \end{pmatrix}$ . Since  $\bar{z}(\cdot)$  is a continuous

function at  $\tau$ , we deduce that  $\eta_i = \tilde{W}_3^{-1} \beta_i$ , with  $\beta_i = \begin{pmatrix} 0^{r-i-1} \\ \sigma_{z_{r-i}}(\tau) \\ 0^{2n-r+i} \end{pmatrix}$  (from (30) (31) either  $\beta_i = 0$  or

$\beta_i = -z_{r-i}(\tau^-) > 0$ ). The result follows because  $\deg(A^T \eta) \leq 1$  if and only if  $(\beta_{n-1}, \dots, \beta_0)^T \in \text{Ker}(A^T)$  (otherwise the product  $\langle A^T \eta, x \rangle$  necessarily involves  $u(\cdot)$  and its derivatives). The last result follows from Propositions 5 and 6. ■

In other words, the optimal input computed along the solution of the distributional problem (17) (18) (19) is given by (5) at a tangential contact if and only if  $\deg(\lambda) \leq 2$ . When  $\deg(\lambda) \geq 3$  one will have to resort to the measure differential problem whose solutions are time functions in  $\mathcal{F}_\infty$  to characterize the optimality, see sections 3.5 and 3.6. Let us note that the constraints that are imposed in Corollary 4 and in Proposition 4 at time  $\tau$ , are of the same nature as the usual constraint  $\bar{z}(\tau^-) = 0$  at an entry time.

### 3.5 An extended integral action $I(u)$

Let us deal now with case (2), i.e.  $\deg(\lambda) \geq 3$ . The case (2.1) is of particular interest because it shows that the optimal controller may be a function of time, while the costate  $\eta$  is a distribution of degree  $\geq 3$ . Thus decreasing the ‘‘regularity’’ of  $u(\cdot)$  while keeping  $\deg(u) \leq 1$  may create serious difficulties in the analysis of this optimisation problem, the costate equation being a distributional differential equation as (39). On the sets  $\text{supp}(d\mathcal{J}_i)$ , the dynamics of the system in (4) (a) (c) becomes algebraic, see (31) (30). Consequently when  $\deg(\lambda) \geq 3$  it is necessary that the optimal control problem in (1) (2), involves an algebraic action in addition to an integral action. One path is the following:  $u^2$  has no meaning if  $u = \delta_0$ . However  $\langle \delta_0, \varphi(t) \rangle^2 = \varphi(0)^2$  has a meaning for any continuous test function  $\varphi(\cdot)$ . Also  $\langle a_1 \delta_0 + a_2 \dot{\delta}_0, \varphi \rangle^2 = a_1^2 \langle \delta_0, \varphi \rangle^2 - 2a_1 a_2 \langle \delta_0, \varphi + \dot{\varphi} \rangle + a_2^2 \langle \delta_0, \dot{\varphi} \rangle^2 = a_1^2 \varphi^2(0) - 2a_1 a_2 (\varphi(0) + \dot{\varphi}(0)) + a_2^2 \dot{\varphi}^2(0)$ . This way of introducing the singular distributions contributions in the quadratic cost, is different from that considered in [16]. Let  $\{\delta_0^n(\cdot)\}$  be a fundamental sequence for the Dirac measure  $\delta_0$  [11, Appendix A]. In [16] it is argued that since  $\int_0^{+\infty} (\delta_0^n)^2(t) dt \rightarrow +\infty$  as  $n \rightarrow +\infty$ , then the cost to be associated to the distributional parts is  $+\infty$ . It is concluded in [16] that only those optimal control problems involving no distributions of degree  $\geq 2$  make sense. Our approach of considering the cost function (or action) when higher degree distributions are present in the optimal control problem, rather follows Moreau’s results in [36]. An important tool in the developments which follow is the measure differential formalism of the HOSP in (22)-(26).

Let us denote  $d\nu_{\bar{z}} = (d\nu_1 \dots d\nu_r)^T \in \mathbb{R}^r$ . Similarly for  $d\mathcal{J}_{\bar{z}}$  which denotes the atomic part of the vector measure  $d\nu_{\bar{z}} = g_{\bar{z}}(t)dt + d\mathcal{J}_{\bar{z}}$  (see (32)). Let us define  $\mathcal{M} = \begin{pmatrix} I_{r-1} & 0^{r-1} \\ 0_{r-1} & M^{(r)} \end{pmatrix} \in \mathbb{R}^{r \times r}$ , and  $\tilde{\mathcal{M}} = \begin{pmatrix} \mathcal{M} & 0_{r \times (2n-r)} \\ 0_{(2n-r) \times r} & 0_{(2n-r) \times (2n-r)} \end{pmatrix} \in \mathbb{R}^{2n \times 2n}$ . From (27) we have  $d\tilde{z} = \tilde{A}\tilde{z}(t)dt + \tilde{\mathcal{M}} \begin{pmatrix} d\nu_{\bar{z}} \\ 0_{2n-r} \end{pmatrix}$ . Thus  $d\nu_{\bar{x}} = \tilde{W}^{-1}\tilde{\mathcal{M}} \begin{pmatrix} d\nu_{\bar{z}} \\ 0_{2n-r} \end{pmatrix} = g_{\bar{x}}(t)dt + d\mathcal{J}_{\bar{x}}$ .

**Lemma 4** *Let  $r^{wu} = 2k$ ,  $k \geq 0$ , and  $m = 1$ . The HOSP jump rule in (30) (31) minimizes the quadratic term  $\frac{1}{2}\langle d\mathcal{J}_{\bar{z}}, \varphi \rangle^T \mathcal{M} \langle d\mathcal{J}_{\bar{z}}, \varphi \rangle$  for any test function  $\varphi(\cdot) \in C^0[\mathbb{R}, \mathbb{R}]$  with support containing  $[0, T_1]$ , under the constraints in (36) and (37).*

**Proof:** One has  $d\nu_i(\{t\}) = \{z_i\}(t^+) - \{z_i\}(t^-)$  for all  $t \in [0, T_1]$ . If  $t$  is not an atom of  $d\nu_i$  then  $d\nu_i(\{t\}) = d\mathcal{J}_i(\{t\}) = 0$ . If  $t$  is an atom of  $d\nu_i$  this implies that all the tangent cones  $T_{\Phi}^k(\{z_1\}(t^-), \dots, \{z_k\}(t^-))$  satisfy  $T_{\Phi}^k(\{z_1\}(t^-), \dots, \{z_k\}(t^-)) = \mathbb{R}^+$  for  $0 \leq k \leq i-1$ . Then  $d\nu_i(\{t\}) = d\mathcal{J}_i(\{t\}) = \operatorname{argmin}_{\sigma_i + \{z_i\}(t^-) \geq 0} \frac{1}{2}\sigma_i^2$ . Indeed from (30) it follows that  $\{z_i\}(t^+) = \operatorname{prox}[T_{\Phi}^{i-1}(z_1(t^-), \dots, \{z_{i-1}\}(t^-)); \{z_i\}(t^-)]$ . We deduce that  $\{z_i\}(t^+) = \operatorname{argmin}_{\sigma \in T_{\Phi}^{i-1}(z_1(t^-), \dots, \{z_{i-1}\}(t^-))} \frac{1}{2}(\sigma - \{z_i\}(t^-))^2$ . In case  $\{z_i\}(t^-) < 0$  and  $T_{\Phi}^{i-1}(z_1(t^-), \dots, \{z_{i-1}\}(t^-)) = \mathbb{R}^+$  (which are the conditions required so that  $d\mathcal{J}_i$  possesses an atom at  $t$ ) the result follows. We have  $\langle d\mathcal{J}_{\bar{z}}, \varphi \rangle^T \mathcal{M} \langle d\mathcal{J}_{\bar{z}}, \varphi \rangle = \sum_{i=1}^{r-1} \langle d\mathcal{J}_i, \varphi \rangle^2 + M^{(r)} \langle d\mathcal{J}_r, \varphi \rangle^2$ , and from Lemma 1  $M^{(r)} > 0$ . The result is proved.  $\blacksquare$

Let us illustrate the second argument of the proof. If  $\{z_i\}(t^-) < 0$  then  $\operatorname{argmin}_{\sigma_i + \{z_i\}(t^-) \geq 0} \frac{1}{2}\sigma_i^2 = -\{z_i\}(t^-)$ . Since  $d\nu_i(\{t\}) = \{z_i\}(t^+) - \{z_i\}(t^-)$  we deduce that  $\{z_i\}(t^+) = 0$  (a plastic jump). The parallel with unilateral Mechanics and Moreau's second order sweeping process can be done, where the post-impact velocity satisfies  $\dot{q}(t^+) = \operatorname{argmin}_{w \in T_{\Phi}(q(t))} \frac{1}{2}(w - \dot{q}(t^-))^T M(q(t))(w - \dot{q}(t^-))$  when the impacts are plastic [11, Remark 5.11].

Lemma 4 extends to odd- $r^{wu}$  systems, provided  $d\mathcal{J}_r = 0$ . We shall find again the specificity of odd- $r^{wu}$  systems in sections 4.2 and 5. Let us propose the following extended action to be minimized:

$$\begin{aligned} \operatorname{minimize}_{u(\cdot) \in \mathbf{U}} \quad & \frac{1}{2} \langle d\mathcal{J}_{\bar{x}}, \varphi \rangle^T \tilde{W}^T \tilde{\mathcal{M}} \tilde{W} \langle d\mathcal{J}_{\bar{x}}, \varphi \rangle + \frac{1}{2} \int_0^{T_1} [\{x\}(t)^T Q \{x\}(t) + \{u\}(t)^T \{u\}(t)] dt \\ & + \frac{1}{2} \{x\}^T(T_1) F \{x\}(T_1) \end{aligned} \quad (41)$$

for any test function  $\varphi(\cdot) \in C^0[\mathbb{R}, \mathbb{R}]$ , with support containing  $[0, T_1]$ ,

**Definition 2** *Let  $\mathbf{U}$  be a set of distributions in  $\mathcal{T}_{r^{wu}}([0, T_1])$  such that the atoms and points of non-analyticity of any  $u \in \mathbf{U}$  are in the set  $\bigcup_{i=1}^r \operatorname{supp}(d\mathcal{J}_i)$ . We denote it as  $\mathbf{U}_{\mathcal{J}}^{r^{wu}}$ .*

The definition of  $\mathbf{U}_{\mathcal{J}}^{r^{wu}}$  is motivated by the fact that the HOSP solutions are constrained to exist in the set  $\mathcal{T}_{r-1}([0, T_1])$ . In particular  $\deg(z_1) \leq 1$ . Having  $u \in \mathbf{U}_{\mathcal{J}}^{r^{wu}}$  means that  $\deg(u) \leq r^{wu} + 1$ , while along solutions of the HOSP  $\deg(\lambda) \leq r + 1 = 2r^{wu} + 1$ . The inspection of the  $z$ -dynamics shows that  $\deg(z_1) \leq 1 \Rightarrow \deg(u) \leq r^{wu} + 1$ . It is consequently useless to look for a larger set (in the set of distributions as defined in section 3.1) of inputs for this optimal control problem. As example 6 shows it is easy to construct cases where  $\bar{x}_1$  is not reachable from  $\bar{x}_0$  in the set  $\mathbf{U}_{\mathcal{J}}^{r^{wu}}$ . This is therefore an important issue.

Let us notice that  $\langle d\mathcal{J}_{\bar{x}}, \varphi \rangle^T \tilde{W}^T \tilde{\mathcal{M}} \tilde{W} \langle d\mathcal{J}_{\bar{x}}, \varphi \rangle$  contains terms of the form  $\varphi^2(t_k)$  for all  $t_k \in \operatorname{supp}(d\mathcal{J}_{\bar{x}})$ , since  $d\mathcal{J}_{\bar{x}}$  is an atomic measure generated by a right continuous jump function  $\mathcal{J}_{\bar{x}}(\cdot)$ . Let us recall that in (41) the state jump times are totally free and are not a priori fixed. The only thing that we impose through the HOSP formalism is that the supports of the atomic measures  $d\mathcal{J}_i$  are orderable and countable sets. It is noteworthy that even in case (2.1), the addition of an algebraic action is necessary. At a time  $t \in \operatorname{supp}(d\nu_i)$  where (2.1) holds, then the value of  $\{u\}(\cdot) = B^T \{\eta\}(\cdot) = B^T \eta$  is irrelevant (and this is the case on the set of all atoms of  $\lambda$  as this set is of zero Lebesgue measure). The optimisation is done at such a  $t$  on the costate  $\eta$ .

**Remark 6** *It is important to notice that problem (1)-(2) is not the same as the optimal control problem (without state constraints) where the control  $u$  may contain a singular distribution [51, 43]. The major discrepancy is that the support of the atomic part satisfies (7). For instance in example 6, jumps in  $x_1(\cdot)$  are allowed only on the line  $x_2 = 0$  according to our framework. It is also quite different from [49] who studies optimal control of a system which undergoes state jumps. One may argue that when  $\deg(u) \geq 2$  then distributional inputs could be applied without the support condition (7), as in [43]. But as seen above an important case is when  $\deg(u) \leq 1$  and  $\deg(\eta) \geq 2$ . The impulsive optimal control problem presented in [24] is also of a different nature, despite it shares some common features with what will be presented in the next two sections (especially the splitting of the intergral action into a “continuous” and a “jump” actions). The idea of separating functions and distributions in the action is not new, see e.g. [40] and [3, equ.(2.5)]. Also in [13, §3.7]  $I(u)$  is split into an integral part and an algebraic part taking into account state jumps. However the works [13, 40] do not apply to the problem considered in this paper as they do not involve state unilateral constraints.*

**Proposition 7** *Suppose that  $\bar{x}_1$  is reachable from  $\bar{x}_0$  in a finite time  $T_1 \geq 0$ , with  $u \in \mathbf{U}_{\mathcal{J}}^{r^{uu}}$ . Let  $r^{uu}$  be even. Then the HOSP measure differential inclusion in (22)–(24) represents the optimality necessary conditions for the extended integral action  $I(u)$  in (41), with  $u = B^T \eta$  for all  $t \in [0, T_1]$  and subject to the constraints in (36) and (37). If  $r^{uu}$  is odd, then the result holds provided the additional constraint  $d\mathcal{J}_r = 0$  holds along the optimal trajectory.*

**Proof:** In view of the material of sections 3.1, 3.2 and 3.4, one sees that if  $\deg(\lambda) \geq 2$  an algebraic action has to be added to  $I(\{u\})$  in (1), taking into account the state  $\{x\}(\cdot)$  and/or costate  $\{\eta\}(\cdot)$  jumps. From the measure differential formalism (40),  $u = B^T \eta$  for all  $t \in [0, T_1]$  means that

$$u(t) = B^T \eta(t) + B^T \int_{[0,t]} d\nu_{\eta}, \quad (42)$$

where  $d\nu_{\eta}$  is defined after (40) and  $\eta(\cdot)$  in (42) is to be understood as the “unconstrained” solution of the ODE  $\dot{\eta}(t) - Qx(t) + A^T \eta(t) = 0$  with  $\dot{x}(t) - Ax(t) - BR^{-1}B^T \eta(t) = 0$ . In other words (42) stems from

$$u(t) = \operatorname{argmax}_{u \in \mathbf{U}_{\mathcal{J}}^{r^{uu}}} H \left( x(t) + \int_{[0,t]} d\nu_x, \eta(t) + \int_{[0,t]} d\nu_{\eta}, u(t) \right) \quad (43)$$

and it is worthwhile noting that the arguments in the Hamiltonian function in (43), are functions of time (the fundamental role played by the measure differential formalism (22)–(24) and Proposition 1 is clear in this context). Let us recall that solutions of the measure differential formalism are  $\mathcal{F}_{\infty}$  functions, and so is  $u(\cdot)$  in (42). From the distributional formalism (39)  $u = B^T \eta$  means that  $\langle u, \varphi \rangle = \langle B^T \eta, \varphi \rangle$  for all  $\varphi \in C_0^{\infty}([0, T_1])$ . Therefore, if the necessary condition system (4) embedded in the distributional HOSP (17)–(19), possesses a distributional solution in  $\mathcal{T}_{r-1}([0, T_1])$ , the corresponding solution of the measure HOSP in (22)–(24) decomposed as in (42) into an unconstrained and a constrained parts, satisfies (42) and (43) (by Proposition 1 there is a bijective correspondence between solutions of the distributional and the measure formalisms of the HOSP). The Hamiltonian function  $H(\{x\}, \{u\}, \{\eta\})$  evaluated along the optimal trajectory (and optimal control) satisfies  $H(\{x\}, \{u\}, \{\eta\}) = \sup_{\{v\} \in \mathbf{U}_{\mathcal{J}}^r} -\frac{1}{2}\{x\}^T Q \{x\} - \frac{1}{2}\{v\}^T R \{v\} + \{\eta\}^T (A\{x\} + B\{v\})$  Lebesgue almost everywhere (equivalently  $\{u\}(\cdot) = B^T \{\eta\}(\cdot)$  almost everywhere on  $[0, T_1]$ ). Also the costate equation is  $\frac{d}{dt}(\{\eta\})(t^+) = Q\{x\}(t^+) - A^T \{\eta\}(t^+) - C^T(0_{2n-1} \ g_r(t^+))^T$  on intervals  $(t_k, t_{k+1}]$  since the solutions are right-continuous. We recall that the values of the functions  $\{x\}(\cdot)$ ,  $\{u\}(\cdot)$ ,  $\{\eta\}(\cdot)$  are irrelevant at the atoms  $t_k$  of the distribution  $\lambda$ . Let us now deal with the algebraic action in (41). One has  $\tilde{z} = \tilde{W}\tilde{x}$ , and consequently  $\langle d\mathcal{J}_{\tilde{z}}, \varphi \rangle^T \mathcal{M} \langle d\mathcal{J}_{\tilde{z}}, \varphi \rangle = \langle d\mathcal{J}_{\tilde{z}}^T, \varphi \rangle \tilde{\mathcal{M}} \langle d\mathcal{J}_{\tilde{z}}, \varphi \rangle = \langle d\mathcal{J}_{\tilde{x}}, \varphi \rangle^T \tilde{W}^T \tilde{\mathcal{M}} \tilde{W} \langle d\mathcal{J}_{\tilde{x}}, \varphi \rangle$ . The solutions are distributions in  $\mathcal{T}_{r-1}([0, T_1])$ , solutions of the HOSP in (17)–(19). From Proposition 1 we can consider equivalently the solutions of the measure formalism in (22)–(24). From Lemma 4, the solution of (22)–(24) minimizes the quadratic term  $\frac{1}{2} \langle d\mathcal{J}_{\tilde{x}}, \varphi \rangle^T \tilde{W}^T \tilde{\mathcal{M}} \tilde{W} \langle d\mathcal{J}_{\tilde{x}}, \varphi \rangle$ . From the definition of  $d\nu_{\tilde{z}}$ , it follows that if contact states are forced to be hypertangential (see Definition 3 below), then  $d\mathcal{J}_{\tilde{x}} = 0$  and the action is

purely integral. All the above applies to even  $r^{wu}$ , and to odd  $r^{wu}$  provided the constraint  $d\mathcal{J}_r = 0$  is added. ■

In case  $u(\cdot)$  is a function, its derivatives may be distributions. This is taken into account by the algebraic action which is a function of  $\eta$ , hence implicitly of  $u(\cdot)$  and its derivatives. Consider example 2. One has at the considered time  $u(t) = \eta_3(t)$ . From (17) the distributional equality  $Du = D\eta_3 = x_3 - \eta_2 = z_3 - \eta_2$  shows that  $\deg(Du) = 2$  since  $\deg(\eta_2) = 2$ . This is the consequence of a jump in  $z_5(\cdot)$  at  $t$  which is the optimal jump according to the action in (41). Most importantly, when  $\lambda$  is a measure, (42) and (43) coincide with [58, equ. (11.42)] and [58, equ. (11.41)], respectively. We note however that the solution and optimal controller regularity conditions as studied in [58, §9, §10, §11] are not relevant in our framework since the data in (1) and (2) are analytic and the solution we are looking for is regular (Definition 1). It is worthwhile pointing out that in case (2.1) the integral term of the action  $I(u)$  is equal to  $\int_0^{T_1} [x(t)^T Qx(t) + u(t)^T u(t)]dt + \frac{1}{2}x^T(T_1)Fx(T_1)$ . However from Corollary 4 one does not have  $\frac{\partial H}{\partial u} = 0$  at atoms of  $d\nu_i$ ,  $1 \leq i \leq r-1$ , but only almost everywhere on  $[0, T_1]$  (see also (42)). The algebraic action is then equal to  $\frac{1}{2}\langle (0_n \ d\mathcal{J}_\eta^T), \varphi \rangle \tilde{W}^T \tilde{\mathcal{M}} \tilde{W} \langle (0_n \ d\mathcal{J}_\eta^T)^T, \varphi \rangle$ .

**Example 3** Following [13, §3.11] let us consider the system with relative degree  $r^{wu} = 2$ ,  $l > 0$ ,

$$\begin{cases} \dot{x}_1(t) = x_2(t) \\ \dot{x}_2(t) = u(t) \\ w(t) = l - x_1(t) \geq 0 \end{cases} \quad (44)$$

with  $Q = 0$  and  $R = 1$ . Then (4) becomes in the  $\tilde{z}$ -canonical representation

$$\begin{cases} \dot{z}_1(t) = z_2(t) \quad (= -x_2(t)) \\ \dot{z}_2(t) = z_3(t) \quad (= -\eta_2(t)) \\ \dot{z}_3(t) = z_4(t) \quad (= \eta_1(t)) \\ \dot{z}_4(t) = \lambda(t) \\ 0 \leq w(t) = l + z_1(t), \quad z_1(t) = -x_1(t) \end{cases} \quad (45)$$

Let us consider the end-point conditions  $z_1(0^-) = z_1(T_1^+) = -l$  and  $z_2(0^-) = z_2(T_1^+) < 0$ . Then the optimal solution is  $\tilde{z}(t) = 0$  for all  $t \in (0, T_1)$  (so that  $u(t) = 0$  for all  $t \in (0, T_1)$  and the integral action is zero along the optimal path), and  $z_2(\cdot)$  jumps at  $t = 0$  and  $t = T_1$  (here  $z_2(\cdot)$  is to be considered as a solution of the measure differential formalism (22) (23) (24)). Thus  $d\nu_2 > 0$  and  $d\nu_1 = d\nu_3 = d\nu_4 = 0$ . The algebraic action is equal to  $\langle d\nu_2, \varphi \rangle^2$ . Along the solutions of the distributional formalism (17) (18) (19) one has  $u = \eta_2 = -z_3$  as an equality of distributions, and  $\deg(u) = 2$ . The optimal solution that is presented in [13, §3.11] is such that  $z_4(t_1^-) < 0$  at an entry time  $t_1$ , so that  $\deg(\lambda) = 2$ . An open issue is to find whether or not admitting jumps in  $z_3(\cdot) = -\eta_2(\cdot)$  (solutions of (22) (23) (24)), hence a discontinuous optimal controller, would allow one to decrease the action value compared to the case when  $\lambda$  is constrained to be a measure. In this case the distributional formalism solution yields  $\deg(\lambda) = 3$ . In this simple example, the test merely amounts to varying the entry and exit time, allowing for  $\deg(\lambda) \geq 3$ .

### 3.6 Hamilton's principle of Mechanics

Let us briefly expose here some facts about Hamilton's principle of Mechanics for Lagrangian systems subjected to complementarity relations and impact laws. As recalled in [11, §3.5], the addition of a unilateral constraint and of an impact law, implies that Hamilton's principle has to be modified in order to keep its meaning (otherwise it cannot represent the nonsmooth feature of the dynamics). The position is supposed to be constrained in  $\Phi = \{q \mid h(q(t)) \geq 0\}$  with  $h : \mathbb{R}^n \rightarrow \mathbb{R}$ , and we take  $T_1 = 1$ . The Lagrangian dynamics is embedded into Moreau's second order sweeping process. Let us first consider the following minimization problem with  $e \in [0, 1]$ .

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2}\sigma_{\dot{q}}(t)^T M(q(t))\sigma_{\dot{q}}(t) + \int_0^1 L(q(t), \dot{q}(t))dt \\ & \begin{cases} q(0) = q_0 \in \Phi \\ q(1) = q_1 \in \Phi \\ \dot{q}(t^+) + e\dot{q}(t^-) \in T_\Phi(q(t)) \end{cases} \end{aligned} \quad (46)$$



The solution of the sweeping process on  $[0, 1]$  which satisfies  $q(0) = q_0$  and  $q(1) = q_1$ , minimizes the action (46) as it satisfies the smooth Euler-Lagrange dynamics outside impacts, and minimizes the quadratic term at impacts. However the writing in (46) is not satisfactory because it explicitly resorts to impact times in the formulation of the algebraic term of the action (notice that if  $t$  is not an impact time then the corresponding algebraic action becomes trivial). Let us define the Stieltjes measure  $dv$  associated to the acceleration, whose atoms are the impact times so that  $d\mathcal{J}_v(\{t\}) = \dot{q}(t^+) - \dot{q}(t^-)$ . Then one may rewrite (46) as

$$\text{minimize}_{\substack{q(0) = q_0 \in \Phi \\ q(1) = q_1 \in \Phi \\ \dot{q}(t^+) + e\dot{q}(t^-) \in T_{\Phi}(q(t))}} \frac{1}{2} \langle d\mathcal{J}_v, \varphi \rangle^T M(q(t)) \langle d\mathcal{J}_v, \varphi \rangle + \int_0^1 L(q(t), \dot{q}(t)) dt \quad (47)$$

with  $\varphi(\cdot) \in C^0$  a test function whose support contains  $[0, 1]$ . As shown in [14] for the case  $e = 1$ , the extremal of the functional

$$\int_0^1 (L(q(t), \dot{q}(t)) dt + h(q(t))\lambda) \quad (48)$$

under the constraints  $h(q(t)) \geq 0$  for all  $t \in [0, 1]$ , some initial data and no kinetic energy loss, is a solution to the bounce problem, i.e.  $\frac{d}{dt} \left( \frac{\partial L(q, \dot{q})}{\partial \dot{q}} \right) (t) - \frac{\partial L(q, \dot{q})}{\partial q} (t) = \nabla h(q(t))\lambda$  on  $[0, 1]$ , and vice-versa. The multiplier  $\lambda$  in (48) is a measure whose support is in the set  $\{t \in [0, 1] \mid h(q(t)) = 0\}$ .

The great advantage of the formulation of the BVP as in (48) compared to (46) and (47), is that it allows one to perform a time-stepping discretization of the integrand in one shot, using the fact that  $\lambda$  is a Stieltjes measure. One then considers the time-discretization of the augmented minimization problem which can be treated as a quadratic programme, thereby extending the direct methods as described in [56]. One sees that in (48) the Lagrangian function is augmented to cope with the unilaterality, and the multiplier is a measure. Mimicking [14] the basic idea is to augment the Lagrangian function  $L(x, u) = \frac{1}{2}x^T Qx + \frac{1}{2}u^T Ru$ , taking into account the specific features of the problem. Let us recall that  $dv_{\bar{z}} = (dv_1 \dots dv_r)^T$ ,  $dv_{\bar{x}} = \tilde{W}^{-1} \tilde{\mathcal{M}} \begin{pmatrix} dv_{\bar{z}} \\ 0_{2n-r} \end{pmatrix}$ , and  $d\bar{x} - \tilde{A}\bar{x}(t)dt - dv_{\bar{x}} = 0$  from (22)–(24). We define the augmented Lagrangian as

$$\bar{L}(\{x\}, \{u\}, dv_{\bar{x}}, dt) = L(\{x\}, \{u\})dt + \{\bar{x}\}(t^+)^T dv_{\bar{x}}$$

The product  $\{\bar{x}\}(t^+)^T dv_{\bar{x}}$  is written here with a strong abuse of notation and requires some care. The function  $\{\bar{x}\}(\cdot)$  is right-continuous and  $dv_{\bar{x}}$  is a measure with possible atoms at times  $t_k$ ,  $k \geq 0$ . The space of functions which are  $dv_{\bar{x}}$ -integrable contains functions continuous at  $t_k$ , and also the functions which are  $dv_{\bar{x}}$ -almost everywhere equal to an integrable and continuous function  $g(\cdot)$ . Since the supports of the atoms are the singletons  $\{t_k\}$ , it is sufficient that  $\{\bar{x}\}(t_k) = g(t_k)$ . Then, denoting the atoms of  $dv_{\bar{x}}$  as  $\delta_{t_k}$  one has

$$\int \{\bar{x}\} d\delta_{t_k} = \int g d\delta_{t_k} = g(t_k) = \{\bar{x}\}(t_k) = \{\bar{x}\}(t_k^+) \quad (49)$$

Equality (49) shows that  $\int \bar{L}(\{x\}, \{u\}, dv_{\bar{x}}, dt)$  is meaningful. We therefore propose an alternative to Proposition 7.

**Proposition 8** *The HOSP measure differential formalism in (22)–(24) represents the necessary optimality conditions for the optimal control problem*

$$\text{minimize}_{u(\cdot) \in \mathbf{U}_{\mathcal{J}}^{r, wu}} I(u) = \int_0^{T_1} \bar{L}(\{x\}, \{u\}, dv_{\bar{x}}, dt) + \frac{1}{2} \{x\}(T_1)^T F \{x\}(T_1) \quad (50)$$

subject to (2) and where the measures  $dv_i$  satisfy the inclusions (28) and  $u = B^T \eta$ .

**Proof:** Let us consider the augmented Hamiltonian measure  $\bar{H}(\{x\}, \{u\}, d\nu_{\bar{x}}, dt) = -\bar{L}(\{x\}, \{u\}, d\nu_{\bar{x}}, dt) + \{\eta\}^T (A\{x\} + B\{u\})dt$ . It is noteworthy that in the framework of (22)-(24) one can drop all brackets as  $\{x\}(\cdot) = x(\cdot)$ , because the solution is considered as a function in  $\mathcal{F}_\infty([0, T_1]; \mathbb{R})$ . Since the set of controllers is  $\mathbf{U}_{\mathcal{J}}^{r,uu}$  and  $u = B^T \eta$  (whose meaning is as in (42)), the support of the measure  $d\nu_{\bar{x}}$  is orderable and countable. Mimicking [14, remark 1.1] we deduce that the extremals of the functional  $I(u)$  satisfy

$$\begin{cases} dx = \frac{\partial}{\partial \eta} \left( -\frac{1}{2}L(x, u) \right) (t)dt - d\nu_x + Ax(t)dt + Budt \\ d\eta = -\frac{\partial}{\partial x} \left( -\frac{1}{2}L(x, u) \right) (t)dt - d\nu_\eta - A^T \eta(t)dt \\ u = B^T \eta \end{cases} \quad (51)$$

as an equality of differential measures for the first two lines of (51). The third equality has the meaning explained in the proof of Proposition 7 and allows one to calculate the optimal input using Proposition 1. One therefore has along solutions of the measure differential inclusion (51)  $u(t) = \operatorname{argmax}_{u \in \mathcal{F}_\infty} \bar{H}(x, u, d\nu_{\bar{x}}, dt)$  Lebesgue almost everywhere (notice that if  $u \in \mathbf{U}_{\mathcal{J}}^{r,uu}$  then  $\{u\}(\cdot) \in \mathcal{F}_\infty$ ). ■

We have therefore extended the measure differential inclusion necessary conditions in (4) (6) to the following

$$\begin{aligned} -d\tilde{x} + \tilde{A}\tilde{x}(t)dt + d\nu_{\tilde{x}} &= 0 \\ d\nu_{\tilde{x}} &= \tilde{W}^{-1} \tilde{\mathcal{M}} \begin{pmatrix} d\nu_{\tilde{z}} \\ 0_{2n-r} \end{pmatrix} \\ d\nu_{\tilde{z}} &= (d\nu_1 \dots d\nu_r)^T \\ d\nu_i &\in -\partial\psi_{T_\Phi^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)) \text{ for all } 1 \leq i \leq r \\ x(0) = \bar{x}_0, \eta(T_1) &= Fx(T_1) + C^T \gamma + \beta = F\bar{x}_1 + C^T \gamma + \beta = \eta_1 \end{aligned}$$

which reduces to (4) (6) when  $\lambda$  is a measure (see (36)). As said above, the objective is to perform a so-called direct approach to discretize the extended action (50) with a specific numerical algorithm that is able to approximate measures. Such a time-stepping algorithm is proposed in [1] for the integration of IVPs.

## 4 Behaviour of the optimal trajectories at junctions with $\partial\Phi$

There are two major points: what happens at an entry time, and what happens after an entry time. In this section we use the partitioning of  $\partial\Phi$  as proposed in [19, 20]. More precisely, we study the *qualitative* behaviour of optimal solutions which attain  $\partial\Phi$ , and possibly leave it or remain on it. The control input  $u(\cdot)$  in (2) is assumed to be piecewise smooth (a piecewise  $C^\infty([0, T_1], \mathbb{R})$  time function). Roughly, the process consists of subdivising  $\partial\Phi$  into subsets and especially trying to select that portion of  $\partial\Phi$  such that no input  $u(\cdot)$  exists that can keep an entry state  $x(\tau)$  inside  $\Phi$ . Then surely an optimal trajectory which possesses a boundary arc, cannot pass through this  $x(\tau)$ .

### 4.1 The admissible domain boundary partitioning

The following assumptions are in order:

- i) The pair  $(A, B)$  is controllable,

- ii)  $\text{Im}(B) \subseteq \text{Ker}(C)$ ,
- iii)  $D = 0$ ,  $C \neq 0$  and  $m = 1$ ,
- iv)  $\Phi \neq \emptyset$ .

One sees that ii) implies that  $CB = 0$ , i.e.  $r^{wu} \geq 2$ . The case  $r^{wu} = 1$  will be treated but deserves special attention. We restrict ourselves to the case  $D = 0$ , because as shown in [19, §VII] it is always possible to transform the dynamics to recover this case. The subsets of  $\partial\Phi$  are invariant under linear state feedback. Three sets are of interest here:  $\chi_{con}$  the subset of  $\partial\Phi$  at which trajectories coming from  $\text{Int}(\Phi)$  can make contact with  $\partial\Phi$ ,  $\chi_{rel}$  the subset of  $\partial\Phi$  at which trajectories starting on  $\partial\Phi$  can leave  $\partial\Phi$  and stay on a positive time interval in  $\Phi$ , and  $\mathcal{V}^*$  the subset of  $\partial\Phi$  such that there exists one  $u(\cdot) = u^*(\cdot)$  that can keep (locally) the trajectory on  $\partial\Phi$  ( $\mathcal{V}^*$  is the largest controlled invariant subspace in  $\partial\Phi$ ). Trajectories that make contact with  $\partial\Phi$  in  $\chi_{con} \setminus \mathcal{V}^*$ , leave  $\Phi$ . Let

$$r(x, u) : \partial\Phi \times \mathbf{U} \rightarrow \mathbb{N} \cup \{+\infty\}$$

and

$$r(x, u) = \min\{i \in \mathbb{N} \mid w^{(i)}(x, u) \neq 0\}$$

with  $r(x, u) = +\infty$  if  $w^{(i)}(x, u) = 0$  for all  $i \in \mathbb{N}$ , where  $\mathbf{U}$  is the set of piecewise smooth (infinitely differentiable) inputs.

The results of [19] hold for a piecewise smooth  $u(\cdot)$ . According to the framework developed in section 3, we are looking for a solution that is regular in the sense of definition 1. Thus  $x(\cdot)$  is analytic on  $[0, T_1] \setminus \{t_k\}$  and so is the optimal control  $u(\cdot)$ . Moreover the times at which the solution is not analytic are the jump times  $t_k \in \text{supp}(\lambda)$  and are entry times. But left accumulations of state  $\tilde{z}$  jumps may occur. As a consequence, to derive most of the results on the qualitative behaviour of optimal trajectories at an entry time  $\tau$ , one has to make the additional assumption that there is no left accumulation at  $\tau$ . In view of Proposition 7, this is equivalent to assuming that the set of admissible controls is restricted to piecewise analytic functions  $\{u\}(\cdot)$  (i.e. the support of the atomic part of  $\lambda$  is finite) and that reachability holds within this set. See also Proposition 13 for some more details on this point. Such an assumption is made for instance in [34, Theorem 1] (in a different context, though), where junction states are then called analytic junctions. Finally most of the results stated below hold for  $m \geq 2$  but assuming that there is only one constraint that is active at the considered time.

**Lemma 5** [19] *The following holds.*

- $\chi_{con} = \{x \in \partial\Phi \mid \exists u \in \mathbf{U} \text{ such that } \{r(x, u) < +\infty \text{ and even, and } w^{(r(x, u))} > 0\}, \text{ or } \{r(x, u) < +\infty \text{ and odd, and } w^{(r(x, u))} < 0\}$ ,
- $\chi_{rel} = \{x \in \partial\Phi \mid \exists u \in \mathbf{U} \text{ such that } r(x, u) < +\infty \text{ and } w^{(r(x, u))} > 0\}$
- $\mathcal{V}^* = \{x \in \partial\Phi \mid \exists u \in \mathbf{U} \text{ such that } r(x, u) = +\infty\}$ .

With some abuse of notation  $u$  may mean that the involved functions depend also on the derivatives of  $u$ . It is noteworthy that as long as  $u(\cdot)$  is smooth in right and left neighborhoods of the considered times of contact with  $\partial\Phi$  and (5) holds, then we can consider the three above sets calculated for (2) or for (4) (a) (c) as being the same sets. It is therefore not worth using different notations for the states of the  $z$ -dynamics and of the  $\tilde{z}$ -dynamics. The following result is obvious from the definition of the three subsets of  $\partial\Phi$ .

**Proposition 9** *Boundary arcs of optimal trajectories exist only in the subset  $\mathcal{V}^*$  of  $\partial\Phi$ . Contact states exist only in  $\chi_{con}$ , and exit states exist only in  $\chi_{rel}$ .*

This shows that an entry time may exist in  $\chi_{con} \setminus \mathcal{V}^*$  only if  $u(\cdot)$  is not smooth. As the following example shows it may even be necessary to apply a distributional input.

**Example 4** Let us choose  $A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ ,  $B = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ ,  $C = (1 \ 0)$ . Then  $\chi_{con} = \{x|x_1 = 0, x_2 \leq 0\}$ ,  $\chi_{rel} = \{x|x_1 = 0, x_2 \geq 0\}$ , and  $\mathcal{V}^* = \{x|x_1 = x_2 = 0\}$ . If  $u(\cdot)$  is restricted to be a function, all trajectories entering  $\partial\Phi$  in  $\chi_{con} \setminus \mathcal{V}^*$  leave  $\Phi$ . All trajectories initialised in  $\chi_{rel} \setminus \mathcal{V}^*$  enter  $\text{Int}(\Phi)$ . Thus boundary arcs are restricted to  $\mathcal{V}^*$  and  $u^* = 0$ . If an entry time  $\tau$  exists in  $\chi_{con} \setminus \mathcal{V}^*$ , then it is necessary to apply a Dirac measure  $u_\tau$  to keep the state in  $\Phi$  and make it jump in  $\chi_{rel}$ .

Let us notice that the case  $x_2(\cdot) \geq 0$  does not satisfy assumption ii). In such a case  $\mathcal{V}^* = \partial\Phi$  [20, Lemma 6.7.1] [19, lemma VII;1].

Algorithms exist [19, Algorithm A.8] [20, Algorithm A.1.6] which allow one to compute these three subsets in a finite number of steps. For instance the sequence  $\mathcal{V}^{k+1} = \text{Ker}(C) \cap A^{-1}(\mathcal{V}^k + \text{Im}(B)) = \{x \in \mathcal{V}^k \mid CA^k x = 0\}$  converges towards  $\mathcal{V}^*$  in at most  $r^{wu}$  steps<sup>(5)</sup>. This is extremely important in view of the development of a numerical analysis of the problem of interest here.

## 4.2 Entry and contact times and states

In this subsection we essentially focus on the qualitative analysis of optimal trajectories and control at entry and contact times, when  $m = n_u = 1$ . We denote  $u^*$  the control inside  $\mathcal{V}^*$ .

**Proposition 10** Let assumptions i)-iv) stand. (i) Let  $r^{wu} = n$ . Then  $\mathcal{V}^* = \{0\}$  and  $u^*(\cdot) \in \text{Ker}(B)$ . If the optimal trajectory has a boundary arc on  $(\tau, \tau + \epsilon)$ ,  $\epsilon \geq 0$ , then  $u^*(t) = 0$  on  $(\tau, \tau + \epsilon)$ . (ii) If  $r^{wu} < n$  then  $\mathcal{V}^* = \{x^*(t) = W^{-1} \begin{pmatrix} 0^{r^{wu}} \\ \xi(t) \end{pmatrix}\} = \{z^*(t) \mid \bar{z}^* = 0\}$ ,  $\dot{\xi}(t) = A_\xi \xi(t)$ , and  $u^*(t) = -(CA^{r^{wu}-1}B)^{-1}CA^{r^{wu}}W^{-1} \begin{pmatrix} 0^{r^{wu}} \\ \xi(t) \end{pmatrix}$ .

**Proof:** (i) From the definition of  $\mathcal{V}^*$ , given a state trajectory  $x(t)$  such that  $Cx(t) = 0$  on a positive time interval, one has to look for a control  $u^*(\cdot)$  that makes  $w^{(i)}(\cdot)$  equal to 0 for all integers  $i \geq 1$ . On  $(\tau, \tau + \epsilon)$  one has  $z_1^{(i)} = w^{(i)} = 0$ . Since  $r^{wu} = n$ ,  $\bar{z} = z$  and thus  $x^* = 0$  since  $z = Wx$  with  $W$  full rank. Since  $w^{(i)}(t) = CA^i x(t) + \sum_{j=1}^i CA^{j-1} B u^{*(i-j)}(t)$ , for all  $i \geq 1$ , we deduce that on  $(\tau, \tau + \epsilon)$  one has  $0 = \sum_{j=1}^i CA^{j-1} B u^{*(i-j)}(t)$  for all  $i \geq 1$ , from which one deduces that  $\frac{d^i u^*}{dt^i}(\cdot) \in \text{Ker}(B)$  for all  $i \geq 0$ . From Proposition 9 and since  $R > 0$ , the optimal control is obviously  $u^* = 0$  on  $[\tau, \tau + \epsilon)$  since otherwise the integral action  $I(u)$  strictly increases. (ii) is proved from the  $z$ -dynamics and the fact that  $\mathcal{V}^* = \lim_{k \rightarrow r^{wu}} \mathcal{V}^k$ , with  $\mathcal{V}^{k+1} = \{x \in \mathcal{V}^k \mid CA^k x = 0\}$  [19]. ■

When assumption ii) is not satisfied,  $r^{wu} = 1$  and it is obvious that  $\mathcal{V}^* = \{z \mid z_1 = 0\} = \partial\Phi$ . Point (ii) of Proposition 10 holds for any  $1 \leq r^{wu} \leq n$ .

**Corollary 5** Let  $r^{wu} = n$  and the above assumptions hold. If the unconstrained optimal trajectory  $x_{un}(\cdot)$  (the solution of (4) (a) (b) with  $C = 0$ ) satisfies  $x_{un}(t) \in \text{Int}(\Phi)$  for all  $t \in [0, T_1] \setminus \{t_k\}_{1 \leq k \leq l}$ ,  $l < +\infty$ , and  $x_{un}(t_k) \in \mathcal{V}^*$ , then the solution  $x(\cdot)$  of (4) is equal to  $x_{un}(\cdot)$  if  $u(\cdot)$  is restricted to be a smooth time function on  $[0, T_1]$ .

**Proof:** We know from Proposition 9 that boundary arcs and contact states belong to  $\mathcal{V}^*$ , and from Proposition 10 that on boundary arcs and contact states  $x = 0$  and  $u = 0$ . Any other state trajectory linking  $\bar{x}_0$  to  $\bar{x}_1$  in  $\Phi$  with a smooth input  $u'$ , will yield an integral action  $I(u') > I(u)$ . ■

If  $D \neq 0$ , one may recover the previous case by a simple transformation, however the new pair  $(\bar{A}, \bar{B})$  is never controllable. The subset  $\mathcal{V}^*$  can nevertheless be computed [20, §6.7] [19, §VII], as the following example shows.

<sup>5</sup>  $A^{-1}\mathcal{V} = \{x \in \mathbb{R}^n \mid Ax \in \mathcal{V}\}$ . It is not required that  $A$  be full rank.

**Example 5** Let  $A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ ,  $B = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ , and  $w(t) = -x_1(t) + 1$ . Then  $\mathcal{V}^* = \{(1, 0)\}$ , and  $u^* = 0$ .

**Proposition 11** If assumption ii) does not hold, i.e.  $r^{wu} = 1$ , then a (possibly discontinuous) time function  $u(\cdot)$  is sufficient to keep  $x(\cdot)$  inside  $\Phi$  at entry times. At such instants  $2 \geq \text{deg}(\lambda)$ .

**Proof:** Since  $\mathcal{V}^* = \partial\Phi = \{z_1 = 0\}$ , it is obvious from the  $z$ -dynamics of  $(A, B, C)$  that a discontinuous  $u(\cdot)$  is sufficient to keep  $x(\cdot)$  in  $\Phi$ . From (14), which is the  $\tilde{z}$ -dynamics of  $(\tilde{A}, \tilde{B}, \tilde{C})$ , we deduce that  $\lambda$  is a Dirac measure. Indeed  $r = 2$  and  $z_1(\cdot) = \{z_1\}(\cdot)$  is a function, and its derivative may be discontinuous at an entry time  $\tau$ . Thus  $Dz_2$  is a Dirac measure at  $\tau$ , and so is  $\lambda$ . We note also that  $x(\cdot)$  is continuous at  $\tau$ , so that  $\eta(\cdot)$  (which is a function) ‘‘concentrates’’ all the state  $\tilde{x}$  jumps. ■

Proposition 11 will be refined when we examine what happens along optimal trajectories for odd- $r^{wu}$  systems, see corollary 10. Proposition 11 suggests that if assumption ii) is true (i.e.  $r^{wu} \geq 2$ ) then at entry times in  $\chi_{con} \setminus \mathcal{V}^*$  a distributional input with degree  $\geq 2$  will be needed to keep  $x(\cdot)$  inside  $\Phi$ . This is a bit more subtle. The following results show that trajectories attain  $\partial\Phi$  in a specific way.

The sets inside  $\chi_{con} \setminus \mathcal{V}^*$  such that the first nonzero derivative of  $w(\cdot) = z_1(\cdot)$  is of order  $2k + 1$  are denoted as

$$\chi_{con}^{2k+1} = \{z \in \mathbb{R}^n \mid z_1^{(i)} = 0 \text{ for all } 0 \leq i \leq 2k, w^{(2k+1)} = z_1^{(2k+1)} < 0\}.$$

In the same way

$$\chi_{con}^{2k} = \{z \in \mathbb{R}^n \mid z_1^{(i)} = 0 \text{ for all } 0 \leq i \leq 2k - 1, w^{(2k)} = z_1^{(2k)} > 0\}.$$

and

$$\chi_{rel}^{2k+1} = \{z \in \mathbb{R}^n \mid z_1^{(i)} = 0 \text{ for all } 0 \leq i \leq 2k, z_1^{(2k+1)} > 0\}.$$

From  $z_1^{(i)}(t) = CA^i x(t) + \sum_{k=r^{wu}-1}^{i-1} CA^k Bu^{(i-1-k)}(t)$  one sees that these sets may depend on  $u(\cdot)$  and its derivatives. A trajectory which attains  $\partial\Phi$  in  $\chi_{con}^{2k}$  detaches immediately from  $\partial\Phi$  if the control  $u(\cdot)$  is smooth, since the first nonzero derivative will keep its positive sign in a right-neighborhood of the contact time. It is noteworthy that a trajectory which comes from  $\text{Int}(\Phi)$  and enters  $\chi_{con}^{2k+1}$  at  $\tau$  satisfies  $z_1^{(2j)}(\tau^-) \geq 0$  and  $z_1^{(2j-1)}(\tau^-) \leq 0$  for all  $0 \leq j \leq k$ . There is a sign inversion at each differentiation order due to the left analyticity in the neighborhood of  $\tau$  and the fact that all derivatives are zero at  $\tau$  up to the order  $2k$ . An illustrative example with  $r^{wu} = 3$  can be found in [11, Example 1.8]. A trajectory that leaves  $\partial\Phi$  at  $\tau$  after a boundary arc has to do it with a non-analytic  $z_1(\cdot)$  (otherwise by a simple Taylor expansion it follows that the trajectory remains in  $\mathcal{V}^*$ ) and such that the first non-zero derivative of  $z_1(\cdot)$  is positive, i.e. in a set  $\chi_{rel}^{2k+1}$ .

**Corollary 6** Let a contact time  $\tau \in [0, T_1]$  be such that  $x(\tau) \in \chi_{con}^{2k+1}$  and  $x(\cdot) \in \Phi$  in a right neighborhood of  $\tau$ . Then if  $2k + 2 \leq r^{wu}$ ,  $k \geq 0$ ,  $u$  has to be a distribution of degree  $2 + r^{wu} - 2k$  with an atom at  $\tau$ . If  $2k + 2 > r^{wu}$  then  $u(\cdot)$  is a function of time in a neighborhood of  $\tau$ , with its first nonzero derivative that is of order  $2k + 1 - r^{wu}$  and is discontinuous at  $\tau$ .

**Proof:** Follows from the decomposition of  $\chi_{con}$  and the  $z$ -dynamics, from which  $z_1^{(r^{wu}-1)}(\cdot) = z_{r^{wu}}(\cdot)$  depends explicitly on  $u(\cdot)$ . In particular  $z_1^{(2k+1)}(\cdot) = \dot{z}_{2k+1}(\cdot) = z_{2k+2}(\cdot)$  and is a function of  $u^{(2k+1-r^{wu})}(\cdot)$ , which has then to be discontinuous to keep the trajectory inside  $\Phi$ . ■

The result of Corollary 6 will be refined next depending on  $r^{wu}$ . In most of the studies (see e.g. [25, 28, 8, 13]) the constraint that the contact with  $\partial\Phi$  occurs ‘‘tangentially’’ is added. This means simply that at a time  $\tau$  at which the optimal trajectory attains  $\partial\Phi$ , then  $\bar{z}(\tau^-) = 0$ , where we take the left limit as it could be that jumps occur (see Proposition 6).

**Definition 3** A contact state at time  $\tau \in [0, T_1]$  is said tangential if  $\bar{z}(\tau^-) = 0$ , and hypertangential if  $\bar{\bar{z}}(\tau^-) = 0$ .

The meaning of the overbar is given after (14). If the contact is hypertangential, it follows from (30) (31) that  $\tilde{z}(\tau^+) = 0$  (which does not mean that  $\tilde{z}(\tau^+) \in \mathcal{V}^*$ ). An hypertangential contact necessarily occurs in a set  $\Xi_{con}^i$  with  $i \geq r$ .

**Proposition 12** *Let  $r^{wu} = 2k$ ,  $k \geq 1$ ,  $\deg(\lambda) \leq 2$ , and  $i \geq 0$ . Then all optimal trajectories which make contact with  $\partial\Phi$  at  $t = \tau$  in the set  $\chi_{con}^{r^{wu}+2i}$  are touch states if  $u^{(2i)}(\cdot)$  is a continuous function at  $t = \tau$ . If contact occurs in  $\chi_{con}^{r^{wu}+2i+1}$ , then necessarily  $i \geq k - 1$  and necessarily  $u^{(2i+1)}(\cdot)$  is discontinuous at  $\tau$ .*

**Proof:** If contact occurs in  $\chi_{con}^{r^{wu}+2i}$ , then  $z_1^{(r^{wu}+2i)} > 0$ . If  $u^{(2i)}(\cdot)$  is a continuous function at  $t = \tau$ , so are all  $u^{(j)}(\cdot)$  with  $0 \leq j \leq 2i - 1$ . Therefore from  $z_1^{(r^{wu}+2i)} = CA^{r^{wu}+2i}BW^{-1}z + \sum_{k=0}^{2i} CA^{r^{wu}-1-k+2i}Bu^{(k)}$ , we deduce that  $z_1^{(r^{wu}+2i)}(\cdot)$  is a continuous function at  $\tau$ . From the definition of  $\chi_{con}^{r^{wu}+2i}$  it results that in an open right neighborhood of  $\tau$ ,  $z_1(t) > 0$ . If contact occurs in  $\chi_{con}^{r^{wu}+2i+1}$ , then  $z_1^{(r^{wu}+2i+1)} < 0$ . If  $u^{(j)}(\cdot)$ ,  $0 \leq j \leq 2i + 1$ , are continuous at  $\tau$ , then in an open right neighborhood of  $\tau$ ,  $z_1(t) < 0$ . Also  $\deg(\lambda) \leq 2$  implies that  $\deg(u^{(r^{wu})}) \leq 2$ . If  $i \leq k - 2$ , then the trajectory escapes from  $\Phi$  in an open right neighborhood of  $\tau$ , because the lower order derivatives of  $u(\cdot)$  must be continuous at  $\tau$ , and so is  $z_1^{(r^{wu}+2i+1)}(\cdot)$ . If  $i \geq k - 1$ , the degree condition on  $\lambda$  can be respected even if  $u^{(2i+1)}(\cdot)$  is discontinuous at  $\tau$ . Moreover if  $u^{(2i+1)}(\cdot)$  is continuous at  $\tau$  then the trajectory escapes from  $\Phi$  in an open right neighborhood of  $\tau$ . ■

It is noteworthy that  $u^{(2i+1)}(\cdot)$  be discontinuous at  $\tau$  when  $i \geq k - 1$  is not only a necessary condition, but this is also sufficient. Indeed the conditions secure that (4) does represent the optimal conditions for (1) (2), and the positivity of  $M^{(2r^{wu})}$  (Lemma 1) assures that Lemma 3 holds. The case  $r^{wu} = 2k + 1$  can be analysed similarly. However it requires more care as Propositions 13 and Corollary 7 show. The case  $r^{wu} = 1$  is treated in Proposition 11 (since then  $\mathcal{V}^* = \partial\Phi$  we get that  $\chi_{con} \setminus \mathcal{V}^* = \emptyset$  so that the condition of Proposition 12 degenerates).

The next step is to study optimal trajectories which make contact with  $\partial\Phi$  in  $\chi_{con} \setminus \mathcal{V}^*$ , and then are sent into  $\mathcal{V}^*$  (where possibly a boundary arc exists).

**Proposition 13** *Let  $m = n_u = 1$ . Let an optimal trajectory make contact with  $\partial\Phi$  at  $t = \tau$ , in the subset  $\chi_{con}^{2r^{wu}-1}$ , and be such that  $z(\tau^+) \in \mathcal{V}^*$ . Then necessarily  $r^{wu}$  is even. The same holds if  $z(\tau^+) \in \chi_{rel}^{2r^{wu}-1}$ .*

**Proof:** From the  $\tilde{z}$ -dynamics and Lemma 1, it follows that if  $\lambda$  is a measure at a time  $t_k$ , then  $d\mathcal{J}_r$  has an atom  $\lambda_\tau = (-1)^{r^{wu}}(CA^{r^{wu}-1}B)^{-2}(z_r(\tau^+) - z_r(\tau^-))\delta_\tau$ . If  $z(\tau^+) \in \mathcal{V}^*$ , then  $z_r(\tau^+) = 0$ . Now since contact is made in  $\chi_{con}^{2r^{wu}-1}$  we have that  $z_1^{(r-1)}(\tau^-) (= z_r(\tau^-)) < 0$ , whereas  $z(\tau^+) \in \mathcal{V}^*$  implies that  $z_1^{(r-1)}(\tau^+) = 0$ . Therefore if  $r^{wu} = 2k + 1$  for some  $k \geq 0$ ,  $\lambda$  cannot be a positive measure. The last point can be proved similarly since then  $z_r(\tau^+) > 0$ . The proof is complete. ■

Let us notice that under the conditions of Proposition 13, a jump in  $z_r(\cdot)$  is equivalent to a jump in  $u^{(r^{wu}-1)}(\cdot)$ , since the derivatives of smaller order of  $u(\cdot)$  are time-continuous functions. Proposition 13 says that under the stated conditions (which in particular imply that  $\deg(\lambda) = 2$ ), only systems with even relative degree possess entry states, which is in accordance with the result of [28]. The second point of Proposition 13 implies in particular that odd- $r^{wu}$  systems do not even possess touch points with  $z(\tau^-) \in \chi_{con}^{2r^{wu}-1}$  and  $z(\tau^+) \in \chi_{rel}^{2r^{wu}-1}$ . Consequently when  $\deg(\lambda) = 2$  it is excluded to get an entry time that is the accumulation of touch points. Proposition 13 concerns conditions under which  $\lambda$  is a singular measure. If the contact occurs in a set  $\chi_{con}^i$  with  $i \geq r$ , then  $\tilde{z}(\tau^-) = 0$  and the problem of existence of a boundary arc is equivalent to the well-posedness of a LCP (see section 5). The next Corollary is a consequence of Proposition 13.

**Corollary 7** *Let  $m = 1$ ,  $\lambda$  be a measure and  $\{\eta\}(\tau^+) - \{\eta\}(\tau^-) = -C^T\lambda_1$ . If  $\lambda_1 > 0$ , then  $r^{wu}$  is an even integer.*

Compiling Corollary 7 and Proposition 11, one deduces that if  $r^{wu} = 1$  and if  $u(\cdot)$  has a discontinuity at an entry time  $\tau$ , then necessarily  $\lambda_1 = 0$ , i.e. the costate  $\eta(\cdot)$  is continuous at  $\tau$ . Hence we retrieve

the result of [8, Lemma 2.4 (2)]. Also from Proposition 4 and Corollaries 2 and 7, it follows that if  $(A^T)^j \sigma_\eta(\tau) \in \text{Ker}(B^T)$  for  $0 \leq j \leq r^{wu} - 2$  and  $u^{r^{wu}-1}(\cdot)$  jumps at  $\tau$ , then  $r^{wu}$  is even. Notice that if  $\deg(\lambda) = 2$ , then contact cannot occur in a set  $\chi_{con}^{2k}$  with  $k < r^{wu}$ , since  $r$  is even. If  $\tilde{z}(\tau^-) \in \chi_{con}^{2k}$ ,  $k < r^{wu}$ , and  $\tilde{z}(\tau^+) \in \mathcal{V}^*$ , then  $\deg(\lambda) \geq 3$ .

**Proposition 14** *Let  $m = n_u = 1$ , and  $r^{wu} = 2k + 1$ ,  $k \geq 0$ . Then  $\tilde{z}(\tau^-)$  is an entry state, only if  $\deg(\lambda) \leq 1$  or  $\deg(\lambda) \geq 3$ .*

**Proof:** One has  $\tilde{z}(\tau^+) \in \mathcal{V}^*$ , hence  $z_1^{(i)}(\tau^+) = 0$  for all  $i \geq 0$ . The fact that  $\tilde{z}(\tau^-)$  is an entry state means that  $\tilde{z}(\tau^-) \in \chi_{con}$ , so that  $z_1(\tau) = 0$ , and the first nonzero  $z_1^{(i)}(\tau^-)$  is either  $> 0$  or  $< 0$ . The higher order derivatives satisfy  $z_1^{(2k)}(\tau^-) \geq 0$  and  $z_1^{(2k+1)}(\tau^-) \leq 0$ . It follows that at  $t = \tau$  the measures  $d\mathcal{J}_i$  in (32) satisfy  $d\mathcal{J}_{2k} \geq 0$  and  $d\mathcal{J}_{2k+1} = 0$ . Proposition 13 says that if  $\deg(\lambda) = 2$  at  $\tau$  (i.e.  $d\mathcal{J}_r > 0$ ), then  $r^{wu}$  is even. A first situation is when  $z_1^{(i)}(\tau^-) \neq 0$  for some  $i < r - 1$ . Depending on the values of the derivatives it is possible that  $d\mathcal{J}_{2k} > 0$  for  $k < r^{wu}$  and that  $d\mathcal{J}_r = 0$  and  $\deg(\lambda) \geq 3$ . A second situation is when the contact is hypertangential (see definition 3), so that  $d\mathcal{J}_i = 0$  for all  $1 \leq i \leq r$  (in this last case the trajectory may either detach from the constraint, or stay on it with  $g_r(t) > 0$  on some nonzero time interval, where  $g_r(\cdot)$  is in (32)). Then  $\deg(\lambda) \leq 1$ , i.e.  $\lambda$  is a function. ■

As an example, let us consider  $r^{wu} = 3$ . If  $z_1(\tau^-) = 0$ ,  $z_2(\tau^-) < 0$ , and  $z_3(\tau^-) = z_3(\tau^+) = z_3(\tau^-) > 0$ , i.e.  $\tilde{z}(\tau^-) \in \chi_{con}^0$ , then  $d\mathcal{J}_2 = z_2(\tau^+) - z_2(\tau^-) = -z_2(\tau^-) > 0$  while all other  $d\mathcal{J}_i$  have no atom at  $\tau$ . Thus  $\deg(\lambda) = 6$ . Another type of approach is detailed in example 2, where this time  $\deg(\lambda) = 3$ . It immediately follows from Proposition 14 that if  $\lambda$  is restricted to be a measure, then systems with  $r^{wu} = 2k + 1$  possess optimal trajectories such that  $\tilde{z}$  and consequently  $\tilde{x}$  are time continuous, because  $\deg(\lambda) \leq 1$ . It is important to recall that there may exist trajectories of (4) with  $x(\cdot)$  and  $B^T \eta(\cdot)$  continuous functions, and with  $\deg(\lambda) \geq 3$ , see Proposition 6 and example 2. Once  $\tilde{z}(\tau^+) \in \mathcal{V}^*$ , then the existence of a boundary arc relies upon the well-posedness of a LCP whose solution is  $\lambda(t)$ ,  $t > \tau$ , see section 5. But if  $r^{wu} = 2k$ ,  $k \geq 1$ , nothing hampers an optimal trajectory to possess left accumulations of state jumps when  $\lambda$  is a measure. In any case, left accumulations at an entry state can exist in the state variables  $\{z_i\}(\cdot)$  with  $i \leq r - 1$ .

**Proposition 15** *Let  $r^{wu} = 1$ ,  $n_u = m = 1$ . Assume that  $\bar{x}_1$  is reachable from  $\bar{x}_0$  with  $\mathbf{U} = \mathbf{U}_{\mathcal{J}}^0$ . Then the optimal trajectory and controller are time-continuous functions on  $[0, T_1]$ .*

**Proof:** The reachability assumption is fundamental, as example 6 shows. Since it holds one can search for optimal inputs which are functions, see (15). From Proposition 14 it follows that  $\deg(\lambda)$  is either  $\leq 1$  or  $\geq 3$ . From (14) it follows that  $\deg(\lambda) \geq 3 \Rightarrow \deg(z_2) \geq 2$ . This is in contradiction with the choice of admissible inputs  $\mathbf{U}$ . Thus  $\deg(\lambda) \leq 1$  so that both  $x(\cdot)$  and  $\eta(\cdot)$  are time-continuous on  $[0, T_1]$ . ■

This result is consistent with the fact that the HOSP solutions are in  $\mathcal{T}_1([0, T_1])$  for  $r = 2$ , which implies that  $\deg(\tilde{z}) \leq 2$ , hence  $\deg(\lambda) \leq 3$ . As example 6 demonstrates, a necessary condition for the framework that is developed in this paper around the HOSP, is that the zero dynamics  $\dot{\xi}(t) = A_\xi \xi(t) + B_\xi z_1(t)$  be reachable from  $\bar{\xi}_0$  to  $\bar{\xi}_1$  with  $z_1$  a function of time.

In [28, 25] it is stated that if at a contact time  $u^{r^{wu}-1}(\tau^+) \neq u^{r^{wu}-1}(\tau^-)$  with all the smaller order derivatives continuous, then odd-relative degree systems do not have a boundary arc. This optimal input jump condition implies that  $\deg(\lambda) = 2$ , as can be easily seen from the  $\tilde{z}$ -dynamics. From Proposition 14 it follows that  $\deg(\lambda) \leq 1$ . However the conclusion that there is no boundary arc is premature, see section 5.

**Remark 7** *The above results can be extended to nonlinear systems as ten Dam's framework extends nonlinear systems which possess a relative degree  $\begin{cases} \dot{x}(t) = f(x(t)) + g(x(t))u(t) \\ w(t) = h(x(t)) \end{cases}$  [20], with  $f(\cdot)$ ,  $g(\cdot)$ ,  $h(\cdot)$  smooth functions of  $x$ . Once the HOSP is extended to this nonlinear case, then all the material in this paper readily extends. These results are of the same nature as results in [25, 28, 34]. For instance compiling Corollary 6 and Proposition 13 one finds that the sum of the order of the first discontinuous*

derivative of  $u(\cdot)$  and of the relative degree  $r^{wu}$ , is always odd at an entry time when  $\deg(\lambda) \leq 2$  (see [33, Theorem 5.1, Corollary 5.2] for similar results). However our study is more general as it relies on a systematic and intrinsic to  $(A, B, C)$  (or  $(\tilde{A}, \tilde{B}, \tilde{C})$ ) partitioning of  $\partial\Phi$  as well as structural properties of the system (4). In contrast the condition which allows to prove a similar result as Proposition 13, namely [28, equ.(81)], does not involve the Markov parameter  $M^{(r)}$ , but is based on suitable Taylor expansions of  $\frac{\partial H}{\partial u}$  and  $w(\cdot)$ . The partitioning of  $\partial\Phi$  also contains constraint qualifications of the form: there exist  $x$  and  $u$  such that  $(Ax + Bu)^T C^T < 0$  for all  $x \in \partial\Phi$  (which is nothing else but  $z_1 \in \chi_{con}^1$ ). Interestingly enough, the work in [25] already defined sets similar to those in Lemma 5, but the developments remained at an embryonic stage. It is sometimes argued that boundary (or constrained) arcs are impossible for odd  $r^{wu}$  [42]. Proposition 13, Corollary 7 and Proposition 14, and the material in section 5 show that this is more subtle. Notice that [34] study the problem with input constraints, so that the necessary conditions are not as in (4) since  $u(t) = -K \operatorname{sgn}(B^T \eta(t))$  for some  $K$ . The necessary conditions can still be written under a complementarity framework since the sign function lends itself to a complementarity formulation:

$$\begin{aligned} \begin{pmatrix} \dot{x}(t) \\ \dot{\eta}(t) \end{pmatrix} &= \begin{pmatrix} A & 0 \\ Q & -A^T \end{pmatrix} \begin{pmatrix} x(t) \\ \eta(t) \end{pmatrix} + \begin{pmatrix} 0 \\ -C^T \end{pmatrix} \lambda(t) + \begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \lambda_1(t) \\ \lambda_2(t) \end{pmatrix} \\ 0 &\leq \begin{pmatrix} w(t) = Cx(t) + D \\ w_1(t) \\ w_2(t) \end{pmatrix} \perp \begin{pmatrix} \lambda(t) \\ \lambda_1(t) \\ \lambda_2(t) \end{pmatrix} \geq 0 \\ \lambda_1(t) + \lambda_2(t) &= 1 \end{aligned} \tag{52}$$

where  $w_1(\cdot)$  and  $-w_2(\cdot)$  are the positive and negative parts of  $B^T \eta(\cdot)$ , respectively. System (52) possesses a non uniform vector relative degree since relay systems lend themselves to a description as relative degree 0 system [15]. The HOSP is presented in [1] for  $r \geq 1$ , and should be extended to  $r = 0$  in order to encompass systems as (52). It is finally noteworthy that when input constraints are considered, then the material in section 3.5 is meaningless, unless these constraints are considered almost everywhere except at contact or entry times.

The following concerns the complementarity conditions of Proposition 2 for an entry state.

**Corollary 8** *Let  $\tilde{z}(\tau^-) \in \chi_{con}^{2k+1}$  and  $\tilde{z}(\tau^+) \in \mathcal{V}^*$ ,  $k \leq \frac{r+1}{2}$ . Then  $0 \leq \{z_i\}(\tau^+) \perp d\nu_i(\{\tau\}) \geq 0$  for all  $1 \leq i \leq 2k+2$ . Moreover  $d\mathcal{J}_{2k+2} > 0$  if  $k \leq r^{wu} - 1$  and  $r^{wu}$  must be even if  $k = r^{wu} - 1$ .*

**Proof:** Using the tangent cones definition in section 3.1, the definition of the sets  $\chi_{con}^{2k+1}$  and  $\mathcal{V}^*$ , and the material of remark 4, it follows that  $T_{\Phi}^i(\{z_1\}(t^-), \dots, \{z_i\}(t^-)) = \mathbb{R}^+$  and  $\partial\psi_{T_{\Phi}^{i-1}(\{z_1\}(t^-), \dots, \{z_{i-1}\}(t^-))}(\{z_i\}(t^+)) = \mathbb{R}^-$  for all  $1 \leq i \leq 2k+2$  (recall that  $z_1^{(2k+1)}(\cdot) = z_{2k+2}(\cdot)$ ). The first result follows from (37). The second statement follows from (30) and (31), and using Lemma 1 and Proposition 13. ■

### 4.3 Exit times and states

Since the results which can be obtained for exit times are similar to those obtained for entry times, only one Proposition is given in this section.

**Proposition 16** *Let  $\deg(\lambda) \leq 2$ . Then the connection between a boundary arc and an interior arc has to occur at an exit state in  $\chi_{rel}^{2k+1}$  with  $k \geq r^{wu} - 1$ .*

We omit the proof which relies on the same arguments as proofs above. Similar results concerning costate and input regularity conditions can be derived for the exit times, and are not given here for the sake of brevity. The following result applies to a contact state that is not an entry state, but a grazing state, and is a direct consequence of Proposition 2.



**Corollary 9** *Let  $\tilde{z}(\tau^-) \in \chi_{con}^{2k+1}$  and  $\tilde{z}(\tau^+) \in \chi_{rel}^{2k+1}$ ,  $k \leq \frac{r+1}{2}$ . Then the following CP is satisfied:*

$$0 \leq \{z_{2k+1}\}(\tau^+) \perp -dv_{2k+1}(\{\tau\}) \in \partial\psi_{T_{\Phi}^{2k}(\{z_1\}(\tau^-), \dots, \{z_{2k}\}(\tau^-))}(\{z_{2k+1}\}(\tau^+)).$$

Corollaries 8, 9 and Proposition 2 show that the additional constraints which may be imposed at contact times, are taken care of with the multipliers  $dv_i$  which in turn satisfy a special set of complementarity conditions, not present in (4). It is worth noting that the complementarity conditions satisfied by the multipliers  $dv_i$  depend both on the pre and post-contact states. This generalizes Mechanics in the framework of which the percussion value (i.e. the magnitude of the Dirac measure at an impact) depends on both the pre and post-impact velocities. To the best of our knowledge this is introduced for the first time in the context of optimal control with state inequalities.

**Remark 8** *The introduction of the measure multipliers  $dv_i$ ,  $1 \leq i \leq r$ , generalizes the case  $r^{wu} = 1$  treated in [5, 30] where a multiplier associated to  $z_2(t)$  is introduced (denoted as  $\dot{\mu}(t)$  in these papers, and which corresponds in our framework to the function  $g_2(t)$ ). The conditions given for instance in [30, equ.(11b) (12b)] are a particular case of conditions in Proposition 2 and especially (34) with  $r = 2$ . The inclusion (34) also includes [31, equ.(25)] which is stated directly along boundary intervals  $(t_k, t_{k+1})$  and  $\{t_k\}$  is assumed to be a finite set in [31].*

## 4.4 Multivariable systems

### 4.4.1 Extension of Proposition 13

If  $m \geq 2$ , some care has to be taken. Indeed as noted in [19], contact can occur on a portion of  $\partial\Phi$  such that for instance  $x(\tau^-) \in \chi_{con,1}^{2k} \cap \mathcal{V}_2^*$ . Then the optimal trajectory has a contact time with  $\partial\Phi_1$  but may have an entry time with  $\partial\Phi_2$ . Let  $n_u = m$ . If  $x(\tau^-) \in \chi_{con,1}^{2k+1} \cap \chi_{con,2}^{2k+5}$  and  $2k+1 < r^{wu} < 2k+5$ , one optimal input will be distributional while the other one will be a function of time at  $\tau$ . It may also happen that  $\chi_{rel,j}^i \notin \partial\Phi$  for some constraint  $j$  and some  $i \geq 1$  (see e.g. [20, Example 6.6.1]). An optimal trajectory with a boundary arc on the constraint  $j$  will possess an exit state with an optimal controller with reduced regularity, depending on  $i$ . Thus the qualitative behaviour of the optimal trajectory that is made above is still possible, but is more involved. It is however noteworthy that the HOSP state reinitialization mapping in (30) (31) still works, and brings an answer to the questioning about the collision map definition in [19, §VI]. In the multivariable case ( $m \geq 2$ , uniform relative degree  $r$ , definite Markov parameter  $M^{(r)}$ ), the conclusions of Proposition 13 need to be refined. For instance, if contact occurs in  $\chi_{con}^{2r^{wu}-1} = \cup_{i=1}^m \chi_{con,i}^{2r_i^{wu}-1}$ , which means that all  $z_{r,i}(\tau^-) < 0$ , and if  $z(\tau^+) \in \mathcal{V}^* = \cap_{i=1}^m \mathcal{V}_i^*$ , then  $r^{wu}$  does not necessarily have to be even. It has to be even if the leading Markov parameter  $M^{(r)}$  is diagonal. But in general if  $M^{(r)}$  is negative definite (without being diagonal), then the reasoning in the proof of Proposition 13 may fail as  $M^{(r)}(-z_r(\tau^-))$  may have positive components despite  $M^{(r)} < 0$  and all components of  $-z_r(\tau^-)$  are positive. Let us recall that a  $m \times m$  strictly semimonotone matrix [18, Definition 3.3.9] (or **E**-matrix) satisfies  $[v \in (\mathbb{R}^+)^m, v \neq 0] \Rightarrow [v_i > 0 \text{ and } (M^{(r)}v)_i > 0 \text{ for some } i]$ . This allows us to propose an extension of Proposition 13.

**Proposition 17** *Let  $m \geq 2$  and the triple  $(\tilde{A}, \tilde{B}, \tilde{C})$  have a uniform relative degree  $r$ . Let  $M^{(r)}$  be an **E**-matrix. Then it is possible that contact occurs in the set  $\chi_{con}^{2r^{wu}-1}$ , and with  $z(\tau^+) \in \mathcal{V}^*$ . If  $-M^{(r)}$  is an **E**-matrix, this is impossible.*

**Proof:** From [18, Theorem 3.10.7 (b) (f)], we get that if  $M^{(r)}$  is an **E**-matrix, then  $[v \geq 0] \Rightarrow [M^{(r)}v \geq 0]$ . The rest of the proof is similar to the proof of Proposition 13. ■

Proposition 17 generalizes Proposition 13, since when  $M^{(r)}$  is a scalar, semi strict monotonicity simply means positivity. Many other generalizations would be possible, like the one indicated above when the leading Markov parameter is diagonal. We do not tackle them here for the sake of brevity.

#### 4.4.2 Continuous dependence of BVP solutions

It is well known from Mechanics that solutions of differential inclusions like the sweeping process, may be discontinuous with respect to initial conditions when  $m \geq 2$  [32, 44]. This creates a fundamental problem for optimal control with unilateral state constraints and  $m \geq 2$ . If one embeds (4) into the HOSP and proves that the solution is optimal with respect to an extended integral action, then the optimal pair  $(x(\cdot), u(\cdot))$  may be very sensitive to the BVP data. It is noteworthy that this phenomenon may already exist when  $\lambda$  is a measure, i.e. when the HOSP is a measure differential inclusion, and should therefore be considered as an important feature. We shall say that the BVP in (4) is continuous with respect to data if small perturbations on the data  $(\bar{x}_0, \bar{x}_1, T_1)$ , result in a small perturbation of the solution  $\tilde{z}$ . More formally, let the sequences  $\{\bar{x}_{0,n}\}$ ,  $\{\bar{x}_{1,n}\}$ , and  $\{T_{1,n}\}$  converge in  $\mathbb{R}$  towards  $\bar{x}_0$ ,  $\bar{x}_1$ , and  $T_1$  respectively. Let us consider the solutions of the BVP (4) as being elements in  $\mathcal{F}_\infty([0, T_1], \mathbb{R})$ , and solutions of the measure differential formalism in (22) (23)(24) (or (25) (26)). Let us denote them as  $\tilde{z}_n(\cdot)$  for the corresponding data. Then weak (resp. uniform) continuity in the data holds if and only if  $\{\tilde{z}_n(\cdot)\}_{n \geq 0}$  converges weakly (resp. uniformly) towards  $\tilde{z}(\cdot)$ . The analysis which follows is essentially qualitative and illustrates why continuity w.r.t. the data may fail.

Let us consider  $A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}$ ,  $B = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}$ ,  $C = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{pmatrix}$ ,  $R = I_2$ ,  $Q = I_4$ . Then

$r = (2 \ 2)^T$  and the leading Markov parameter of  $(\tilde{A}, \tilde{B}, \tilde{C})$  is  $M^{(4)} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} > 0$ . The  $\tilde{z}$ -dynamics is (in the  $\tilde{z}$  coordinates, one has  $z_1 = (w_1 \ w_2)^T = (x_1 \ x_1 + x_3)^T$ )

$$\begin{cases} \dot{z}_1^1(t) = z_2^1(t) \\ \dot{z}_2^1(t) = z_3^1(t) \\ \dot{z}_3^1(t) = z_4^1(t) \\ \dot{z}_4^1(t) = \eta_2(t) + x_1(t) + \lambda_1 + \lambda_2 \\ w_1(t) = z_1^1(t) \end{cases} \quad \begin{cases} \dot{z}_1^2(t) = z_2^2(t) \\ \dot{z}_2^2(t) = z_3^2(t) \\ \dot{z}_3^2(t) = z_4^2(t) \\ \dot{z}_4^2(t) = \eta_2(t) - x_3(t) + \eta_4(t) - x_1(t) + \lambda_1 + 2\lambda_2 \\ w_2(t) = z_1^2(t) \end{cases} \quad (53)$$

with  $z^1 = (z_1^1, z_2^1, z_3^1, z_4^1)^T$ ,  $z^2 = (z_1^2, z_2^2, z_3^2, z_4^2)^T$  (these two vectors are needed to define the tangent cone as a product, see remark 3),  $z_i = (z_i^1 \ z_i^2)^T$ ,  $1 \leq i \leq 4$ . Let us now consider a solution  $\tilde{x}(\cdot)$  of the BVP (4), with initial condition  $(\bar{x}_0, \eta(0^-))$  (point  $A$  on figure 1). Let us assume that  $\tilde{x}(\cdot)$  hits the boundary  $w_2 = 0$  at a time  $t_0 < T_1$  in a neighborhood of the origin (point  $B$  on figure 1), with the following data:  $(z_1^1(t_0^-), z_2^1(t_0^-)) = (\bar{z}_1^1, 0)$ ,  $\bar{z}_1^1 > 0$ ,  $\dot{w}_2(t_0^-) = z_2^2(t_0^-) = \dot{x}_1(t_0^-) + \dot{x}_3(t_0^-) < 0$ ,  $\dot{x}_1(t_0^-) < 0$ ,  $\dot{x}_3(t_0^-) < 0$  and  $\dot{x}_1(t_0^-) = \dot{x}_3(t_0^-)$  (a normal contact),  $z_2^2(t_0^-) < 0$ ,  $z_3^2(t_0^-) = 0$ ,  $z_4^2(t_0^-) < 0$ , whereas  $z_2^1(t_0^-) \in \mathbb{R}$ ,  $z_3^1(t_0^-) \in \mathbb{R}$ ,  $z_4^1(t_0^-) \in \mathbb{R}$ . This way we assume that the contact is made in the set  $\chi_{con,2}^1$ . The tangent cones can be computed as indicated in section 3.1 (see in particular remarks 3 and 4), and we find  $T_\Phi^0 = \Phi = \mathbb{R}^+ \times \mathbb{R}^+$ ,  $T_\Phi^1(z_1(t_0^-)) = T_\Phi(z_1(t_0^-)) = \mathbb{R} \times \mathbb{R}^+$ ,  $T_\Phi^2(z_1(t_0^-), z_2(t_0^-)) = T_{\mathbb{R}}(z_2^1(t_0^-)) \times T_{\mathbb{R}}(z_2^2(t_0^-)) = \mathbb{R} \times \mathbb{R}^+$ ,  $T_\Phi^3(z_1(t_0^-), z_2(t_0^-), z_3(t_0^-)) = \mathbb{R} \times \mathbb{R}^+$ . It follows from (30) and (31) that  $z_1^1(t_0^+) = z_1^1(t_0^-)$ ,  $z_1^2(t_0^+) = z_1^2(t_0^-)$ ,  $z_2^2(t_0^+) = 0$ ,  $z_2^1(t_0^+) = z_2^1(t_0^-)$ ,  $z_3^1(t_0^+) = z_3^1(t_0^-)$ ,  $z_3^2(t_0^+) = 0$ . Though  $M^{(4)}$  doesn't satisfy the condition in remark 5 i), there are two possibilities:  $z_4^1(t_0^+) \in \mathbb{R}$  and  $z_4^2(t_0^+) = 0$  or  $z_4^1(t_0^+) \in \mathbb{R}$  and  $z_4^2(t_0^+) > 0$ . Since it is assumed that  $B$  is in a neighborhood of the origin, and since the trajectory is analytic in a right neighborhood of  $t_0$ , there exists a junction time  $t_1 > t_0$  with contact either at the origin (the trajectory "slides" on the boundary  $w_2 = 0$ ) or with the boundary  $w_1 = 0$  at some point  $C$  that is also in a neighborhood of the origin. In any case the trajectory will either continue to evolve along the boundary  $w_1 = 0$ , or hit it and leave it with a large tangential velocity. If the optimal trajectory detaches at  $B$  and hits the boundary  $w_1 = 0$  at  $t_1$  at some point  $C$ , then the "velocity" is reinitialized to  $\dot{x}_1(t_1^+) = 0$  and  $\dot{x}_3(t_1^+) = \dot{x}_1(t_1^-) + \dot{x}_3(t_1^-)$ . Since we can assume that the trajectory evolves arbitrarily close to the corner at the origin, it follows that  $\dot{x}_3(t_1^+)$  has a magnitude that is of the same order as that of  $\dot{x}_1(t_0^+) = \dot{x}_1(t_0^-)$  (because  $z_2^1(t_0^+) = z_2^1(t_0^-)$ ). Therefore given any  $\epsilon > 0$ , it is possible to initialize the trajectory in  $A$  with a large enough  $|\dot{x}_1(t_0^-)|$ , so that at a time  $t_1 + \delta t_1 > t_1 > t_0$ , it enters the domain

$D_\epsilon = \{0 \leq x_1 \leq \epsilon, x_3 \geq \epsilon\}$ . The same reasoning for the trajectory starting at  $A + \delta A$  can be done to conclude that it enters  $\bar{D}_\epsilon = \{x_1 \geq \epsilon, 0 \leq x_1 + x_3 \leq \epsilon\}$  after a finite time. It is noteworthy that  $\delta \tilde{x}_0$  may be arbitrarily small while  $\epsilon$  in  $D_\epsilon$  and  $\bar{D}_\epsilon$  keeps its value. Whatever the behaviour at the corner may be, the continuous dependence on data as defined above cannot hold.

We conclude that provided the initial data satisfies some magnitude constraint, the unperturbed optimal trajectory starting at  $(\bar{x}_0, \eta(0))$  reaches an end-point  $(\bar{x}_1, \eta_1)$  in a time  $T_1$ , but the perturbed trajectory will not reach a neighborhood of  $(\bar{x}_1, \eta_1)$  in a time  $T_1 + \delta T_1$ . In order that it does, the BVP solution starting at  $\bar{x}_0 + \delta \bar{x}_0$  has to be initialized with a costate that is not a small perturbation of  $\eta(0)$ . Therefore the discontinuous dependence behaviour of the IVP solution, may create a jump in the initial condition of the BVP solution. Consequently, the optimal controller will also suffer from a ‘‘discontinuity’’.

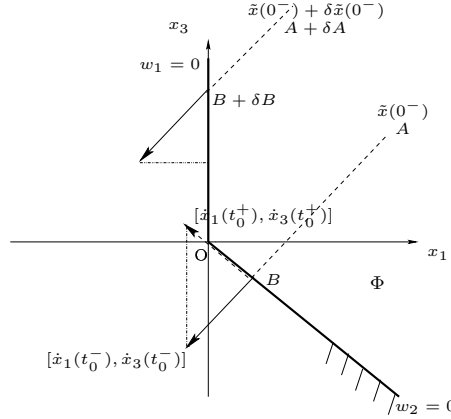


Figure 1: Sensitivity of the BVP solutions w.r.t. initial data  $\bar{x}_0$ .

## 5 Boundary arcs, relative degree $r^{wu} = 2k + 1$ , $k \geq 0$

In the previous section nothing has been said about the influence of the sign of the leading Markov parameter on the behaviour of optimal trajectories on  $\partial\Phi$  (i.e. on boundary arcs in  $\mathcal{V}^*$ ). We study this aspect now, i.e. we study what happens when a solution of the system in (4) has a boundary arc. From Lemma 1 one sees that systems in (2) with  $r^{wu} = 2k + 1$ ,  $k \geq 0$  yield  $(\tilde{A}, \tilde{B}, \tilde{C})$  with a leading Markov parameter that is always (semi) negative definite. Let us investigate the consequences on the well-posedness of the dynamics on boundary arcs, i.e. those portions of optimal trajectories that satisfy  $Cx(t) + D = 0$  for  $t \in (\tau, \tau + \epsilon) \subseteq [0, T_1]$ ,  $\epsilon > 0$ . On  $(\tau, \tau + \epsilon)$  one has  $u(t) = R^{-1}B^T\eta(t)$  and  $\dot{\eta}(t) = Qx(t) - A^T\eta(t) - C^T\lambda(t)$ , with  $\lambda(t)$  the solution of the LCP in (38) (equal to  $g_r(t)$  in this case, where  $g_r(\cdot)$  is in (32) and proposition 2). One notices that the complementarity relations  $0 \leq \{z_i\}(t^+) \perp g_i(t) \geq 0$  are trivially satisfied for all  $1 \leq i \leq r - 1$  (see (33)).

**Remark 9** Following [13, §3.10, 3.11] let us define the Hamiltonian function  $\bar{H}(x, u, \eta, \alpha) = H(x, u, \eta) + \alpha^T z_{r^{wu}+1}$ . On a boundary arc one has  $z_{r^{wu}+1}(t) = \dot{z}_{r^{wu}+1}(t) = CA^{r^{wu}}x(t) + CA^{r^{wu}-1}Bu(t) = 0$ . Thus we deduce that  $u(t) = -(CA^{r^{wu}-1}B)^{-1}CA^{r^{wu}}x(t)$  (and consequently  $B^T\eta(t) = -B^T(CA^{r^{wu}-1}B)^{-1}CA^{r^{wu}}x(t)$ ). Let us now compute the multiplier  $\alpha(\cdot)$  from  $\frac{\partial \bar{H}}{\partial u} = 0$ . One finds  $\alpha(t) = (-1)^{r^{wu}+1}[M^{(r)}]^{-1}(CA^{r^{wu}}x(t) + CA^{r^{wu}-1}BB^T\eta(t))$ . Inserting the value for  $B^T\eta(t)$  one finds that  $\alpha(t) = 0$  which is consistent with (33) since  $\alpha(\cdot) = g_{r^{wu}+1}(\cdot)$ , the multiplier associated with the coordinate  $z_{r^{wu}+1}$ .

Let us examine the case  $r^{wu} = 1$ . On boundary arcs one has  $w(t) = \tilde{C}\tilde{x}(t) + D = 0$  and  $\dot{w}(t) = \tilde{C}\tilde{A}\tilde{x}(t) = 0$ , since  $r = 2$ , and  $\ddot{w}(t) = \tilde{C}\tilde{A}^2\tilde{x}(t) + M^{(2)}\lambda(t) = 0$ . The fact that the trajectory keeps on evolving on  $\partial\Phi$  or detaches from  $\partial\Phi$  is monitored by the LCP  $0 \leq \lambda(t) \perp \ddot{w}(t) = \tilde{C}\tilde{A}^2\tilde{x}(t) + M^{(2)}\lambda(t) \geq 0$ .

The fact that the dynamics on boundary arcs is well-posed or not is closely linked to the well-posedness of this LCP, which holds on the interval  $(\tau, \tau + \epsilon)$ . The following is true since  $M^{(2)} < 0$ :

**Lemma 6** *Let  $m = 1$ . Consider the LCP  $0 \leq \lambda(t) \perp \tilde{C}\tilde{A}^2\tilde{x}(t) + M^{(2)}\lambda(t) \geq 0$ . Then:*

- *If  $\tilde{C}\tilde{A}^2\tilde{x}(t) < 0$ , the LCP has no solution,*
- *If  $\tilde{C}\tilde{A}^2\tilde{x}(t) = 0$ , the LCP has one solution  $\lambda(t) = 0$ ,*
- *If  $\tilde{C}\tilde{A}^2\tilde{x}(t) > 0$ , the LCP has two solutions  $\lambda(t) = 0$  and  $\lambda(t) = -(M^{(2)})^{-1}\tilde{C}\tilde{A}^2\tilde{x}(t)$ .*

Let the state satisfy  $\tilde{C}\tilde{A}^2\tilde{x}(\tau) > 0$  at an entry time  $\tau \in [0, T_1]$ . The last item shows that there may not be uniqueness of the solution to the optimal control problem, if the optimal trajectory grazes the constraint boundary  $\partial\Phi$ . Either the trajectory “detaches” from  $\partial\Phi$  ( $\lambda(t) = 0 \Rightarrow \ddot{w}(t) > 0$ ), or remains on  $\partial\Phi$  ( $\lambda(t) > 0 \Rightarrow \ddot{w}(t) = 0$ ). This is not surprising since the leading Markov parameter is negative. The problem loses its convexity and the system in (4) is no longer well-posed as an IVP. It seems that little attention has been paid in the literature to the fact that the multiplier  $\lambda$  is the solution of a LCP on boundary arcs. The boundary arcs input satisfy  $\dot{w}(t) = CAx(t) + CBu(t)$ , however this may not be the optimal controller. We shall see later what may happen when  $\tilde{C}\tilde{A}^2\tilde{x}(\tau) < 0$  at an entry time  $\tau$ . Lemma 6 can be extended to the general case  $r^{wu} = 2k + 1$  and  $m \geq 1$ , where the LCP in (38) has to be studied. In view of Propositions 13 and 14 and Corollary 7, there are severe restrictions for the existence of entry times for systems with odd  $r^{wu}$ . The problem of interest here is to investigate what may happen after a contact time.

**Proposition 18** *Let  $r^{wu} = 2k + 1$ ,  $k \geq 0$ , and let  $m = 1$ . Let  $x(\tau^+) \in \partial\Phi$ . Then:*

- *If  $\tilde{C}\tilde{A}^r\tilde{x}(\tau^+) < 0$ , the trajectory leaves  $\Phi$  in a right neighborhood of  $\tau$ .*
- *If  $\tilde{C}\tilde{A}^r\tilde{x}(\tau^+) = 0$ ,  $\tau$  may be an entry time followed by a grazing trajectory.*
- *If  $\tilde{C}\tilde{A}^r\tilde{x}(\tau^+) > 0$ ,  $\tau$  may be either a touch time ( $\lambda(\tau^+) = 0$ ) or an entry time ( $\lambda(\tau^+) > 0$ ).*

Consequently, optimal trajectories satisfy  $\tilde{C}\tilde{A}^r x(\tau^+) \geq 0$ .

**Proof :** First let us notice that it is possible that  $\bar{z}(\tau^-) = 0$ , which implies from (30) (31) that  $\bar{z}(\tau^+) = 0$ , but that  $\tilde{C}\tilde{A}^r x(\tau^+) < 0$ , because of the value of the zero-dynamics state  $\xi(\tau)$ . In other words, odd- $r^{wu}$  systems can possess entry times with  $\deg(\lambda) \leq 1$  (see Proposition 14). However  $\xi(\tau)$  has to be such that  $\tilde{C}\tilde{A}^r x(\tau^+) \geq 0$  for a boundary arc to exist on the right of  $\tau$ . Clearly for an even- $r^{wu}$  system, the trajectory can be kept in  $\Phi$  as there always exists a multiplier  $g_r(\tau^+)$  solution of the LCP (38). ■

Proposition 18 suggests that the zero dynamics plays a major role in the well-posedness of the LCP and whether or not odd- $r^{wu}$  systems possess boundary arcs. For instance if  $r^{wu} = n$ , the first item of Proposition 18 becomes irrelevant at an hypertangential contact state, because necessarily  $\tilde{C}\tilde{A}^r x(\tau^+) = 0$ .

**Corollary 10** *Let  $r^{wu} = m = n_u = 1$ . If an entry time  $\tau$  exists on  $[0, T_1]$ , then uniqueness of the solution of the necessary conditions system holds if and only if the trajectory along the boundary arc is grazing, that is  $(CA^2 + CBB^T Q)x(t) + (CABB^T - CBB^T A^T)\eta(t) = 0$  for all  $t \in [\tau, \tau + \epsilon]$  where  $\tau + \epsilon$  is the exit time.*

**Proof :** Follows from Proposition 18 and the calculation of  $\tilde{C}\tilde{A}^2$ . ■

**Definition 4** *A boundary arc with entry time  $\tau$  and exit time  $\tau + \epsilon$  is well-posed if the LCP  $0 \leq w^{(r)}(t) \perp \lambda(t) \geq 0$  possesses at least one solution for all  $t \in (\tau, \tau + \epsilon)$ .*

On a boundary arc  $\bar{z}^T = (z_1, z_2, \dots, z_r) = (0, \dots, 0)$ . Definition 4 applies to systems with  $m \geq 1$ . The next Corollary concerns scalar systems as  $\dot{x}(t) = ax(t) + bu(t)$ ,  $w(t) = cx(t) + d$ ,  $Q = R = 1$ . Then the Hamiltonian dynamics is  $\begin{pmatrix} \dot{x}(t) \\ \dot{\eta}(t) \end{pmatrix} = \begin{pmatrix} a & b^2 \\ 1 & -a \end{pmatrix} \begin{pmatrix} x(t) \\ \eta(t) \end{pmatrix} + \begin{pmatrix} 0 \\ -1 \end{pmatrix} \lambda$ . The LCP to be solved on boundary arcs is  $0 \leq \lambda \perp -b^2c^2\lambda - d(a^2 + b^2) \geq 0$ .

**Corollary 11 i)** If  $r^{wu} = n = 1$  and  $cx_0 + d > 0$ ,  $cx_1 + d = 0$ , then dynamics on boundary arcs is well-posed if and only if  $d \leq 0$ . Moreover  $\partial\Phi = \{x = -\frac{d}{c}\}$  is attained only at  $t = T_1$  on the optimal trajectory. **ii)** Let us now consider  $n \geq 2$ ,  $D = 0$  and  $m = 1$  in (1). If  $r^{wu} < n$ , then the well-posedness of the dynamics along boundary arcs implies that  $\tilde{C}\tilde{A}^2\tilde{W}^{-1} \begin{pmatrix} 0 \\ \xi(t) \end{pmatrix} \geq 0$ ,  $\forall t \in (\tau, \tau + \epsilon)$ .

**Proof :** Since ii) is an easy consequence of Lemma 6, we just prove i). Integrating the dynamics without constraint one finds the optimal trajectory on  $[0, T_1]$ :

$$x(t) = \frac{x_0 \exp(-\sqrt{a^2+b^2}T_1) - x_1}{\exp(-\sqrt{a^2+b^2}T_1) - \exp(\sqrt{a^2+b^2}T_1)} \exp(\sqrt{a^2+b^2}t) + \frac{x_1 - x_0 \exp(\sqrt{a^2+b^2}T_1)}{\exp(-\sqrt{a^2+b^2}T_1) - \exp(\sqrt{a^2+b^2}T_1)} \exp(-\sqrt{a^2+b^2}t) \quad (54)$$

The optimal control can be calculated with  $\eta(t) = \frac{1}{b^2}(\dot{x}(t) - ax(t))$  and  $u(t) = b\eta(t)$ . It can be checked that  $x(t) > -\frac{d}{c}$  for all  $t \in [0, T_1]$ . Consequently the boundary  $\partial\Phi$  is attained only at  $t = T_1$ , and  $T_1$  is a *contact* time. The LCP for  $\lambda$  on  $\partial\Phi$  is given by  $0 \leq \lambda \perp -d(a^2 + b^2) - c^2b^2\lambda \geq 0$ . We deduce that if  $-d < 0$  the LCP has no solution, if  $d = 0$  then  $\lambda = 0$  and if  $-d > 0$  then two solutions  $\lambda = 0$  and  $\lambda = -\frac{d(a^2+b^2)}{c^2b^2}$  are possible. However the behaviour on the boundary has no consequence on the value of the action integral  $I(u)$  since the optimal control and state trajectory are Lebesgue measurable functions. ■

The last condition in Corollary 11 can be written also as  $\tilde{C}\tilde{A}^2\tilde{W}^{-1} \begin{pmatrix} 0 \\ \xi(\tau) \exp(\tilde{A}_\xi(t - \tau)) \end{pmatrix} \geq 0$ . The matrix  $\tilde{A}_\xi$  is characterized from the zero dynamics of  $(A, B, C)$ , see Lemma 2. In case i) one sees that  $\mathcal{V}^* = \{-\frac{d}{c}\}$  and  $u^* = \frac{ad}{cb}$ . From Proposition 9 boundary arcs and contact times belong to  $\mathcal{V}^*$ , which is consistent with Corollary 11 i).

Let us now consider the planar system  $A = \begin{pmatrix} 0 & 1 \\ a & b \end{pmatrix}$ ,  $B = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ ,  $C = (c \ d)$ ,  $D = e$ . This system has transfer function  $G(s) = \frac{c+ds}{s(s-b)-a}$  and relative degree  $r^{wu} = 1$  and  $CB = d \neq 0$ . Therefore  $M^{(2)} = -d^2 < 0$ . Let us take  $Q = \text{diag}(q_1, q_2)$ . To this system we associate its optimality conditions as in (4) and we check that  $\tilde{C}\tilde{A}\tilde{B} = -d^2$ . We get  $H_{w\lambda}(s) = \frac{(c+ds)(c-ds)}{(s^2+bs-a)(s^2-bs-a)-s^2+1}$ .

**Lemma 7** Assume that  $-ad^2 + c^2 + bcd \neq 0$  and that  $c \neq 0$ . The above planar system has a well-posed boundary arc on the time interval  $[\tau, \tau + \epsilon] \subseteq [0, T_1]$  if and only if

$$\begin{cases} \dot{\eta}_1(t) = \alpha\eta_2(t) - c\eta_1(t) + \delta \\ \dot{\eta}_2(t) = -(1+d)\eta_1(t) + \beta\eta_2(t) \\ \eta(\tau) = \begin{pmatrix} e_3^T \\ e_4^T \end{pmatrix} \exp(\tilde{A}\tau)\tilde{x}(0) \\ \gamma\eta_2(t) - d\eta_1(t) + \epsilon \geq 0 \text{ for all } t \in [\tau, \tau + \epsilon] \end{cases} \quad (55)$$

where  $\alpha = \frac{q_1d^2+a^2d^2-ac^2-abcd}{-ad^2+c^2+bcd} - \frac{c}{d} \frac{abd^3-cd^2(q_2+a+b^2+c^3)}{ad^2-c^2-bcd}$ ,  $\beta = \frac{-q_2dc+abd^2-bc^2-b^2cd+abd^3-cd^2(q_2+a+b^2)+c^3}{-ad^2+c^2+bcd}$ ,  $\gamma = \frac{abd^3-cd^2(q_2+a+b^2)+c^3}{-ad^2+c^2+bcd}$ ,  $\epsilon = -e\frac{a}{c}(c+bd)[ad^2(-ad^2+c^2+bcd)+1]$ ,  $\delta = -q_1\frac{ad^2e}{c}(-ad^2+c^2+bcd) - q_1\frac{e}{c}$ .

**Proof :** On the boundary arc one has  $\tilde{C}\tilde{x}(t) + D = 0$  and  $\tilde{C}\tilde{A}\tilde{x}(t) = 0$ , which respectively yield  $cx_1(t) + dx_2(t) + e = 0$  and  $adx_1(t) + (c+bd)x_2(t) + d\eta_2(t) = 0$ . We deduce that  $x_2(t) = \frac{-cd}{-ad^2+c^2+bcd}\eta_2(t) + ade(-ad^2+c^2+bcd)$ , and  $x_1(t) = \frac{d^2}{-ad^2+c^2+bcd}\eta_2(t) - \frac{ad^2e}{c}(-ad^2+c^2+bcd) - \frac{e}{c}$ . Moreover from the values of  $\tilde{A}$  and  $\tilde{C}$  we get that

$$\begin{cases} \dot{\eta}_1(t) = q_1x_1(t) - a\eta_2(t) - c\lambda(t) \\ \dot{\eta}_2(t) = q_2x_2(t) - \eta_1(t) - b\eta_2(t) - d\lambda(t) \\ 0 \leq cx_1(t) + dx_2(t) + e \perp \lambda(t) \geq 0 \end{cases} \quad (56)$$

Along the boundary arc we have from Lemma 6,  $\lambda(t) = \frac{1}{d^2} \tilde{C} \tilde{A}^2 \tilde{x}(t)$  provided  $\tilde{C} \tilde{A}^2 \tilde{x}(t) \geq 0$ . Now  $\tilde{C} \tilde{A}^2 \tilde{x}(t) = a(c+bd)x_1(t) + (cb+d(q_2+a+b^2))x_2(t) - d\eta_1(t) + c\eta_2(t)$ . Combining the above values we get that  $\lambda(t) = \frac{abd^3 - cd^2(q_2+a+b^2) + c^3}{-d^2(ad^2 - c^2 - bcd)} \eta_2(t) - \frac{1}{d} \eta_1(t)$ . Injecting all the calculated values for  $x_1(t)$ ,  $x_2(t)$  and  $\lambda(t)$  into (56), the result follows. The conditions are sufficient but also necessary for if  $\gamma\eta_2(t) - d\eta_1(t) + \epsilon < 0$  the boundary arc does not exist at time  $t$ . ■

Since the system in (55) is linear of order two, it is quite possible in practice to test whether or not (55) holds. This provides conditions on the entry values  $\eta_1(\tau)$  and  $\eta_2(\tau)$  such that a boundary arc is well-posed or not. Clearly Lemma 7 does not answer the question whether or not the optimal trajectory possesses boundary arcs. Notice that  $\tilde{x}(0)$  is not known as  $\eta(0)$  is an unknown of the BVP.

## 6 Numerical solution of the BVP

Contrary to the unconstrained case [54, Lemma 8.2.9], it is not clear here which IVP corresponds to the BVP in (4). The time-stepping scheme proposed in [1] to integrate the IVP may be used in a multiple-shooting algorithm (at least in brute-force algorithms). However when facing a minimisation problem as in (1)-(2) (which is a BVP), one doesn't know a priori how many times the optimal trajectory will reach the boundary  $\partial\Phi$ , how many times it will detach from  $\partial\Phi$ , and how the derivatives of  $w(\cdot)$  will evolve along the optimal path. In particular whether or not distributional inputs and costate will be needed is also an unknown of the minimisation problem. From a numerical point of view, this is exactly the difference between time-stepping schemes and event-driven schemes [12]. One could say that time-stepping discretisations are time-discretisations of the measure differential formalism (because *measures* are approximated numerically), whereas event-driven algorithms are time-discretisations of systems of the form  $\dot{x}(t) = f(x, t)$  if  $t \neq t_k$ ,  $x(t_k^+) = \mathcal{F}[x(t_k^-)]$  for  $t = t_k$ . Time-stepping schemes do not require the knowledge of state jump times but work with constant time steps. Moreover they accommodate accumulations of events (thus regular solutions as in definition 1 $\hat{A}$  can be approximated), which event-driven and multiple shooting schemes do not [42]. When solving a BVP with a time-stepping scheme, one does not need to guess the jump state times (or more generally the existence of any junction time), but runs the simulation in one shot. This is thought to be a great advantage over event-driven and multiple shooting algorithms. A time-stepping algorithm solving the BVP in (4), starting from the extended action in (50) and using the IVP solver presented in [1], will be the object of a future work. This belongs to the class of so-called direct methods [56] in which the optimal control problem is transformed into a nonlinear programming problem [13, §7.11]. However the crucial point is to take into account possible state  $\{x\}(\cdot)$  and costate  $\{\eta\}(\cdot)$  jumps in the time discretization algorithm, without a priori assumption on the number of junction times.

## 7 Conclusions

This paper concerns the optimal control of linear invariant systems, with inequality (or unilateral) constraints on the state. This is a major topic which has continuously attracted the interest of researchers since more than 70 years. The main contributions are the following. The Bolyanskii-Pontryagin necessary conditions are embedded into a distributional differential inclusion framework (the higher order Moreau's sweeping process) which allows us to clearly understand their dynamics and provides a clear formalism for the derivation of an extended action to be minimized (allowing for costate and/or inputs and/or system's state, to belong to a specific set of Schwarz' distributions). The powerful geometric tools introduced by ten Dam in the context of unilaterally constrained controlled dynamical systems, are used to improve the qualitative analysis of optimal trajectories and controllers. They allow us to generalize results concerning the specificity of odd relative degree systems (one of the results being that odd relative degree systems can possess boundary arcs with entry times such that the multiplier associated to the unilateral constraint is a distribution of degree  $\geq 3$ ; then the action to be minimized has to be extended to incorporate the state and/or costate jumps). The theory of complementarity problems is

shown to be quite useful in order to better understand the behaviour of optimal path on the boundary of the admissible domain. Most of the tools which are used in this paper (distributional inclusion, ten Dam's geometrical study, complementarity problems) have never been introduced before in the context of optimal control with unilateral state constraints. The proposed framework paves the way to many extensions, and towards the design of time-stepping algorithms which may constitute a nice alternative to multiple shooting algorithms to solve the boundary value problem.

## A Some mathematical definitions

Given  $A \in \mathbb{R}^{n \times n}$  and  $B \in \mathbb{R}^n$ , the problem of finding  $x \in \mathbb{R}^n$  and  $y \in \mathbb{R}^n$  satisfying

$$y = Ax + B \geq 0, x \geq 0, x^T y = 0 \quad (57)$$

is called a Linear Complementarity Problem (LCP). It can be equivalently written as

$$0 \leq x \perp y = Ax + B \geq 0 \quad (58)$$

Roughly the LCP has a unique solution  $x^*$  whatever  $B$  if and only if  $A$  satisfies some positivity conditions, see [18]. Positive definiteness of  $A$  is sufficient. The so-called **P**-property of matrices is necessary and sufficient for the LCP to have a unique solution for any  $B$ .

The next notions may be found in [1, 38, 32] Let  $I$  denote a non-degenerate real interval (not empty nor reduced to a singleton).

• By  $z \in BV(I; \mathbb{R}^n)$  it is meant that  $z$  is a  $\mathbb{R}^n$ -valued function of bounded variation if there exists a constant  $C > 0$  such that for all finite sequences  $t_0 < t_1 < \dots < t_N$  ( $N$  arbitrary) of points of  $I$ , we have

$$\sum_{i=1}^N \|z(t_i) - z(t_{i-1})\| \leq C.$$

Let  $J$  be a subinterval of  $I$ . The real number

$$\text{var}(z, J) := \sup \sum_{i=1}^N \|z(t_i) - z(t_{i-1})\|,$$

where the supremum is taken with respect to all the finite sequences  $t_0 < t_1 < \dots < t_N$  ( $N$  arbitrary) of points of  $J$ , is called the variation of  $z$  in  $J$ .

Any BV function has a countable set of discontinuity points and is almost everywhere differentiable. A BV function defined on  $[a, b] \subset I$  possesses left-limits in  $]a, b]$  and right-limits in  $[a, b[$ . Moreover, the functions  $t \mapsto z(t^+) := \lim_{s \rightarrow t, s > t} z(s)$  and  $t \mapsto z(t^-) := \lim_{s \rightarrow t, s < t} z(s)$  are both BV functions.

• We denote by  $LBV(I; \mathbb{R}^n)$  the space of functions of locally bounded variation, i.e. of bounded variation on every compact subinterval of  $I$ .

• We denote by  $RCLBV(I; \mathbb{R}^n)$  the space of right-continuous functions of locally bounded variation. It is known that if  $z \in RCLBV(I; \mathbb{R}^n)$  and  $[a, b]$  denotes a compact subinterval of  $I$ , then  $z$  can be represented in the form:

$$z(t) = \mathcal{J}_z(t) + [z](t) + \zeta_z(t), \forall t \in [a, b],$$

where  $\mathcal{J}_z$  is a jump function,  $[z]$  is an absolutely continuous function and  $\zeta_z$  is a singular function. Here  $\mathcal{J}_z$  is a jump function in the sense that  $\mathcal{J}_z$  is right-continuous and given any  $\varepsilon > 0$ , there exist finitely many points of discontinuity  $t_1, \dots, t_N$  of  $\mathcal{J}_z$  such that  $\sum_{i=1}^N \|\mathcal{J}_z(t_i) - \mathcal{J}_z(t_i^-)\| + \varepsilon > \text{var}(\mathcal{J}_z, [a, b])$ ,  $[z]$  is an absolutely continuous function in the sense that for every  $\varepsilon > 0$ , there exists  $\delta > 0$  such that  $\sum_{i=1}^N \|[z](\beta_i) - [z](\alpha_i)\| < \varepsilon$ , for any collection of disjoint subintervals  $]\alpha_i, \beta_i] \subset [a, b]$  ( $1 \leq i \leq N$ ) such

that  $\sum_{i=1}^N (\beta_i - \alpha_i) < \delta$ , and  $\zeta_u$  is a singular function in the sense that  $\zeta_z$  is a continuous and bounded variation function on  $[a, b]$  such that  $\dot{\zeta}_z = 0$  almost everywhere on  $[a, b]$ .

• By  $z \in RC SLBV(I; \mathbb{R}^n)$  it is meant that  $z$  is a right-continuous function of special locally bounded variation, i.e.  $z$  is of bounded variation and can be written as the sum of a jump function and an absolutely continuous function on every compact subinterval of  $I$ . So, if  $z \in RC SLBV(I; \mathbb{R}^n)$  then

$$z = [z] + \mathcal{J}_z \tag{59}$$

where  $[z]$  is a locally absolutely continuous function called the absolutely continuous component of  $z$  and  $\mathcal{J}_z$  is uniquely defined up to a constant by

$$\mathcal{J}_z(t) = \sum_{t \geq t_n} z(t_n^+) - z(t_n^-) = \sum_{t \geq t_n} z(t_n) - z(t_n^-) \tag{60}$$

where  $t_1, t_2, \dots, t_n, \dots$  denote the countably many points of discontinuity of  $z$  in  $I$ .

**Stieltjes measure.** Let  $z \in LBV(I; \mathbb{R}^n)$  be given. We denote by  $du$  the Stieltjes measure generated by  $z$ . Recall that for  $a \leq b$ ,  $a, b \in I$ :

$$\begin{aligned} dz([a, b]) &= z(b^+) - z(a^-), \\ dz([a, b]) &= z(b^-) - z(a^-), \\ dz(]a, b]) &= z(b^+) - z(a^+), \\ dz(]a, b]) &= z(b^-) - z(a^+). \end{aligned}$$

In particular, we have

$$dz(\{a\}) = z(a^+) - z(a^-).$$

## B Example of a function in $\mathcal{F}_\infty$ and distributions in $\mathcal{T}_n$

This is taken from [1]. Set  $I = (0, +\infty)$  and let  $z : \tilde{I} \rightarrow \mathbb{R}$  be the function given by

$$z(t) = |\sin(t)|, \quad \forall t \geq 0.$$

It is clear that

$$\hat{z}^{(0)} := z \in RC SLBV(\tilde{I}; \mathbb{R})$$

since  $z(\cdot)$  is Lipschitz-continuous. Then we obtain

$$\hat{z}^{(1)}(t) := \frac{d^+}{dt} [\hat{z}^{(0)}](t) = \frac{d^+ z}{dt}(t) = \cos(t - k\pi) \text{ if } t \in [k\pi, (k+1)\pi), (k \in \mathbb{N}).$$

We see that  $E_0(\hat{z}^{(1)}) = \{k\pi; k \in \mathbb{N} \setminus \{0\}\}$  and

$$\hat{z}^{(1)}(\cdot) = [\hat{z}^{(1)}](\cdot) + J(\cdot),$$

where

$$[\hat{z}^{(1)}](t) = -2k + \cos(t - k\pi) \text{ if } t \in [k\pi, (k+1)\pi], (k \in \mathbb{N})$$

and

$$J(t) = 2k \text{ if } t \in [k\pi, (k+1)\pi], (k \in \mathbb{N}).$$

Thus

$$\hat{z}^{(1)}(\cdot) \in RC SLBV(\tilde{I}; \mathbb{R})$$

Then

$$\hat{z}^{(2)}(t) := \frac{d^+}{dt} [\hat{z}^{(1)}](t) = -|\sin(t)|$$



so that

$$\hat{z}^{(2)}(\cdot) \in RC SLBV(\tilde{I}; \mathbb{R}).$$

And so on, we see that

$$\hat{z}^{(k)}(t) = \begin{cases} (-1)^m \hat{z}^{(0)}(t) & \text{if } k = 2m \\ (-1)^m \hat{z}^{(1)}(t) & \text{if } k = 2m + 1 \end{cases}, \quad (m \in \mathbb{N}),$$

so that  $\hat{z}^{(k)} \in RC SLBV(\tilde{I}; \mathbb{R})$ ,  $\forall k \in \mathbb{N}$ , and thus

$$z(\cdot) \in \mathcal{F}_\infty(\tilde{I}; \mathbb{R}).$$

Let us now consider the distribution  $T$  defined by

$$\langle T, \varphi \rangle = \int_{-\infty}^{\infty} |\sin(t)| \varphi(t) dt, \quad \forall \varphi \in C_0^\infty(I).$$

Then for a given function  $\varphi \in C_0^\infty(I)$ , we see that:

$$\begin{aligned} \langle DT, \varphi \rangle &= \int_{-\infty}^{\infty} \hat{z}^{(1)}(t) \varphi(t) dt = \sum_{k \in \mathbb{N} \cap \text{supp}\{\varphi\}} \int_{k\pi}^{(k+1)\pi} \cos(t - k\pi) \varphi(t) dt, \\ \langle D^2T, \varphi \rangle &= \int_{-\infty}^{\infty} \hat{z}^{(2)}(t) \varphi(t) dt + \sum_{k \in \mathbb{N} \setminus \{0\} \cap \text{supp}\{\varphi\}} (\hat{z}^{(1)}(k\pi^+) - \hat{z}^{(1)}(k\pi^-)) \langle \delta_{k\pi}, \varphi \rangle = \\ &= - \int_{-\infty}^{\infty} |\sin(t)| \varphi(t) dt + 2 \sum_{k \in \mathbb{N} \setminus \{0\} \cap \text{supp}\{\varphi\}} \langle \delta_{k\pi}, \varphi \rangle, \end{aligned}$$

and so on. We have:

$$\begin{aligned} T &\equiv \{T\} = \hat{z}^{(0)} = |\sin(\cdot)|, \quad \text{deg}(T) = 0, \\ DT &\equiv \{T^{(1)}\} = \{DT\} = \hat{z}^{(1)} = \cos(\cdot - k\pi) \text{ on } [k\pi, (k+1)\pi) \quad (k \in \mathbb{N}), \quad \text{deg}(DT) = 1, \\ D^2T &\equiv \ll D^2T \gg = -|\sin(\cdot)| + 2 \sum_{k \in \mathbb{N} \setminus \{0\}} \delta_{k\pi}, \quad \text{deg}(D^2T) = 2, \\ \{T^{(2)}\} &= \{D^2T\} = \hat{z}^{(2)} = -|\sin(\cdot)|, \end{aligned}$$

and

$$d \ll D^2T \gg = d\hat{z}^{(1)}.$$

## C Hints on the controllability of (2)

**Example 6** *Let us consider the system*

$$\begin{cases} \dot{x}_1(t) = x_2(t) \\ \dot{x}_2(t) = u(t) \\ x_2(t) \geq 0 \end{cases} \quad (61)$$

Clearly if  $\mathbf{U}$  is restricted to the set of bounded functions of time, then  $(A, B)$  is not controllable in the set  $\Phi = \{x \mid x_2 \geq 0\}$ . However if distributions of degree 3 (derivatives of Dirac measures in the sense of Schwartz' distributions) are allowed in  $\mathbf{U}$ , then controllability holds (but, the quadratic term in  $u$  in (1) is meaningless for such inputs). Indeed negative jumps in  $x_1(\cdot)$  are then possible to reach a state  $x(T_1)$  from  $x(0)$  with  $x_1(T_1) < x_1(0)$ ,  $T_1 > 0$ . In particular  $x(T_1) = 0$  belongs to the reachable space from  $x(0^-) > 0$ . If the constraint is replaced by  $x_1(t) + d \geq 0$ ,  $d \in \mathbb{R}$ , then controllability holds in  $\Phi$  with  $\mathbf{U}$  a set of functions.

Let us consider the triple  $(A, B, C)$  in (2) with  $D = 0$ , and its  $z$ -dynamics. Let us consider  $u(t) = (CA^{r^{wu}-1}B)^{-1}[-CA^{r^{wu}}Wz(t) + v(t)]$ . We obtain

$$\begin{cases} z_1^{(r^{wu})}(t) = v(t) \\ \dot{\xi}(t) = A_\xi \xi(t) + B_\xi z_1(t) \\ w(t) = z_1(t) \geq 0 \end{cases} \quad (62)$$

Let us consider that all variables in (62) are functions (i.e. we exclude inputs  $v(\cdot)$  with degree  $\geq 2$ ). Complete controllability of (2) in  $\Phi$  means that for all  $\bar{x}_0, \bar{x}_1$  in  $\Phi$ , there exists  $u(\cdot)$  which steers  $x(\cdot)$  from  $x(0) = \bar{x}_0$  to  $x(T_1) = \bar{x}_1$  for some  $T_1 \geq 0$ . One sees from (62) that the system is completely controllable in  $\Phi = \{z \mid z_1 \geq 0\}$  only if the pair  $(A_\xi, B_\xi)$  is completely controllable with positive inputs. From [10, Corollary 3.7] the following holds:

**Lemma 8** *The system in (2) is completely controllable only if the pair  $(A_\xi, B_\xi)$  is controllable ( $\Leftrightarrow$  the associated Kalman matrix has rank  $n - r^{wu}$ ) and there is no real eigenvector  $\mu$  of  $A_\xi^T$  satisfying  $\mu^T B_\xi z_1 \leq 0$  for all  $z_1 \geq 0$ .*

It is easy to see that the system in (61) fails to satisfy the necessary condition as  $A_\xi = 0$  and  $B_\xi = 1$ . Lemma 8 can be used to determine which systems may necessitate distributional inputs to reach  $\bar{x}_1$  from  $\bar{x}_0$  (of course there may also be a large set of boundary conditions such that a function of time input is sufficient, as (61) shows: we are dealing here with *controllability* in  $\Phi$ ).

**Corollary 12** *Let the system in (2) be given by  $x_1^{(n)}(t) = u(t)$ ,  $x_1^{(k-1)}(t) = x_k(t) \geq 0$ . If  $2 \leq k \leq n$  the system is not completely controllable in  $\Phi$  with  $u(\cdot)$  a function.*

**Proof:** Since  $k \geq 2$  the  $z$ -dynamics reads  $z_1^{(r^{wu})}(t) = u(t)$ ,  $\xi^{(n-r^{wu})}(t) = z_1(t)$ , with  $r^{wu} = n - k + 1$ . A real eigenvector  $\mu$  of  $A_\xi^T$  with eigenvalue  $\alpha$  satisfies  $A_\xi^T \mu = \alpha \mu$ . It can be computed that this equality

is  $\begin{pmatrix} 0 \\ \mu_1 \\ \dots \\ \mu_{k-2} \end{pmatrix} = \alpha \begin{pmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_{k-1} \end{pmatrix}$ . Also  $\mu^T B_\xi z_1 = \mu_{k-1} z_1$ . It follows that either  $\alpha \neq 0$  and  $\mu = 0^{k-1}$ , or  $\alpha = 0$  and  $\mu_i = 0$  for  $1 \leq i \leq k - 2$ , whereas  $\mu_{k-1} \in \mathbb{R}$ . In any case the inequality  $\mu_{k-1} z_1 > 0$  for all  $z_1 \geq 0$  cannot be satisfied. From Lemma 8 the system cannot be completely controllable in  $\Phi$  with function-of-time inputs. ■

Another necessary condition for complete controllability of (2) stemming from [10, Corollary 3.7] is that  $A_\xi$  possesses only purely imaginary eigenvalues. Otherwise controllability of  $\dot{\xi}(t) = A_\xi \xi(t) + B_\xi z_1(t)$  with  $z_1(t) \geq 0$  fails. These necessary conditions may be considered as caution flags when formulating the optimal control problem (1) (2). In the framework of the HOSP one always has  $\deg(z_1) \leq 1$  so Corollary 12 and Lemma 8 continue to hold if (62) is embedded into (17) (18)(19).

## References

- [1] V. Acary, B. Brogliato, D. Goeleven, 2006 "Higher order Moreau's sweeping process: Mathematical formulation and numerical simulation", Mathematical Programming, in press. Preliminary version available at <http://hal.inria.fr>
- [2] V. Acary, B. Brogliato, 2003 "Higher order Moreau's sweeping process", Colloquium in the honour of the 80-th birthday of J.J. Moreau, 17/18 November 2003, Montpellier, France. In *Progresses in Nonsmooth Mechanics and Analysis* (P. Alart, O. Maissenneuve, R.T. Rockafellar, Eds.), Advances in Mathematics and Mechanics, Kluwer (2005), pp.261-277.

- 
- [3] K.M. Anstreicher, 1983 “Generation of feasible descent directions in continuous time linear programming”, Systems Optimization Laboratory, Dept. of Operations Research, Stanford university, technical report SOL 83-18, September 1983.
- [4] A.V. Arutyunov, S.M. Aseev, 1995 “State constraints in optimal control. The degeneracy phenomenon”, Systems and Control Letters, vol.26, pp.267-273.
- [5] D. Augustin, H. Maurer, 2001 “Second order sufficient conditions and sensitivity analysis for the optimal control of a container crane under state constraints”, Optimization, vol.49, pp.351-368.
- [6] R. Bellman, 1953 “Bottleneck problems and dynamical programming”, Proc. National Academy of Sciences, vol.39, pp.947-951.
- [7] P. Berkman, H.J. Pesch, 1995 “Abort landing in windshear: Optimal control problem with third-order state constraint and varied switching structure”, Journal of Optimization Theory and Applications, vol.85, no 1, pp.21-57, April.
- [8] B. Bonnard, L. Faubourg, G. Launay, E. Trélat, 2003 “Optimal control with state constraints and the space shuttle re-entry problem”, Journal of Dynamical and Control Systems, vol.9, no 2, pp.155-200.
- [9] B. Bonnard, E. Trélat, 2002 “Une approche géométrique de contrôle optimal de l’arc atmosphérique de la navette spatiale”, ESAIM: Control, Optimisation and Calculus of Variations, vol.7, pp.179-222.
- [10] R.F. Brammer, 1972 “Controllability in linear autonomous systems with positive controllers”, SIAM J. Control, vol.10, no 2, pp.339-353, May.
- [11] B. Brogliato, 1999 *Nonsmooth Mechanics*, 2nd edition, Springer Verlag London, Communications and Control Engineering Series. (erratum and addendum at <http://www.inrialpes.fr/bipop/people/brogliato/brogli.html> )
- [12] B. Brogliato, A.A. ten Dam, L. Paoli, F. Génot, M. Abadie, 2002 “Numerical simulation of finite dimensional multibody nonsmooth mechanical systems”, ASME Applied Mechanics Reviews, vol.55, no 2, pp.107-150, March.
- [13] A.E. Bryson Jr, Y.C. Ho, 1975 *Applied Optimal Control; Optimization, Estimation, and Control*, Taylor and Francis, revised printing.
- [14] G. Buttazzo, D. Percivale, 1983 “On the approximation of the elastic bounce problem on Riemannian manifolds”, Journal of Differential Equations, vol.47, pp.227-245.
- [15] M.K. Çamlıbel, 2001 *Complementarity Methods in the Analysis of Piecewise Linear Dynamical Systems*, PhD thesis ISBN: 90-5668-079 X, Tilburg University, NL.
- [16] D. Cobb, 1983 “Descriptor variable systems and optimal state regulation”, IEEE Transactions on Automatic Control, vol.28, no 5, pp.601-611.
- [17] D. Cobb, C.J. Wang, 2003 “A characterization of bounded-input bounded-output stability for linear time-varying systems with distributional inputs”, SIAM J. Control Optimization, vol.42, no 4, pp.1222-1243.
- [18] R.W. Cottle, J.S. Pang, R.E. Stone, 1992 *The Linear Complementarity Problem*, Academic Press, Computer Science and Scientific Computing.
- [19] A.A. ten Dam, K.F. Dwarshuis, J.C. Willems, 1997 “The contact problem for linear continuous-time dynamical systems: a geometric approach”, IEEE Transactions on Automatic Control, vol.42, no 4, pp.458-472.

- 
- [20] A.A. ten Dam, 1997 *Unilaterally Constrained Dynamical Systems*, Ph.D. Thesis, Rijksuniversiteit Groningen, NL, available at <http://irs.ub.rug.nl/ppn/159407869>
- [21] V.A. Dykhta, 1999 “Optimal pulse control in models of economics and quantum electronics”, *Automation and Remote Control*, vol.60, no 11, pp.1603-1613.
- [22] J. Campos Ferreira, 1997 *Introduction to the Theory of Distributions*, Pitman Monographs and Surveys in Pure and Applied Mathematics, 87, Addison Wesley Longman Limited.
- [23] D. Goeleven, D. Motreanu, Y. Dumont, M. Rochdi, 2003 *Variational and Hemivariational Inequalities: Theory, Methods and Applications; Volume I: Unilateral Analysis and Unilateral Mechanics*, Kluwer Academic Publishers, Nonconvex Optimization and its Applications.
- [24] W.M. Haddad, V. Chellaboina, S.G. Nersesov, 2006 *Impulsive and Hybrid Dynamical Systems. Stability, Dissipativity, and Control*, Princeton Series in Applied Mathematics, Princeton University Press.
- [25] W.E. Hamilton Jr, 1972 “On nonexistence of boundary arcs in control problems with bounded state variables”, *IEEE Transactions on Automatic Control*, vol.17, no 3, pp.338-343.
- [26] R.F. Hartl, S.P. Sethi, R.G. Vickson, 1995 “A survey of the maximum principles for optimal control problems with state constraints”, *SIAM Review*, vol.37, no 2, pp.181-218.
- [27] J.B. Hiriart-Urruty, C. Lemaréchal, 2000 *Fundamentals of Convex Analysis*, Springer Grundlehren Text Editions, Berlin.
- [28] D.H. Jacobson, M.M. Lele, J.L. Speyer, 1971 “New necessary conditions of optimality for control problems with state-variable inequality constraints”, *Journal Math. Anal. Appl.*, vol.35, pp.255-284.
- [29] P. Lancaster, M. Tismenetsky, 1985 *The Theory of Matrices*, 2nd edition, Academic Press London.
- [30] K. Malanovski, H. Maurer, S. Pickenhain, 2004 “Second-order sufficient conditions for state-constrained optimal control problems”, *J. of Optimization Theory and Applications*, vol.123, no 3, pp.595-617.
- [31] K. Malanovski, H. Maurer, 2001 “Sensitivity analysis fo optimal control problems subject to higher order state constraints”, *Annals of Operations Research*, vol.101, pp.43-73.
- [32] M.D.P. Monteiro Marques, 1993 *Differential Inclusions in Nonsmooth Mechanical Problems, Shocks and Dry Friction*, Birkhauser, Progress in Nonlinear Differential Equations and Their Applications.
- [33] H. Maurer, 1977 “On optimal control problems with bounded state variables and control appearing linearly”, *SIAM Journal on Control and Optimization*, vol.15, no 3, pp.345-362.
- [34] J.P. McDanell, W.F. Powers, 1971 “Necessary conditions for joining singular and nonsingular subarcs”, *SIAM J. Control and Optimisation*, vol.9, pp.161-173.
- [35] B.M. Miller, E.Y. Rubinovich, 2003 *Impulsive Control in Continuous and Discrete-Continuous Systems*, Kluwer Academic Publishers, Dordrecht.
- [36] J.J. Moreau, 1989 “An expression of classical mechanics”, *Ann. Inst. H. Poincaré Anal. Non Linéaire*, vol.6, suppl.1-48.
- [37] J.J. Moreau, 1988 “Unilateral contact and dry friction in finite freedom dynamic”, *CISM Courses and Lectures no 302*, International Centre for Mechanical Sciences, J.J. Moreau and P.D. Panagiotopoulos (Eds.), Springer-Verlag, pp.1-82.
- [38] J.J. Moreau, 1988 “Bounded variation in time”, in *Topics in Nonsmooth Dynamics*, J.J. Moreau, P.D. Panagiotopoulos, G. Strang (Eds.), pp.1-74, Birkhauser, Basel.

- [39] J.J. Moreau, 1977 “Evolution problem associated with a moving convex set”, *Journal of Differential Equations*, vol.26, pp.347-374.
- [40] J.M. Murray, 1986 “Existence theorems for optimal control and calculus of variations problems where the state can jump”, *SIAM J. Control and Optimisation*, vol.24, no 3, pp.412-438, May.
- [41] D.O. Norris, 1973 “Nonlinear programming applied to state-constrained optimization problem”, *J. of Mathematical Analysis and Applications*, vol.43, pp.261-272.
- [42] H.J. Oberle, W. Grimm, 2001 “BNDSCO, A program for the Numerical Solution of Optimal Control Problems”, Report no 515 der DFVLR, Deutsche Forschungs und Versuchsanstalt für Luft und Raumfahrt e. V., 1989. *Hamburger Beiträge zur Angewandten Mathematik*.
- [43] Y. Orlov, 2002 “Schwartz’ distributions in nonlinear setting: Applications to differential equations, filtering and optimal control”, *Mathematical Problems in Engineering*, vol.8, no 4-5, pp.367-387.
- [44] L. Paoli, 2002 “A numerical scheme for impact problems with inelastic shocks: a convergence result in the multi-constraint case”, in *Nonsmooth/Nonconvex Mechanics with Applications in Engineering*, Proc. of the Int. Conf. in memoriam of P.D. Panagiotopoulos, 5-6 July 2002, C.C. Baniotopoulos (Ed.), Editions Ziti, Thessaloniki, Greece, pp.269-274.
- [45] F.L. Pereira, G.N. Silva, 2000 “Necessary conditions of optimality for vector-valued impulsive control problems”, *Systems and Control Letters*, vol.40, pp.205-215.
- [46] A.F. Perold, 1978 “Fundamentals of a continuous time simplex method”, *Systems Optimization Laboratory, Dept. of Operations Research, Stanford university*, technical report SOL 78-26, December 1978.
- [47] M.D.R. De Pinho, M.M.A. Ferreira, F.A.C.C. Fontes, 2002 “An Euler-Lagrange inclusion for optimal control problems with state constraints”, *Journal of Dynamical and Control Systems*, vol.8, no 1, pp.23-45, January.
- [48] P. Sannuti, 1983 “Direct singular perturbation analysis of high-gain and cheap control problems”, *Automatica*, vol.19, no 1, pp.41-51.
- [49] P. Sattayatham, 2004 “Strongly nonlinear impulsive evolution equations and optimal control”, *Nonlinear Analysis*, vol.57, pp.1005-1020.
- [50] A.J. van der Schaft, J.M. Schumacher, 2000 *An Introduction to Hybrid Dynamical Systems*, Springer Verlag Lecture Notes in Control and Information Sciences 251.
- [51] W.W. Schmaedeke, 1965 “Optimal control theory for nonlinear vector differential equations containing measures”, *J. SIAM Control, Ser. A*, vol.3, no 2, pp.231-280.
- [52] J.M. Schumacher, 2004 “Complementarity systems in optimization”, *Mathematical Programming, Ser. B*, vol.101, pp.263-295. .
- [53] G.N. Silva, R.B. Vinter, 1997 “Necessary conditions for optimal impulsive control problems”, *SIAM J. Control Optimization*, vol.35, no 6, pp.1829-1846.
- [54] E.D. Sontag, 1998 “*Mathematical Control Theory; Deterministic Finite Dimensional Systems*”, Springer Texts in Applied Mathematics 6, 2nd edition.
- [55] A. Steindl, W. Steiner, H. Troger, 2005 “Optimal control of a retrieval of a tethered subsatellite”, *IUTAM Symp. on Chaotic Dynamics and Control of Systems and Processes in Mechanics*, pp.441-450, G. Rega and F. Vestroni (Eds.), Springer Dordrecht, *Solid Mechanics and its Applications*.

- [56] O. von Stryk, R. Burlisch 1992 “Direct and indirect methods for trajectory optimization”, *Annals of Operations Research*, vol.37, pp.357-373.
- [57] F.A. Valentine, 1937 “The problem of Lagrange with differential inequalities as added side conditions”, in *Contributions to the Calculus of Variations*, Chicago, Chicago University Press, pp.407-448.
- [58] R. Vinter, 2000 *Optimal Control*, Birkhauser, Systems and Control: Foundations and Applications, Boston.

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>3</b>  |
| <b>2</b> | <b>Some properties of the necessary conditions system</b>   | <b>6</b>  |
| <b>3</b> | <b>The higher order Moreau's sweeping process</b>   | <b>8</b>  |
| 3.1      | Presentation of the differential inclusion . . . . .  | 8         |
| 3.2      | Solutions of the necessary conditions IVP . . . . .   | 14        |
| 3.3      | Motivation example . . . . .  | 14        |
| 3.4      | Meaning of the costate $\eta(\cdot)$ jump condition . . . . .   | 15        |
| 3.5      | An extended integral action $I(u)$ . . . . .  | 18        |
| 3.6      | Hamilton's principle of Mechanics . . . . .   | 21        |
| <b>4</b> | <b>Behaviour of the optimal trajectories at junctions with <math>\partial\Phi</math></b>                        | <b>23</b> |
| 4.1      | The admissible domain boundary partitioning . . . . .   | 23        |
| 4.2      | Entry and contact times and states . . . . .  | 25        |
| 4.3      | Exit times and states . . . . .   | 29        |
| 4.4      | Multivariable systems . . . . .   | 30        |
| 4.4.1    | Extension of Proposition 13 . . . . .   | 30        |
| 4.4.2    | Continuous dependence of BVP solutions . . . . .  | 31        |
| <b>5</b> | <b>Boundary arcs, relative degree <math>r^{wu} = 2k + 1, k \geq 0</math></b>                                    | <b>32</b> |
| <b>6</b> | <b>Numerical solution of the BVP</b>  | <b>35</b> |
| <b>7</b> | <b>Conclusions</b>  | <b>35</b> |
| <b>A</b> | <b>Some mathematical definitions</b>  | <b>36</b> |
| <b>B</b> | <b>Example of a function in <math>\mathcal{F}_\infty</math> and distributions in <math>\mathcal{T}_n</math></b> | <b>37</b> |
| <b>C</b> | <b>Hints on the controllability of (2)</b>  | <b>38</b> |



---

Unité de recherche INRIA Rhône-Alpes  
655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Unité de recherche INRIA Futurs : Parc Club Orsay Université - ZAC des Vignes  
4, rue Jacques Monod - 91893 ORSAY Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

Unité de recherche INRIA Sophia Antipolis : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399