# An explicit description of null controllable regions of linear systems with saturating actuators

Tingshu Hu[a,*,1], Zongli Lin[a,1], Li Qiu[b,2]

[a] *Department of Electrical Engineering, University of Virginia, Charlottesville, VA 22903, USA*
[b] *Department of Electrical & Electronic Engineering, Hong Kong University of Science & Technology, Clear Water Bay, Kowloon, Hong Kong*

## Abstract

We give simple exact descriptions of the null controllable regions for general linear systems with saturating actuators. The description is in terms of a set of extremal trajectories of the anti-stable subsystem. For lower order systems or systems with only real eigenvalues, this description is further simplified to result in explicit formulae for the boundaries of the null controllable regions. © 2002 Elsevier Science B.V. All rights reserved.

## 1. Introduction

One of the most fundamental issues associated with the control of a system is its controllability. Since all practical control inputs are bounded (due to actuator saturation), the constrained controllability was formulated earlier than the unconstrained one. While the unconstrained controllability has been well understood for several decades, there have been continual efforts towards full understanding of the constrained controllability (see, e.g., [1–3,6–8,10,14–17] and the references therein).

For a linear system with a constrained input, the null controllable region at a time $T \in (0, \infty)$, denoted as $\mathscr{C}(T)$, is defined to be the set of states that can be steered to the origin in time $T$ with a constrained control. The union of $\mathscr{C}(T)$ for all $T \in (0, \infty)$, denoted as $\mathscr{C}$, is called the null controllable region. In the earlier studies, the null controllable region, also called the controllable set or the reachable set (of the time reversed system), was closely related to the time optimal control (e.g., [2,6,9,11]): for a given initial state $x_0$, the time optimal control problem has a solution if and only if $x_0 \in \mathscr{C}$; If $x_0$ is on the boundary of $\mathscr{C}(T)$, then the minimal time to steer $x_0$ to the origin is $T$. It is well-known that the time optimal controls are bang–bang controls. For discrete-time systems, the time optimal control can be computed through linear programming and $\mathscr{C}(T)$ can be exactly obtained, although the computational burden increases as $T$ increases. Also closely related to the time optimal control is the model predictive control or the receding horizon control. The model predictive

---

* Corresponding author.
 *E-mail addresses:* th7f@virginia.edu (Tingshu Hu), zl5y@virginia.edu (Zongli Lin), eeqiu@ee.ust.hk (Li Qiu).

control has been extensively studied and has found wide applications in slow processes (see, e.g., [12,13] for a survey). The development of the model predictive control also contributes to the characterization of $\mathscr{C}(T)$ for discrete-time systems. On the contrary, for continuous-time systems, since the time optimal control is generally impossible to compute except by numerical approximation, there has been no result on the explicit or analytical characterization of $\mathscr{C}(T)$ or $\mathscr{C}$ of exponentially unstable systems. There are however numerical algorithms available to obtain approximations of $\mathscr{C}$ based on some partial properties about the boundary of $\mathscr{C}$ for second-order systems (e.g., [16,17]). There are also numerical methods for testing if a particular point in the state space is inside $\mathscr{C}$ (e.g., [5]). In this paper, we will focus on the analytic characterization of $\mathscr{C}$ for general linear systems.

We recall that a linear system is said to be anti-stable if all its poles are in the open right-half plane and semi-stable if all its poles are in the closed left-half plane.

For a semi-stable linear system, it is well-known [11,14,15] that the null controllable region is the whole state-space as long as the system is controllable in the usual linear system sense. For a general system with exponentially unstable modes, there exists a nice decomposition result concerning the null controllable region [4]. Suppose such a system is decomposed into the sum of a controllable semi-stable subsystem and an anti-stable subsystem, then the null controllable region of the whole system is the Cartesian product of the null controllable region of the first subsystem, which is its whole state space, and that of the second subsystem, which is a bounded convex open set.

However, little was known about the null controllable region of an anti-stable system. This paper is dedicated to solving this problem. We will give simple exact descriptions of the null controllable region of a general anti-stable linear system with saturating actuators in terms of a set of extremal trajectories of its time reversed system. This set of extremal trajectories is particularly easy to describe for low order systems or systems with only real eigenvalues. For example, for a second-order system, the boundary of its null controllable region is covered by at most two extremal trajectories; and for a third-order system, the set of extremal trajectories can be described in terms of parameters in a real interval.

The remainder of the paper is organized as follows. Section 2 contains some preliminaries and definitions of notation. Section 3 gives a simple exact description of the null controllable regions of anti-stable linear systems with bounded controls. Section 4 draws a brief conclusion to this paper.

## 2. Preliminaries and notation

Consider a linear system

$$\dot{x}(t) = Ax(t) + Bu(t), \tag{1}$$

where $x(t) \in \mathbf{R}^n$ is the state and $u(t) \in \mathbf{R}^m$ is the control. Let

$$\mathscr{U}_a = \{u: u \text{ is measurable and } \|u(t)\|_\infty \leqslant 1, \, \forall t \in \mathbf{R}\}, \tag{2}$$

where $\|u(t)\|_\infty = \max_i |u_i(t)|$. A control signal $u$ is said to be *admissible* if $u \in \mathscr{U}_a$. In this paper, we are interested in the control of system (1) by using admissible controls. Our concern is the set of states that can be steered to the origin by admissible controls.

**Definition 2.1.** A state $x_0$ is said to be null controllable if there exist a $T \in [0, \infty)$ and an admissible control $u$ such that the state trajectory $x(t)$ of the system satisfies $x(0) = x_0$ and $x(T) = 0$. The set of all null controllable states is called the null controllable region of the system and is denoted by $\mathscr{C}$.

With the above definition, we see that $x_0 \in \mathscr{C}$ if and only if there exist $T \in [0, \infty)$ and a $u \in \mathscr{U}_a$ such that

$$0 = x(T) = e^{AT}x_0 + \int_0^T e^{A(T-\tau)}Bu(\tau)\,d\tau = e^{AT}\left(x_0 + \int_0^T e^{-A\tau}Bu(\tau)\,d\tau\right).$$

It follows that

$$\mathscr{C} = \bigcup_{T \in [0,\infty)} \left\{x = -\int_0^T e^{-A\tau}Bu(\tau)\,d\tau : u \in \mathscr{U}_a\right\}. \tag{3}$$

The minus sign "$-$" before the integral can be removed since $\mathscr{U}_a$ is a symmetric set. In what follows we recall from the literature some existing results on the characterization of the null controllable region.

**Proposition 2.1.** *Assume that $(A, B)$ is controllable.*
(a) *If $A$ is semi-stable, then $\mathscr{C} = \mathbf{R}^n$.*
(b) *If $A$ is anti-stable, then $\mathscr{C}$ is a bounded convex open set containing the origin.*
(c) *If*

$$A = \begin{bmatrix} A_1 & 0 \\ 0 & A_2 \end{bmatrix}$$

*with $A_1 \in \mathbf{R}^{n_1 \times n_1}$ anti-stable and $A_2 \in \mathbf{R}^{n_2 \times n_2}$ semi-stable, and $B$ is partitioned as*

$$\begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

*accordingly, then $\mathscr{C} = \mathscr{C}_1 \times \mathbf{R}^{n_2}$ where $\mathscr{C}_1$ is the null controllable region of the anti-stable system $\dot{x}_1(t) = A_1 x_1(t) + B_1 u(t)$.*

Statement (a) is well-known [11,14,15]. Statements (b) and (c) are proven in [4]. Because of this proposition, we can concentrate on the study of null controllable regions of anti-stable systems. For this kind of systems,

$$\bar{\mathscr{C}} = \left\{x = \int_0^\infty e^{-A\tau}Bu(\tau)\,d\tau : u \in \mathscr{U}_a\right\}, \tag{4}$$

where $\bar{\mathscr{C}}$ denotes the closure of $\mathscr{C}$. We will use "$\partial$" to denote the boundary of a set. In [5], a nonlinear programming based algorithm is proposed to test if a point in the state space belongs to $\mathscr{C}$. In Section 3, we will derive a method for explicitly describing $\partial\mathscr{C}$. To this end, we will need some more preliminaries.

If $B = [b_1 \ b_2 \ \dots \ b_m]$ and the null controllable region of the system $\dot{x}(t) = Ax(t) + b_i u_i(t)$ is $\mathscr{C}_i$, $i = 1, \dots, m$, then

$$\mathscr{C} = \sum_{i=1}^m \mathscr{C}_i = \{x_1 + x_2 + \cdots + x_m : x_i \in \mathscr{C}_i, \ i = 1, 2, \dots, m\}. \tag{5}$$

In view of (5) and Proposition 2.1, in the study of the null controllable regions we will assume, without loss of generality, that $(A, B)$ is controllable, $A$ is anti-stable, and $m = 1$. For clarity, we rename $B$ as $b$.

For a general system

$$\dot{x} = f(x, u), \tag{6}$$

its time reversed system is

$$\dot{z} = -f(z, v). \tag{7}$$

It is easy to see that $x(t)$ solves (6) with $x(0) = x_0, x(t_1) = x_1$, and certain $u$ if and only if $z(t) = x(t_1 - t)$ solves (7) with $z(0) = x_1$, $z(t_1) = x_0$, and $v(t) = u(t_1 - t)$. The two systems have the same curves as trajectories, but traverse in opposite directions.

Consider the time reversed system of (1):

$$\dot{z}(t) = -Az(t) - bv(t). \tag{8}$$

**Definition 2.2.** A state $z_f$ is said to be reachable if there exist $T \in [0, \infty)$ and an admissible control $v$ such that the state trajectory $z(t)$ of system (8) satisfies $z(0) = 0$ and $z(T) = z_f$. The set of all reachable states is called the reachable region of system (8) and is denoted by $\mathscr{R}$.

It is known that $\mathscr{C}$ of (1) is the same as $\mathscr{R}$ of (8) (see, e.g., [11]). To avoid confusion, we will continue to use the notation $x$, $u$ and $\mathscr{C}$ for the original system (1), and $z$, $v$ and $\mathscr{R}$ for the time-reversed system (8).

## 3. Null controllable regions

In Section 3.1, we show that the boundary of the null controllable region of a general anti-stable linear system with saturating actuator is composed of a set of extremal trajectories of the time reversed system. The descriptions of this set are further simplified for systems with only real poles and for systems with complex poles in Sections 3.2 and 3.3, respectively.

### 3.1. Description of the null controllable regions

We will characterize the null controllable region $\mathscr{C}$ of system (1) through studying the reachable region $\mathscr{R}$ of its time reversed system (8).

Since $A$ is anti-stable, we have

$$\bar{\mathscr{R}} = \left\{ z = \int_0^\infty \mathrm{e}^{-A\tau} b v(\tau) \, \mathrm{d}\tau : v \in \mathscr{U}_a \right\} = \left\{ z = \int_{-\infty}^0 \mathrm{e}^{A\tau} b v(\tau) \, \mathrm{d}\tau : v \in \mathscr{U}_a \right\}.$$

The change of integration interval from $[0, \infty]$ to $[-\infty, 0]$ is crucial to our development, as will be clear from Eq. (17). Noticing that $\mathrm{e}^{A\tau} = \mathrm{e}^{-A(0-\tau)}$, we see that a point $z$ in $\bar{\mathscr{R}}$ is a state of the time-reversed system (8) at $t = 0$ by applying an admissible control $v$ from $-\infty$ to $0$.

**Theorem 3.1.**

$$\partial\mathscr{R} = \left\{ z = \int_{-\infty}^0 \mathrm{e}^{A\tau} b \, \mathrm{sgn}(c' \mathrm{e}^{A\tau} b) \, \mathrm{d}\tau : c \in \mathbf{R}^n \backslash \{0\} \right\}. \tag{9}$$

$\bar{\mathscr{R}}$ *is strictly convex. Moreover, for each* $z^* \in \partial\mathscr{R}$, *there exists a unique admissible control* $v^*$ *such that*

$$z^* = \int_{-\infty}^0 \mathrm{e}^{A\tau} b v^*(\tau) \, \mathrm{d}\tau. \tag{10}$$

**Remark 3.1.** We give some simple facts about convex sets in this remark. Consider a closed set $S$. If $S$ is convex and $z^* \in \partial S$, then by separation theorem, there exists a hyperplane $c'z = k$ that is tangential to $\partial S$ at $z^*$ and the set $S$ lies completely to one side of the hyperplane, i.e.,

$$c'z \leqslant k = c'z^*, \quad \forall z \in S.$$

A set $S$ is said to be strictly convex if it is convex and for any two points $z_1, z_2 \in \partial S$, $\alpha z_1 + (1-\alpha)z_2 \notin \partial S$ for all $\alpha \in (0, 1)$. This is equivalent to saying that any hyperplane that is tangential to $\partial S$ has only one intersection point with $\partial S$, or, for any $c \neq 0$, there exists a unique $z^* \in \partial S$ such that $c'z^* = \max_{z \in S} c'z$.

**Proof of Theorem 1.** First, the convexity of $\bar{\mathscr{R}}$ can easily be verified by definition. Let $z^* \in \partial \bar{\mathscr{R}}$. Then, there exists a nonzero vector $c \in \mathbf{R}^n$ such that

$$c'z^* = \max_{z \in \bar{\mathscr{R}}} c'z = \max_{v \in \mathscr{U}_a} \int_{-\infty}^{0} c' \mathrm{e}^{A\tau} bv(\tau)\, \mathrm{d}\tau. \tag{11}$$

Since $c \neq 0$ and $(A, b)$ is controllable, $c' \mathrm{e}^{At} b \not\equiv 0$. Since $c' \mathrm{e}^{At} b$ has a finite number of zeros in any finite interval,

$$\mu(\{t: c'\mathrm{e}^{At}b = 0\}) = 0, \tag{12}$$

where $\mu(\cdot)$ denotes the measure of a set.

It is easy to see that

$$v^*(t) = \mathrm{sgn}(c'\mathrm{e}^{At}b)$$

maximizes the right-hand side of (11). We maintain that $v^*$ is the unique optimal solution of (11). To verify this, we need to show that for any $v \in \mathscr{U}_a, v \neq v^*$,

$$\int_{-\infty}^{0} c'\mathrm{e}^{A\tau}bv^*(\tau)\,\mathrm{d}\tau > \int_{-\infty}^{0} c'\mathrm{e}^{A\tau}bv(\tau)\,\mathrm{d}\tau. \tag{13}$$

Since $v \neq v^*$, there are a set $E_1 \subset [-\infty, 0]$ with nonzero measure, i.e., $\mu(E_1) = \delta_1 > 0$, and a number $\varepsilon_1 > 0$ such that

$$|v(t) - v^*(t)| \geqslant \varepsilon_1, \quad \forall t \in E_1.$$

By (12), there exist a set $E \subset E_1$, with $\mu(E) = \delta > 0$, and a positive number $\varepsilon > 0$ such that

$$|c'\mathrm{e}^{At}b| \geqslant \varepsilon, \quad \forall t \in E.$$

Noting that $v \in \mathscr{U}_a$, we have

$$c'\mathrm{e}^{At}b\,(v^*(t) - v(t)) \geqslant 0, \quad \forall t \in [-\infty, 0].$$

It then follows that

$$\int_{-\infty}^{0} c'\mathrm{e}^{A\tau}b(v^*(\tau) - v(\tau))\,\mathrm{d}\tau$$

$$\geqslant \int_{E} c'\mathrm{e}^{A\tau}b(v^*(\tau) - v(t))\,\mathrm{d}\tau = \int_{E} |c'\mathrm{e}^{A\tau}b|\,|v^*(\tau) - v(\tau)|\,\mathrm{d}\tau \geqslant \delta\varepsilon\varepsilon_1 > 0.$$

This shows that $v^*(t)$ is the unique optimal solution of (11) and hence the unique admissible control satisfying

$$z^* = \int_{-\infty}^{0} \mathrm{e}^{A\tau}bv^*(\tau)\,\mathrm{d}\tau. \tag{14}$$

On the other hand, if

$$z^* = \int_{-\infty}^{0} \mathrm{e}^{A\tau}b\,\mathrm{sgn}(c'\mathrm{e}^{A\tau}b)\,\mathrm{d}\tau$$

for some nonzero $c$, then obviously

$$c'z^* = \max_{z \in \bar{\mathscr{R}}} c'z.$$

This shows that $z^* \in \partial \mathscr{R}$ and we have (9).

Since for each $c \neq 0$, the optimal solution $v^*(t)$ and $z^*$ of (11) is unique, we see that $\bar{\mathscr{R}}$ is strictly convex. $\quad\square$

Theorem 3.1 says that for $z^* \in \partial \mathcal{R}$, there is a unique admissible control $v^*$ satisfying (10). From (9), this implies $v^*(t) = \text{sgn}(c'e^{At}b)$ for some $c \neq 0$ (such $c$, $\|c\| = 1$ may be nonunique, where $\| \cdot \|$ is the Euclidean norm). So, if $v$ is an admissible control and there is no $c$ such that $v(t) = \text{sgn}(c'e^{At}b)$ for $t \leq 0$, then

$$\int_{-\infty}^{0} e^{A\tau} b v(\tau) \, d\tau \notin \partial \mathcal{R}$$

and must be in the interior of $\mathcal{R}$.

Since $\text{sgn}(kc'e^{A\tau}b) = \text{sgn}(c'e^{A\tau}b)$ for any positive number $k$, Eq. (9) shows that $\partial \mathcal{R}$ can be determined from the surface of a unit ball in $\mathbf{R}^n$. In what follows, we will simplify (9) and describe $\partial \mathcal{R}$ in terms of a set of trajectories of the time-reversed system (8).

Denote

$$\mathscr{E} := \{v(t) = \text{sgn}(c'e^{At}b), \ t \in \mathbf{R} : c \in \mathbf{R}^n \backslash \{0\}\} \tag{15}$$

and for an admissible control $v$, denote

$$\Phi(t, v) := \int_{-\infty}^{t} e^{-A(t-\tau)} b v(\tau) \, d\tau. \tag{16}$$

Since $A$ is anti-stable, the integral in (16) exists for all $t \in \mathbf{R}$, so $\Phi(t, v)$ is well defined.

If $v(t) = \text{sgn}(c'e^{At}b)$, then

$$\Phi(t, v) = \int_{-\infty}^{t} e^{-A(t-\tau)} b v(\tau) \, d\tau = \int_{-\infty}^{0} e^{A\tau} b \, \text{sgn}(c'e^{At}e^{A\tau}b) \, d\tau \in \partial \mathcal{R} \tag{17}$$

for any $t \in \mathbf{R}$, i.e., $\Phi(t, v)$ lies entirely on $\partial \mathcal{R}$. An admissible control $v$ such that $\Phi(t, v)$ lies entirely on $\partial \mathcal{R}$ is said to be *extremal* and such $\Phi(t, v)$ an *extremal trajectory*. On the other hand, given an admissible control $v(t)$, if there exists no $c$ such that $v(t) = \text{sgn}(c'e^{At}b)$ for all $t \leq 0$, then by Theorem 3.1, $\Phi(0, v) \notin \partial \mathcal{R}$ and must be in the interior of $\mathcal{R}$. By the time invariance property of the system, if there exists no $c$ such that $v(t) = \text{sgn}(c'e^{At}b)$ for all $t \leq t_0$, $\Phi(t, v)$ must be in the interior of $\mathcal{R}$ for all $t \geq t_0$. Consequently, $\mathscr{E}$ is the set of extremal controls.

The following lemma shows that $\partial \mathcal{R}$ is covered by the set of extremal trajectories.

**Lemma 3.1.**

$$\partial \mathcal{R} = \{\Phi(t, v) : t \in \mathbf{R}, v \in \mathscr{E}\}. \tag{18}$$

**Proof.** For any fixed $t \in \mathbf{R}$, it follows from (9) that

$$\partial \mathcal{R} = \left\{ \int_{-\infty}^{t} e^{-A(t-\tau)} b \, \text{sgn}(c'e^{-At}e^{A\tau}b) \, d\tau : c \right\} = \left\{ \int_{-\infty}^{t} e^{-A(t-\tau)} b \, \text{sgn}(c'e^{A\tau}b) \, d\tau : c \right\},$$

i.e., $\partial \mathcal{R} = \{\Phi(t, v) : v \in \mathscr{E}\}$, for any fixed $t \in \mathbf{R}$. So $\partial \mathcal{R}$ can be viewed as the set of extremal trajectories at any frozen time. Now let $t$ vary, then each point on $\partial \mathcal{R}$ moves along a trajectory but the whole set is invariant. So we can also write $\partial \mathcal{R} = \{\Phi(t, v) : v \in \mathscr{E}, \ t \in \mathbf{R}\}$, which is equivalent to (18). $\square$

Unlike (9), Eq. (18) shows that $\partial \mathcal{R}$ is covered by extremal trajectories. It, however, introduces redundancy by repeating the same set $\{\Phi(t, v) : v \in \mathscr{E}\}$ for all $t \in \mathbf{R}$. This redundancy can be removed by a careful examination of the set $\mathscr{E}$. Indeed, the set $\{\Phi(t, v) : t \in \mathbf{R}\}$ can be identical for a class of $v \in \mathscr{E}$.

**Definition 3.1.**
(a) Two extremal controls $v_1, v_2 \in \mathscr{E}$ are said to be equivalent, denoted by $v_1 \sim v_2$, if there exists an $h \in \mathbf{R}$ such that $v_1(t) = v_2(t - h)$ for all $t \in \mathbf{R}$.

(b) Two vectors $c_1, c_2 \in \mathbf{R}^n$ are said to be equivalent, denoted by $c_1 \sim c_2$, if there exist a $k > 0$ and an $h \in \mathbf{R}$ such that $c_1 = k e^{A'h} c_2$.

Noting that a shift in time of the control corresponds to the same shift of the state trajectory, we see that, if $v_1 \sim v_2$, then $\{\Phi(t, v_1): t \in \mathbf{R}\} = \{\Phi(t, v_2): t \in \mathbf{R}\}$; and if $c_1 \sim c_2$, then $\mathrm{sgn}(c_1' e^{At} b) \sim \mathrm{sgn}(c_2' e^{At} b)$.

**Definition 3.2.**
(a) A set $\mathscr{E}_{\min} \subset \mathscr{E}$ is called a minimal representative of $\mathscr{E}$ if for any $v \in \mathscr{E}$, there exists a unique $v_1 \in \mathscr{E}_{\min}$ such that $v \sim v_1$.
(b) A set $M \subset \mathbf{R}^n$ is called a minimal representative of $\mathbf{R}^n$ if for any $c \in \mathbf{R}^n$, there exists a unique $c_1 \in M$ such that $c \sim c_1$.

With this definition, there will be no pair of distinct elements in $\mathscr{E}_{\min}$ or in $M$ that are equivalent. It should be noted that the minimal representative of $\mathscr{E}$ or $\mathbf{R}^n$ is unique up to equivalence and $\mathscr{E}_{\min}$ and $M$ always exist. An immediate consequence of these definitions and Lemma 3.1 is

**Theorem 3.2.** *If $\mathscr{E}_{\min}$ is a minimal representative of $\mathscr{E}$, then*

$$\partial \mathscr{R} = \{\Phi(t, v): t \in \mathbf{R}, v \in \mathscr{E}_{\min}\};$$

*If $M$ is a minimal representative of $\mathbf{R}^n$, then*

$$\partial \mathscr{R} = \{\Phi(t, \mathrm{sgn}(c' e^{At} b)): t \in \mathbf{R}, c \in M \backslash \{0\}\}.$$

It turns out that for some classes of systems, $\mathscr{E}_{\min}$ can be easily described. For second order systems, $\mathscr{E}_{\min}$ contains only one or two elements, so $\partial \mathscr{R}$ can be covered by no more than two trajectories; and for third-order systems, $\mathscr{E}_{\min}$ corresponds to some real intervals. We will see later that for systems of different eigenvalue structures, the descriptions of $\mathscr{E}_{\min}$ can be quite different.

*3.2. Systems with only real eigenvalues*

It follows from, for example, [11, p. 77], that if $A$ has only real eigenvalues and $c \neq 0$, then $c' e^{At} b$ has at most $n - 1$ zeros. This implies that an extremal control can have at most $n - 1$ switches. We will show that the converse is also true.

**Theorem 3.3.** *For system* (8), *assume that $A$ has only real eigenvalues, then*
(a) *an extremal control has at most $n - 1$ switches*;
(b) *any bang–bang control with $n - 1$ or less switches is an extremal control.*

**Proof.** See Appendix A. $\square$

By Theorem 3.3, the set of extremal controls can be described as follows:

$$\mathscr{E} = \left\{ \pm v: v(t) = \begin{cases} 1 & -\infty \leqslant t < t_1, \\ (-1)^i & t_i \leqslant t < t_{i+1}, \quad -\infty < t_1 < t_2 \leqslant \cdots \leqslant t_{n-1} \leqslant \infty \\ (-1)^{n-1}, & t_{n-1} \leqslant t < \infty, \end{cases} \right\} \cup \{v(t) \equiv \pm 1\},$$

where $t_i, i = 1, \ldots, n - 1$, are the switching times. If $v(t)$ has a switch, then the first switch occurs at $t = t_1$. Here we allow $t_i = t_{i+1}$ $(i \neq 1)$ and $t_{n-1} = \infty$, so the above description of $\mathscr{E}$ consists of all bang–bang controls with $n - 1$ or less switches.

To obtain a minimal representative of $\mathscr{E}$, we can simply set $t_1 = 0$, that is,

$$\mathscr{E}_{\min} = \left\{ \pm v\colon v(t) = \begin{cases} 1, & -\infty \le t < t_1, \\ (-1)^i, & t_i \le t < t_{i+1}, \quad 0 = t_1 < t_2 \le \cdots \le t_{n-1} \le \infty \\ (-1)^{n-1}, & t_{n-1} \le t < \infty. \end{cases} \right\} \cup \{v(t) \equiv \pm 1\}.$$

For every $v \in \mathscr{E}_{\min}$, we have $v(t) = 1$ (or $-1$) for all $t < 0$. Hence, for $t \le 0$,

$$\Phi(t,v) = \int_{-\infty}^{t} e^{-A(0-\tau)} b \, d\tau = A^{-1}b \quad (\text{or } -A^{-1}b).$$

Afterwards, $v(t)$ is a bang–bang control with $n - 2$ or less switches. Denote $z_e^+ = -A^{-1}b$ and $z_e^- = A^{-1}b$, then from Theorem 3.2 we have,

**Observation 3.1.** *$\partial \mathscr{R} = \partial \mathscr{C}$ is covered by two bunches of trajectories. The first bunch consists of trajectories of* (8) *whose initial state is $z_e^+$ and the input is a bang–bang control that starts at $t = 0$ with $-1$ and has $n - 2$ or less switches. The second bunch consists of the trajectories of* (8) *whose initial state is $z_e^-$ and the input is a bang–bang control that starts at $t = 0$ with $1$ and has $n - 2$ or less switches.*

Furthermore, $\partial \mathscr{R}$ can be simply described in terms of the open-loop transition matrix. Note that for a fixed $t \ge 0$,

$$\{\Phi(t,v)\colon v \in \mathscr{E}_{\min}\}$$

$$= \left\{ \pm \left[ e^{-At} z_e^+ - \sum_{i=1}^{n-1} \int_{t_i}^{t_{i+1}} e^{-A(t-\tau)} b(-1)^i \, d\tau \right] \colon 0 = t_1 < t_2 < \cdots \le t_{n-1} \le t_n = t \right\} \cup \{\pm z_e^+\}$$

$$= \left\{ \pm \left[ \sum_{i=1}^{n-1} 2(-1)^i e^{-A(t-t_i)} + (-1)^n I \right] A^{-1}b\colon 0 = t_1 < t_2 < \cdots \le t_{n-1} \le t \right\} \cup \{\pm z_e^+\}.$$

Hence,

$$\partial \mathscr{R} = \{\Phi(t,v)\colon t \in \mathbf{R}, v \in \mathscr{E}_{\min}\}$$

$$= \left\{ \pm \left[ \sum_{i=1}^{n-1} 2(-1)^i e^{-A(t-t_i)} + (-1)^n I \right] A^{-1}b\colon 0 = t_1 \le t_2 \le \cdots \le t_{n-1} \le t \le \infty \right\}.$$

Here, we allow $t_1 = t_2$ to include $\pm z_e^+$. For second-order systems,

$$\partial \mathscr{R} = \left\{ \pm \left[ e^{-At} z_e^- - \int_0^t e^{-A(t-\tau)} b \, d\tau \right] \colon t \in [0,\infty] \right\} = \{\pm(-2e^{-At} + I)A^{-1}b\colon t \in [0,\infty]\}. \tag{19}$$
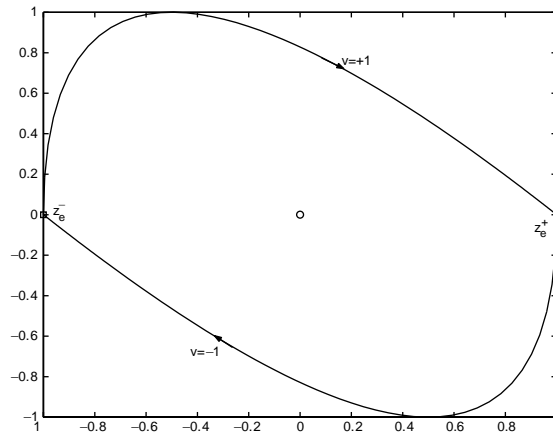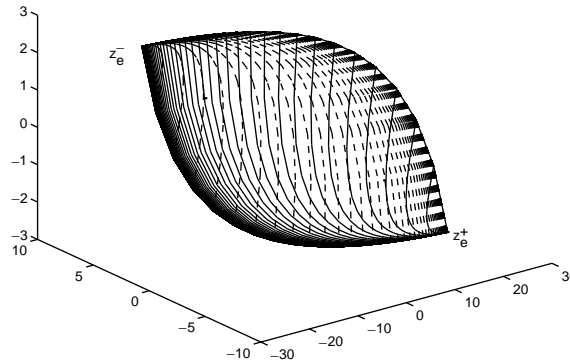
Plotted in Fig. 1 is the $\partial \mathscr{R}$ of a second-order system with

$$A = \begin{bmatrix} 0 & -0.5 \\ 1 & 1.5 \end{bmatrix}, \quad B = \begin{bmatrix} 0 \\ -1 \end{bmatrix}.$$

We see that $\partial \mathscr{R}$ consists of a trajectory from $z_e^+$ to $z_e^-$ under the constant control $v = -1$ and a trajectory from $z_e^-$ to $z_e^+$ under the constant control $v = 1$.

If $n = 3$, then one half of $\partial \mathscr{R} = \partial \mathscr{C}$ can be formed by the trajectories of (8) starting from $z_e^+$ with the control initially being $-1$ and then switching at any time to $1$. So the trajectories go toward $z_e^-$ at first then turn back toward $z_e^+$. The other half is just symmetric to the first half. That is

$$\partial \mathscr{R} = \left\{ \pm \left[ e^{-At} z_e^+ + \int_0^{t_2} e^{-A(t-\tau)} b \, d\tau - \int_{t_2}^t e^{-A(t-\tau)} b \, d\tau \right] \colon 0 \le t_2 \le t \le \infty \right\}. \tag{20}$$

Fig. 1. $\partial\mathcal{R}$ of a second order system.



Fig. 2. $\partial\mathcal{R}$ of a third-order system.

Plotted in Fig. 2 is the $\partial\mathcal{R}$ of a third-order system with

$$A = \begin{bmatrix} 0.2 & 1 & 0 \\ 0 & 0.2 & 0 \\ 0 & 0 & 0.4 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

Since the trajectories of the original system and those of the time-reversed system are the same but traverse in opposite directions, we can also say that $\partial\mathcal{R} = \partial\mathcal{C}$ is covered by a set of trajectories of the original system. While all the trajectories of the time-reversed system start at $z_e^+$ or $z_e^-$ and are very easy to generate by simulation, it is impossible to obtain the same trajectories from the original system. For example, when $n=2$, one half of $\partial\mathcal{R}$ is formed by the trajectory of the time-reversed system that starts at $z_e^-$ under a constant control $v = +1$. The trajectory goes from $z_e^-$ toward $z_e^+$ asymptotically but never reaches $z_e^+$ at a finite time. It seems that if we apply $u = +1$ at $z_e^+$ to the original system, the trajectory will go from $z_e^+$ to $z_e^-$ along the same trajectory of the time-reversed system. However, this is not the case. The trajectory of the original system will stay at $z_e^+$ under the constant control $u = +1$. The boundary $\partial\mathcal{R}$ can only be partially generated from the original system if we know one point on it other than $\pm z_e^+$. But this point is not easy to determine.

### 3.3. Systems with complex eigenvalues

For a system with complex eigenvalues, the minimal representative set $\mathscr{E}_{\min}$ is harder to determine. In what follows, we consider two important cases.

*Case 1. $A \in \mathbf{R}^{2 \times 2}$ has a pair of complex eigenvalues $\alpha \pm j\beta$, $\alpha, \beta > 0$.*

In order to arrive at an explicit formula for $\partial\mathscr{C}$, we need to simplify $c'e^{At}b$. To this end, let $V$ be the nonsingular matrix such that

$$A = V \begin{bmatrix} \alpha & -\beta \\ \beta & \alpha \end{bmatrix} V^{-1}$$

and let

$$\begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = V'c, \quad \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = V^{-1}b,$$

then

$$c'e^{At}b = [c_1 \quad c_2] \begin{bmatrix} \cos(\beta t) & -\sin(\beta t) \\ \sin(\beta t) & \cos(\beta t) \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} e^{\alpha t}$$

$$= [\cos(\beta t) \quad \sin(\beta t)] \begin{bmatrix} b_1 & b_2 \\ -b_2 & b_1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} e^{\alpha t}.$$

Since

$$\begin{bmatrix} b_1 & b_2 \\ -b_2 & b_1 \end{bmatrix}$$

is nonsingular, it follows that

$$\left\{ \begin{bmatrix} b_1 & b_2 \\ -b_2 & b_1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} : c \neq 0 \right\} = \left\{ r \begin{bmatrix} \sin(\theta) \\ \cos(\theta) \end{bmatrix} : r \neq 0, \theta \in [0, 2\pi) \right\}.$$

Hence

$$\{\operatorname{sgn}(c'e^{At}b) : c \neq 0\} = \{\operatorname{sgn}(\sin(\beta t + \theta)) : \theta \in [0, 2\pi)\},$$

and the set of extremal controls is

$$\mathscr{E} = \{v(t) = \operatorname{sgn}(\sin(\beta t + \theta)), t \in \mathbf{R} : \theta \in [0, 2\pi)\}.$$

It is easy to see that

$$\mathscr{E}_{\min} = \{v(t) = \operatorname{sgn}(\sin(\beta t)), t \in \mathbf{R}\}$$

contains only one element. Denote $T_p = \pi/\beta$, then $e^{-AT_p} = -e^{-\alpha T_p}I$. Let

$$z_s^- = (I + e^{-AT_p})^{-1}(I - e^{-AT_p})A^{-1}b = \frac{1 + e^{-\alpha T_p}}{1 - e^{-\alpha T_p}} z_e^-. \tag{21}$$

It can be verified that the extremal trajectory corresponding to $v(t) = \operatorname{sgn}(\sin(\beta t))$ is periodic with period $2T_p$ and,

$$\partial\mathscr{R} = \left\{ \pm \left[ e^{-At} z_s^- - \int_0^t e^{-A(t-\tau)}b \, d\tau \right] : t \in [0, T_p) \right\}$$

$$= \{\pm[e^{-At} z_s^- - (I - e^{-At})A^{-1}b] : t \in [0, T_p)\}. \tag{22}$$

*Case* 2: $A \in \mathbf{R}^{3 \times 3}$ *has eigenvalues* $\alpha \pm j\beta$ *and* $\alpha_1$, *with* $\alpha, \beta, \alpha_1 > 0$.

(a) $\alpha = \alpha_1$. Then similar to Case 1,

$$\mathcal{E} = \{v(t) = \mathrm{sgn}(k + \sin(\beta t + \theta)), t \in \mathbf{R}: k \in \mathbf{R}, \theta \in [0, 2\pi)\}.$$

Since $\mathrm{sgn}(k + \sin(\beta t + \theta))$ is the same for all $k \geqslant 1$ (or $k \leqslant -1$), we have

$$\mathcal{E}_{\min} = \{v(t) = \mathrm{sgn}(k + \sin(\beta t)), t \in \mathbf{R}: k \in [-1, 1]\}.$$

Each $v \in \mathcal{E}_{\min}$ is periodic with period $2T_p$, but the lengths of positive and negative parts vary with $k$. $\Phi(t, v)$ can be easily determined from simulation.

(b) $\alpha_1 \neq \alpha$. Then

$$\mathcal{E} = \{v(t) = \mathrm{sgn}(k_1 e^{(\alpha_1 - \alpha)t} + k_2 \sin(\beta t + \theta)), t \in \mathbf{R}: (k_1, k_2) \neq (0, 0), \theta \in [0, 2\pi)\}.$$

$\mathcal{E}$ can be decomposed as $\mathcal{E} = \mathcal{E}_1 \cup \mathcal{E}_2 \cup \mathcal{E}_3$, where

$$\mathcal{E}_1 = \{v(t) \equiv \pm 1\} \quad (k_2 = 0),$$

$$\mathcal{E}_2 = \{v(t) = \pm \mathrm{sgn}(\sin(\beta t + \theta)): \theta \in [0, 2\pi)\} \quad (k_1 = 0),$$

$$\mathcal{E}_3 = \{v(t) = \pm \mathrm{sgn}(k e^{(\alpha_1 - \alpha)t} + \sin(\beta t + \theta)): k > 0, \theta \in [0, 2\pi)\}.$$

We will show that a minimal representative of $\mathcal{E}_3$ is

$$\mathcal{E}_{3\,\min} = \{v(t) = \pm \mathrm{sgn}(e^{(\alpha_1 - \alpha)t} + \sin(\beta t + \theta)): \theta \in [0, 2\pi)\}. \tag{23}$$

Let

$$v(t) = \mathrm{sgn}(k e^{(\alpha_1 - \alpha)t} + \sin(\beta t + \theta)) \in \mathcal{E}_3.$$

Since $k > 0$, there is a number $h \in \mathbf{R}$ such that $e^{(\alpha_1 - \alpha)h} = k$. So

$$v(t) = \mathrm{sgn}(e^{(\alpha_1 - \alpha)(t + h)} + \sin(\beta(t + h) - \beta h + \theta)) = v_1(t + h)$$

for some $v_1(t) \in \mathcal{E}_{3\,\min}$. On the other hand, given $v_1, v_2 \in \mathcal{E}_{3\,\min}$, suppose

$$v_1(t) = \mathrm{sgn}(e^{(\alpha_1 - \alpha)t} + \sin(\beta t + \theta_1))$$

$$v_2(t) = \mathrm{sgn}(e^{(\alpha_1 - \alpha)t} + \sin(\beta t + \theta_2))$$

and $v_1 \sim v_2$, i.e., $v_1(t) = v_2(t - h)$ for some $h$, then

$$\mathrm{sgn}(e^{(\alpha_1 - \alpha)t} + \sin(\beta t + \theta_1)) = \mathrm{sgn}(e^{(\alpha_1 - \alpha)(t - h)} + \sin(\beta(t - h) + \theta_2)).$$

If $\alpha_1 < \alpha$ (or $\alpha_1 > \alpha$), both $e^{(\alpha_1 - \alpha)t}$ and $e^{(\alpha_1 - \alpha)(t - h)}$ go to zero as $t$ goes to $\infty$ (or $-\infty$). For $v_1(t)$ and $v_2(t - h)$ to change signs at the same time, we must have $\beta t + \theta_1 = \beta(t - h) + \theta_2 + l\pi$, for some integer $l$. Since at any switching time of $v_1(t)$ and $v_2(t)$, $\sin(\beta t + \theta_1) < 0$, $\sin(\beta(t - h) + \theta_2) < 0$, we conclude that $\sin(\beta t + \theta_1) = \sin(\beta(t - h) + \theta_2)$ and $e^{(\alpha_1 - \alpha)t} = e^{(\alpha_1 - \alpha)(t - h)}$. So we get $h = 0, \theta_1 = \theta_2$. These shows that $\mathcal{E}_{3\,\min}$ is a minimal representative of $\mathcal{E}_3$.

The minimal representative of $\mathcal{E}_2$ is the same as $\mathcal{E}_{\min}$ in Case 1. It follows that

$$\mathcal{E}_{\min} = \{v(t) \equiv \pm 1\} \cup \{v(t) = \mathrm{sgn}(\sin(\beta t))\} \cup \mathcal{E}_{3\,\min}.$$

If $\alpha_1 < \alpha$, for each $v \in \mathcal{E}_{3\,\min}$, $v(t) = 1$(or $-1$) for all $t \leqslant 0$, so the corresponding extremal trajectories stay at $z_e^+ = -A^{-1}b$ or $z_e^-$ before $t = 0$. And after some time, they go toward the periodic trajectory since as $t$ goes

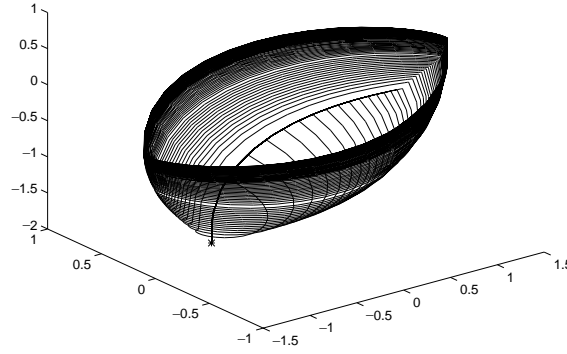Fig. 3. Extremal trajectories on $\partial \mathcal{R}$, $\alpha_1 < \alpha$.

to infinity, $v(t)$ becomes periodic; When $\alpha_1 > \alpha$, for each $v \in \mathcal{E}_{3\,\min}$, $v(t) = 1$(or $-1$) for all $t \geqslant 0$, and the corresponding extremal trajectories start from near periodic and go toward $z_e^+$ or $z_e^-$.

Plotted in Fig. 3 are some extremal trajectories on $\partial \mathcal{R}$ of the time-reversed system (8) with

$$
A = \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0.8 & -2 \\ 0 & 2 & 0.8 \end{bmatrix}; \quad B = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix},
$$

which has two complex poles.

For higher order systems, the relative location of the eigenvalues are more diversified and the analysis will be technically much more involved. It can, however, be expected that in the general case, the number of parameters used to describe $\mathcal{E}_{\min}$ is $n - 2$.

## 4. Conclusions

We gave a clear understanding of the null controllable regions of general linear systems with saturating actuators. We showed that the boundary of the null controllable region of an anti-stable linear system is composed of a set of extremal trajectories of its time-reversed system. The description of the boundary of the null controllable region is further simplified for lower-order systems and systems that have only real eigenvalues.

## Appendix A. Proof of Theorem 3.3

First we present a lemma. Let us use $\mathcal{P}_k$ to denote the set of real polynomials with degree less than integer $k$. The number 0 is considered a polynomial with arbitrary degree or with degree $-1$.

**Lemma A.1.** *Given $N$ positive integers, $k_1, k_2, \ldots, k_N$, define a set of functions*

$$
\mathcal{G}_N := \left\{ g(t) = \sum_{i=1}^{N} e^{a_i t} f_i(t) \colon f_i \in \mathcal{P}_{k_i}, a_i \in \mathbf{R}, g(t) \not\equiv 0 \right\}.
$$

*Then $g(t) \in \mathcal{G}_N$ has at most $\sum_{i=1}^{N} k_i - 1$ zeros.*

**Proof.** We prove this lemma by induction. It is easy to see that the statement is true when $N = 1$. Now assume that it is true when $N$ is replaced by $N - 1$. Let $g(t) \in \mathscr{G}_N$. Suppose on the contrary that $g$ has $\sum_{i=1}^{N} k_i$ or more zeros. Then $\tilde{g}(t) = g(t)e^{-a_N t}$ also has $\sum_{i=1}^{N} k_i$ or more zeros. Therefore, the $k_N$th derivative of $\tilde{g}$,

$$[\tilde{g}(t)]^{(k_N)} = \left[ \sum_{i=1}^{N-1} e^{(a_i - a_N)t} f_i(t) + f_N(t) \right]^{(k_N)} = \left[ \sum_{i=1}^{N-1} e^{(a_i - a_N)t} f_i(t) \right]^{(k_N)} \in \mathscr{G}_{N-1},$$

has at least $\sum_{i=1}^{N-1} k_i$ zeros, which is a contradiction. $\quad\square$

**Proof of Theorem 3.3.** The proof of (a) was sketched in [11, p. 77], where an additional assumption of normality was required. This assumption is satisfied for system (8) since it is single input and $(A, b)$ is controllable. To show (b), assume that $A$ has $N$ distinct real eigenvalues $\lambda_i, i = 1, 2, \ldots, N$, each with a multiplicity of $k_i$ ($\sum_{i=1}^{N} k_i = n$). It is well-known that $c'e^{At}b = \sum_{i=1}^{N} e^{\lambda_i t} f_i(t)$ for some $f_i \in \mathscr{P}_{k_i}$. If $c \neq 0$, then $c'e^{At}b \not\equiv 0$ by the controllability of $(A, b)$. (Thus (a) follows from Lemma A.1). To complete the proof of (b), we first show that any bang–bang control $v$ with $n - 1$ switches is an extremal control.

Let $t_1, t_2, \ldots, t_{n-1} \in \mathbf{R}$ be distinct switching times of $v$. From the following $n - 1$ linear equations

$$c'e^{At_i}b = 0, \quad i = 1, 2, \ldots, n - 1$$

at least one nonzero vector $c \in \mathbf{R}^n$ can be solved. With any such $c$, (a) implies $g(t) = c'e^{At}b \not\equiv 0$ has no other zeros than the $n - 1$ zeros at $t_i, i = 1, 2, \ldots, n - 1$.

Now the question is whether $g(t)$ indeed changes signs at each $t_i$. If it does, then $v(t) = \text{sgn}(c'e^{At}b)$ (or $\text{sgn}(-c'e^{At}b)$) and $v$ is an extremal control.

We now show that $g$ does change signs at each $t_i$. If $g$ does not change sign at a certain $t_i$, then $g(t)$ must have a local extremum at $t_i$, so $\dot{g}(t_i) = 0$. We argue that there is at most one $t_i$ such that $\dot{g}(t_i) = 0$, otherwise $\dot{g}$ will have at least $n$ zeros, counting the at least $n - 2$ ones lying within the intervals $(t_i, t_{i+1})$, which is impossible by Lemma A.1, since $\dot{g}$ has the same structure as $g$.

We further conclude that $g$, however, cannot have a local extremum at any of these $t_i$'s.

Let $g(t) = \sum_{i=1}^{N} e^{\lambda_i t} f_i(t)$. Assume without loss of generality that $f_N(t) \not\equiv 0$. Suppose on the contrary that $g$ has a local minimum (or maximum) at $t_1$, then $\tilde{g}(t) = g(t)e^{-\lambda_N t}$ also has a local minimum (or maximum) at $t_1$, furthermore, $\tilde{g}(t_i) = 0, \dot{\tilde{g}}(t_i) \neq 0, i = 2, 3, \ldots, n - 1$. Hence, there exists an $\varepsilon > 0$ (or $\varepsilon < 0$) such that $\tilde{g}(t) - \varepsilon = \sum_{i=1}^{N-1} e^{(\lambda_i - \lambda_N)t} f_i(t) + f_N(t) - \varepsilon$ has $n$ zeros, which contradicts with Lemma A.1. Therefore, $g$ changes signs at all $t_i$. This shows that $v(t) = \text{sgn}(c'e^{At}b)$ (or $\text{sgn}(-c'e^{At}b)$) is an extremal control.

Now consider the case that $v$ has less than $n - 1$ switches, say $n - 1 - j$ switches, $t_i, i = 1, 2, \ldots, n - 1 - j$. For simplicity and without loss of generality, assume that $A$ is in the Jordan canonical form (the state transformation matrix can be absorbed in $c'$ and $b$. Partition $A, b$ as

$$A = \begin{bmatrix} \star & \star \\ 0 & A_1 \end{bmatrix}, \quad b = \begin{bmatrix} \star \\ b_1 \end{bmatrix}$$

where $A_1$ is of size $n - j$. It is easy to see that $A_1$ is also of the Jordan canonical form and $(A_1, b_1)$ is controllable. Furthermore,

$$e^{At} = \begin{bmatrix} \star & \star \\ 0 & e^{A_1 t} \end{bmatrix}.$$

Partition $c = [0 \quad c_1']'$ accordingly, then $c'e^{At}b = c_1'e^{A_1 t}b_1$. By the foregoing proof for the full dimensional case, we see that there exists $c_1$ such that $v(t) = \text{sgn}(c_1'e^{A_1 t}b_1)$ is a bang–bang control with switching times exactly at $t_i, i = 1, 2, \ldots, n - 1 - j$.

Therefore, we conclude that any bang–bang control with less than $n - 1$ switches is also extremal. $\quad\square$

# References

[1] D.S. Bernstein, A.N. Michel, A chronological bibliography on saturating actuators, Internat. J. Robust Nonlinear Control 5 (1995) 375–380.

[2] C.A. Desoer, J. Wing, An optimal strategy for a saturating sampled data system, IRE Trans. Automat. Control AC-6 (1961) 5–15.

[3] M.E. Fisher, J.E. Gayek, Estimating reachable sets for two-dimensional linear discrete systems, J. Opt. Theory Appl. 56 (1987) 67–88.

[4] O. Hájek, Control Theory in the Plane, Springer, Berlin, 1991.

[5] T. Hu, L. Qiu, Controllable regions of linear systems with bounded inputs, Systems Control Lett. 33 (1998) 55–61.

[6] R.E. Kalman, Optimal nonlinear control of saturating systems by intermittent action, IRE Wescon Convention Record Part 4, 1957, pp. 130–135.

[7] S.S. Keerthi, E.G. Gilbert, Computation of minimum-time feedback control laws for discrete-time systems with state-control constraints, IEEE trans. Automat. Control 32 (1987) 432–435.

[8] J.B. Lasserre, Reachable, controllable sets and stabilizing control of constrained linear systems, Automatica 29 (1993) 531–536.

[9] E. Lee, L. Markus, Foundations of Optimal Control, Wiley, New York, 1967.

[10] J.N. Lin, Determination of reachable set for a linear discrete system, IEEE Trans. Automat. Control AC-15 (1970) 339–342.

[11] J. Macki, M. Strauss, Introduction to Optimal Control, Springer, Berlin, 1982.

[12] D.Q. Mayne, Control of constrained dynamic systems, European J. Control 7 (2001) 87–99.

[13] D.Q. Mayne, J.B. Rawlings, C.V. Rao, P.O.M. Scokaert, Constrained model predictive control: stability and optimality, Automatica 36 (2000) 789–814.

[14] W.E. Schmitendorf, B.R. Barmish, Null controllability of linear systems with constrained controls, SIAM J. Control Optim. 18 (1980) 327–345.

[15] E.D. Sontag, An algebraic approach to bounded controllability of linear systems, Internat. J. Control 39 (1984) 181–188.

[16] J. Stephan, M. Bodson, J. Lehocsky, Properties of recoverable sets for input and state constrained systems, Proceedings of the American Control Conference, Seattle, 1995, pp. 3912–3913.

[17] J. Stephan, M. Bodson, J. Lehocsky, Calculation of recoverable sets for 2-dimensional systems with input and state constraints, Proceedings of the Conference on Decision and Control, New Orleans, 1995, pp. 631–636.