

Global stabilization of linear discrete-time systems with bounded feedback

Yudi Yang^{a,*}, Eduardo D. Sontag^{b,1}, Hector J. Sussmann^{b,2}

^a IBM, MD60, 500 Mamaroneck Ave., Harrison, NY 10528, USA

^b Department of Mathematics, Rutgers University, New Brunswick, NJ 08903, USA

Received 10 September 1996; received in revised form 19 November 1996; accepted 20 February 1997

Abstract

This paper deals with the problem of global stabilization of linear discrete time systems by means of bounded feedback laws. The main result proved is an analog of one proved for the continuous time case by the authors, and shows that such stabilization is possible if and only if the system is stabilizable with arbitrary controls and the transition matrix has spectral radius less than or equal to one. The proof provides in principle an algorithm for the construction of such feedback laws, which can be implemented either as cascades or as parallel connections (“single hidden layer neural networks”) of simple saturation functions. © 1997 Elsevier Science B.V.

Keywords: Linear discrete-time systems; Saturated feedback; Global stabilization

1. Introduction

This paper is concerned with the global stabilization to the origin $x = 0$ of the state $x(t)$ of a linear discrete-time system

$$\Sigma: \quad x(t+1) = Ax(t) + Bu(t), \quad (1.1)$$

when the control values $u(t)$ are constrained to lie in a bounded subset \mathcal{U} of \mathbb{R}^m which contains zero in its interior. (As usual, $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$.) The study of stabilization under such constraints is not only a natural mathematical problem, but also arises often in many applied areas.

The *open loop* question is well-understood. Call a system (1.1) *asymptotically null control-*

lable with bounded controls (ANCBC) if there is some \mathcal{U} with the above properties such that, for each initial state $x(0) \in \mathbb{R}^n$, there exists a sequence $u(\cdot) = u(0), u(1), \dots$, with all values $u(t) \in \mathcal{U}$, which steers the solution $x(t)$ asymptotically to the origin, that is, so that the solution of (1.1) converges to zero. (It turns out, and in fact follows also from the results to be given, that if this property holds for some such \mathcal{U} then it also holds for every \mathcal{U} which contains the origin in its interior.) Now, it is known (cf. [3]) that a system is ANCBC if and only if (1) the pair (A, B) is stabilizable or “asycontrollable” in the usual unconstrained sense (equivalently, the rank of $[\lambda I - A, B]$ is n for all complex λ with $|\lambda| \geq 1$, cf. e.g. [5, Exercise 4.4.7]) and (2) the spectral radius of A is less or equal to one. This provides an elegant algebraic solution of the open-loop question. What is proved in this paper is that, under exactly the same conditions, there is in fact a simple *feedback synthesis* that achieves closed-loop stabilization. The

* Corresponding author. E-mail: yudiy@vnet.ibm.com.

¹ Supported in part by US Air Force Grant AFOSR-94-0343.

² Supported in part by NSF Grant DMS-8902994 and by US Air Force Grant AFOSR-94-0343.

feedback laws that achieve this goal can be optionally of a form that involves series (cascade) connections of linear functions and saturation devices or, alternatively, if desired, of a parallel form involving such saturations.

The results in this paper are in no way surprising or unexpected, since they are closely analogous to similar results presented by the authors, and by A. Teel, for continuous time systems, in the sequence of papers [6–9]. Although the organization of the current work is tightly patterned after that of [8], and many of the arguments – but not all – are, conceptually, straightforward generalizations of the corresponding arguments in that continuous time paper, it seems appropriate to present the discrete time results, because there are many technical estimates that have to be carefully established for this particular case and which are not totally obvious.

To simplify the presentation, we present a result that is weaker than the complete analog of the result in [8]: we restrict the saturations to be used when implementing feedback laws to be of a special kind, while in the continuous time result we showed that rather arbitrary saturation functions could be used as the building blocks. However, for applications, it would appear that our choice of primitive saturation functions is sufficient.

The organization of the paper is as follows. In Section 2 we introduce notations as well as state the main results; this is almost a verbatim translation of the corresponding continuous time material. In Section 3 we provide a technical lemma on changing to a suitable canonical form, while another technical lemma, dealing with an ultimate boundedness result, is given in Section 4. The result in this section is not proved in a manner analogous to the corresponding result in [8], since doing so would require first obtaining the discrete time analogues of the finite gain results given in [2]; a direct proof is given instead. Finally, in Section 5 we give the proof of the main result, with arguments that are again quite similar to those used for continuous time.

The results in this paper are extracted from Chapter 6 of the doctoral thesis [11]. Other references to closely related problems are [1, 10]: the former gave a result on semi-global stabilizability (feedback laws that are guaranteed to work on any given compact, though not necessarily globally) using a simple saturated linear feedback, and the latter provided partial results on global stabilizability for some special systems.

2. Statement of the main results

We start by introducing notations for the classes of functions which will be used to describe the feedback laws to be synthesized. (These definitions and notations are essentially the same as in the paper [8], except that they are built out of a special saturation function, defined next, instead of the far more general saturations used in that paper.) We let \mathcal{S} consist of the saturations at various levels $\delta > 0$, that is, the set of all functions $\mathbb{R} \rightarrow \mathbb{R}$ of the type

$$\sigma(s) = \delta \operatorname{sat}(s/\delta),$$

where $\delta > 0$, and

$$\operatorname{sat}(s) = \operatorname{sign}(s) \cdot \min\{|s|, 1\}.$$

Next we introduce, for each nonnegative integer k and each finite sequence $\sigma = (\sigma_1, \dots, \sigma_k)$ of functions in \mathcal{S} , a set of functions from \mathbb{R}^n to \mathbb{R} , denoted $\mathcal{F}_n(\sigma)$, which consist of “cascades” of saturations. By induction on k , we define these sets as follows:

- when $k = 0$ (which we can interpret as corresponding to the “empty sequence” σ), $\mathcal{F}_n(\sigma)$ consists of just one element, namely, the zero function from \mathbb{R}^n to \mathbb{R} ;
- when $k = 1$, we define $\mathcal{F}_n(\sigma_1)$ as the set of all the functions $h: \mathbb{R}^n \rightarrow \mathbb{R}$ of the form $h(x) = \sigma_1(g(x))$, where $g: \mathbb{R}^n \rightarrow \mathbb{R}$ is a linear function;
- for $k > 1$, $\mathcal{F}_n(\sigma_1, \dots, \sigma_k)$ is the set of all those functions $h: \mathbb{R}^n \rightarrow \mathbb{R}$ that are of the form $h(x) = \sigma_k(f(x) + cg(x))$, for some linear $f: \mathbb{R}^n \rightarrow \mathbb{R}$, some $g \in \mathcal{F}_n(\sigma_1, \dots, \sigma_{k-1})$, and some $c \geq 0$.

A second family of sets of functions $\mathcal{G}_n(\sigma)$, corresponding to “parallel combinations” of saturations, is defined as follows: for each nonnegative integer k and each finite sequence $\sigma = (\sigma_1, \dots, \sigma_k)$ of functions in \mathcal{S} , $\mathcal{G}_n(\sigma)$ is the class of functions $h: \mathbb{R}^n \rightarrow \mathbb{R}$ of the form

$$h(x) = \sigma_1(f_1(x)) + \sigma_2(f_2(x)) + \dots + \sigma_k(f_k(x)),$$

where f_1, \dots, f_k are linear functions.

Finally, given any m -tuple $\mathbf{l} = (l^1, \dots, l^m)$ of non-negative integers, and any finite sequence $\sigma = (\sigma_1, \dots, \sigma_{|\mathbf{l}|})$ of functions in \mathcal{S} , where $|\mathbf{l}| = l^1 + \dots + l^m$, we define the following classes of vector functions built out of the classes of scalar functions which were just defined. We write in partitioned form $\sigma = (\sigma_1^1, \dots, \sigma_{l^1}^1, \dots, \sigma_1^m, \dots, \sigma_{l^m}^m)$, and let $\mathcal{F}_n^{\mathbf{l}}(\sigma)$ (respectively, $\mathcal{G}_n^{\mathbf{l}}(\sigma)$) be the set of all functions $h: \mathbb{R}^n \rightarrow \mathbb{R}^m$ that are of the form (h_1, \dots, h_m) , where

$h_i \in \mathcal{F}_n(\sigma_1^i, \dots, \sigma_{\mu_i}^i)$ (respectively, $h_i \in \mathcal{G}_n(\sigma_1^i, \dots, \sigma_{\mu_i}^i)$) for $i = 1, 2, \dots, m$. (So $\mathcal{F}_n^l(\sigma) = \mathcal{F}_n(\sigma)$, $\mathcal{G}_n^l(\sigma) = \mathcal{G}_n(\sigma)$ when $m = 1$.) For a sequence of saturations as here, we denote as $\|\sigma\|$ the maximum bound (the “ δ ”s in their definition) among all the σ_i ’s. (We use $|x|$ for the Euclidean norm of a vector x , in order to avoid confusion.)

Let $\delta > 0$. We say that a function $\xi: \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}^n$ is eventually bounded by δ (and write $|\xi| \leq_{ev} \delta$), if there exists $T > 0$ such that $|\xi(t)| \leq \delta$ for all $t \geq T$. Given an n -dimensional system $\mathcal{E}: x(t+1) = f(x(t))$, we say that \mathcal{E} is IICS (integrable-input converging-state) if, whenever $\{e(t)\}_0^\infty \in l_1$, every solution $t \rightarrow x(t)$ of $x(t+1) = f(x(t)) + e(t)$ converges to zero as $t \rightarrow \infty$. (We need this concept in order to be able to state a result which can be used in an induction proof.) For a system $x(t+1) = f(x(t), u(t))$, we say that a feedback $u = k(x)$ is stabilizing if 0 is a globally asymptotically stable equilibrium of the closed-loop system $x(t+1) = f(x(t), k(x(t)))$. If, in addition, this closed-loop system is IICS, then we will say that k is IICS-stabilizing.

For an $n \times n$ real matrix A , let $N(A)$ be the number of eigenvalues z of A such that $|z| = 1$ and $\text{Im } z \geq 0$, counting multiplicities.

This is the explicit version of our main result:

Theorem 1. Assume that Σ is an ANCBC linear system $x(t+1) = Ax(t) + Bu(t)$ with state space \mathbb{R}^n and input space \mathbb{R}^m . Let $N = N(A)$. Then, for every $\varepsilon > 0$, there exist a sequence $\sigma = (\sigma_1, \dots, \sigma_N)$ of functions belonging to \mathcal{S} with $\|\sigma\| \leq \varepsilon$ and an m -tuple $\mathbf{l} = (l^1, \dots, l^m)$ of nonnegative integers such that $|\mathbf{l}| = l^1 + \dots + l^m = N$, for which there are IICS-stabilizing feedbacks

$$u = k_{\mathcal{F}}(x), \tag{2.1}$$

$$u = k_{\mathcal{G}}(x) \tag{2.2}$$

such that $k_{\mathcal{F}} \in \mathcal{F}_n^l(\sigma)$, $k_{\mathcal{G}} \in \mathcal{G}_n^l(\sigma)$.

We will say that (2.1), (2.2) are “feedbacks of Type \mathcal{F} ” and “of Type \mathcal{G} ”, respectively.

A linear discrete-time system Σ is bounded feedback stabilizable (BFS) if there exists a bounded locally Lipschitz feedback k that stabilizes Σ . A linear discrete-time system Σ is small feedback stabilizable (SFS) if for every $\varepsilon > 0$ there exists a stabilizing feedback k for Σ such that $|k(x)| \leq \varepsilon$ for all x .

The following is an easy corollary of Theorem 1, and conveys the main conclusions in a simplified form.

Theorem 2. Let Σ be a linear discrete-time system. Then the following conditions are equivalent:

1. Σ is SFS,
2. Σ is BFS,
3. Σ is ANCBC.

Note that the implication $3 \Rightarrow 1$ follows from Theorem 1, while $1 \Rightarrow 2$ and $2 \Rightarrow 3$ are trivially true.

3. A useful change of coordinates

In this section we present a technical lemma which is needed in the proof of Theorem 1. It follows the lines of the analogous continuous-time result, Lemma 3.1, in [8].

Lemma 3.1. Consider an n -dimensional linear single-input system

$$\Sigma: \quad x(t+1) = Ax(t) + bu(t). \tag{3.1}$$

Suppose that (A, b) is a controllable pair and that all the eigenvalues of A have magnitude 1.

- (i) If $\lambda = 1$ or $\lambda = -1$ is an eigenvalue of A , then there is a linear change of coordinates $Tx = (y_1, \dots, y_n)' = (\bar{y}', y_n)'$ of \mathbb{R}^n that transforms Σ into the form

$$\begin{aligned} \bar{y}(t+1) &= A_1 \bar{y}(t) + b_1(y_n(t) + u(t)), \\ y_n(t+1) &= \lambda(y_n(t) + u(t)), \end{aligned} \tag{3.2}$$

where the pair (A_1, b_1) is controllable and y_n is a scalar variable.

- (ii) If A has an eigenvalue of the form $\alpha + \beta i$, with $\beta \neq 0$, then there is a linear change of coordinates $Tx = (y_1, \dots, y_n)' = (\bar{y}', y_{n-1}, y_n)'$ of \mathbb{R}^n that transforms Σ into the form

$$\begin{aligned} \bar{y}(t+1) &= A_1 \bar{y}(t) + b_1(y_n(t) + u(t)), \\ y_{n-1}(t+1) &= \alpha y_{n-1}(t) - \beta(y_n(t) + u(t)), \\ y_n(t+1) &= \beta y_{n-1}(t) + \alpha(y_n(t) + u(t)), \end{aligned} \tag{3.3}$$

where the pair (A_1, b_1) is controllable and y_{n-1}, y_n are scalar variables.

Proof. We first prove (i). If $\lambda = 1$ or $\lambda = -1$ is an eigenvalue of A , then there exists a nonzero n -dimensional row vector v such that $vA = \lambda v$. It follows

from the Hautus condition for controllability (see e.g. [5, Lemma 3.3.7]) that $vb \neq 0$; thus, we may normalize v so that $vb = \lambda$, which we assume from now on. We apply a preliminary linear change of coordinates $Tx = (\bar{z}', z_n)'$, where the matrix T is picked so that $z_n = vx$; in the new coordinates, the system equations take the following block form:

$$\bar{z}(t+1) = A_1 \bar{z}(t) + z_n(t) \tilde{b}_1 + u(t) \tilde{b}_2,$$

$$z_n(t+1) = \lambda z_n(t) + \lambda u(t).$$

We now apply a second coordinate change, letting $\bar{y} = \bar{z} + z_n \tilde{b}_3$, $y_n = z_n$, where the vector \tilde{b}_3 will be specified below. The system equations now become:

$$\begin{aligned} \bar{y}(t+1) &= A_1 \bar{y}(t) + y_n(t) (\tilde{b}_1 + (\lambda I - A_1) \tilde{b}_3) \\ &\quad + u(t) (\tilde{b}_2 + \lambda \tilde{b}_3), \end{aligned}$$

$$y_n(t+1) = \lambda (y_n(t) + u(t)).$$

We pick \tilde{b}_3 to be any solution of $\tilde{b}_2 + \lambda \tilde{b}_3 = \tilde{b}_1 + (\lambda I - A_1) \tilde{b}_3$, i.e., $A_1 \tilde{b}_3 = \tilde{b}_1 - \tilde{b}_2$. (This is possible because A_1 is nonsingular; note that all its eigenvalues are in the unit circle.) With $b_1 = \tilde{b}_1 + (\lambda I - A_1) \tilde{b}_3$, the equations have the desired form (3.2).

We next prove part (ii). Let $\lambda = \alpha + \beta i$, $\beta \neq 0$, be an eigenvalue of A . Let v be a left eigenvector associated to λ , i.e. $vA = \lambda v$, $v \neq 0$. Again by Hautus' condition, $vb \neq 0$. Write $v = v_1 + iv_2$, with v_1 and v_2 real. We may assume that $v_1 b \neq 0$ (otherwise, use iv in place of v), and, hence, normalizing, that $v_1 b = -\beta$. Let $\tau = v_2 b$ and consider the following real 2×2 matrix:

$$P = \frac{1}{\beta^2 + \tau^2} \begin{pmatrix} \beta^2 + \alpha\tau & \beta(\alpha - \tau) \\ \beta(\tau - \alpha) & \beta^2 + \alpha\tau \end{pmatrix}.$$

Make a linear change of coordinates $Tx = (\bar{z}', z_{n-1}, z_n)'$ so that $(z_{n-1}, z_n)' \equiv P(v_1 x, v_2 x)'$. In the new coordinates, the system equations become:

$$\begin{aligned} \bar{z}(t+1) &= A_1 \bar{z}(t) + z_{n-1}(t) \tilde{b}_1 \\ &\quad + z_n(t) \tilde{b}_2 + u(t) \tilde{b}_3, \\ z_{n-1}(t+1) &= \alpha z_{n-1}(t) - \beta(z_n(t) + u(t)), \\ z_n(t+1) &= \beta z_{n-1}(t) + \alpha(z_n(t) + u(t)), \end{aligned} \quad (3.4)$$

and every eigenvalue of A_1 has magnitude 1. Finally, we change coordinates once more, by letting $\bar{y} = \bar{z} + z_{n-1} \tilde{b}_4 + z_n \tilde{b}_5$, $y_{n-1} = z_{n-1}$, $y_n = z_n$, where the vectors \tilde{b}_4 , \tilde{b}_5 will be chosen below. Then the last two equations of (3.4) are as desired, and the equation of

\bar{y} becomes

$$\begin{aligned} \bar{y}(t+1) &= A_1 \bar{y}(t) + y_{n-1}(t) (\tilde{b}_1 - A_1 \tilde{b}_4 \\ &\quad + \alpha \tilde{b}_4 + \beta \tilde{b}_5) \\ &\quad + y_n(t) (\tilde{b}_2 - A_1 \tilde{b}_5 + \alpha \tilde{b}_5 - \beta \tilde{b}_4) \\ &\quad + u(t) (\tilde{b}_3 - \beta \tilde{b}_4 + \alpha \tilde{b}_5). \end{aligned} \quad (3.5)$$

If we could choose \tilde{b}_4 , \tilde{b}_5 such that

$$\tilde{b}_1 - A_1 \tilde{b}_4 + \alpha \tilde{b}_4 + \beta \tilde{b}_5 = 0 \quad (3.6)$$

and

$$\tilde{b}_3 - \beta \tilde{b}_4 + \alpha \tilde{b}_5 = \tilde{b}_2 - A_1 \tilde{b}_5 - \beta \tilde{b}_4 + \alpha \tilde{b}_5, \quad (3.7)$$

then we could let

$$b_1 = \tilde{b}_2 - A_1 \tilde{b}_5 - \beta \tilde{b}_4 + \alpha \tilde{b}_5 \quad (3.8)$$

and the system equations would become (3.3) as desired. To prove the existence of \tilde{b}_4 and \tilde{b}_5 , we rewrite (3.7) as $A_1 \tilde{b}_5 = \tilde{b}_2 - \tilde{b}_3$, from which we get \tilde{b}_5 because A_1 is nonsingular. Then from (3.6), we have $(A_1 - \alpha I) \tilde{b}_4 = \tilde{b}_1 + \beta \tilde{b}_5$. Since the eigenvalues of A_1 have magnitude 1 and $\alpha \neq \pm 1$, the matrix $A_1 - \alpha I$ is nonsingular, and so \tilde{b}_4 exists as well. \square

4. An ultimate boundedness result

The main technical lemma needed for the proof of our main result is given in this section. Though its conclusions are similar to Lemma 3.2 in [8], the proof that we provide is quite different. Because we restricted attention to a special type of saturation functions, the argument is substantially simpler than that in the cited paper.

Lemma 4.1. *Let a, b be two real constants such that $a^2 + b^2 = 1$ and $b \neq 0$. Let $e_j = (e_j(0), e_j(1), e_j(2), \dots)$, $j = 1, 2$, be two elements of l_1 . Pick any $\delta > 0$ and any $\varepsilon \in (0, \delta/4)$. Suppose that $v: \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ is so that $|v| \leq_{\text{ev}} \varepsilon$. Then, if $\gamma = (x(\cdot), y(\cdot)): \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}^2$ is any solution of the system*

$$\begin{aligned} x(t+1) &= ax(t) - by(t) + bu(t) + e_1(t), \\ y(t+1) &= bx(t) + ay(t) - au(t) + e_2(t), \end{aligned} \quad (4.1)$$

where

$$u(t) = \sigma(y(t) + \xi v(t)) + \eta v(t), \quad (4.2)$$

and $\xi + \eta = 1$, $\xi, \eta \geq 0$, and $\sigma(s) = \delta \text{sat}(s/\delta)$, it follows that

$$\limsup_{t \rightarrow +\infty} |\gamma(t)| < r = \frac{1}{|b|}(7|a| + 4)\varepsilon + 7\varepsilon. \quad (4.3)$$

Proof. Without loss of generality, we assume $b > 0$ (if this were not the case, the result can be proved for the negatives $-a, -b$, etc., substituted for the original data; note that the assumptions hold for these, and the conclusions involve only absolute values). Let $\theta = \arctan(b/a)$, $0 < \theta < \pi$, if $a \neq 0$, and $\theta = \pi/2$ if $a = 0$. Then $a + ib = e^{i\theta}$. Let $z(t) = x(t) + iy(t)$, $e(t) = e_1(t) + ie_2(t)$. Then

$$z(t + 1) = e^{i\theta}(z(t) - iu(t)) + e(t). \quad (4.4)$$

Again, without loss of generality, we assume that $\|e\|_1 < \varepsilon$, (otherwise we can find $T > 0$ such that $\sum_{t \geq T} |e(t)| < \varepsilon$, and then we only need to consider the solution for $t \geq T$). Similarly, we assume $|v(t)| \leq \varepsilon$ for all t . So

$$\begin{aligned} |z(t + 1)| &\leq |z(t) - iu(t)| + |e(t)| \\ &= \sqrt{x(t)^2 + (y(t) - u(t))^2} + |e(t)| \\ &= \sqrt{|z(t)|^2 - u(t)(2y(t) - u(t))} + |e(t)| \\ &= |z(t)| + w(t) + |e(t)|, \end{aligned} \quad (4.5)$$

where

$$w(t) = \frac{-u(t)(2y(t) - u(t))}{|z(t)| + \sqrt{|z(t)|^2 - u(t)(2y(t) - u(t))}}. \quad (4.6)$$

If t is so that $|y(t)| \geq 3\varepsilon$, then from (4.2) it follows that

$$2\varepsilon \leq |u(t)| \leq \frac{4}{3}|y(t)|.$$

and $u(t)$ has the same sign as $y(t)$. So

$$w(t) \leq -\frac{2\varepsilon \cdot \frac{2}{3}|y(t)|}{2|z(t)|} \leq -\frac{2\varepsilon^2}{|z(t)|}.$$

Thus, from (4.5), we have

$$|z(t + 1)| \leq |z(t)| - \frac{2\varepsilon^2}{|z(t)|} + |e(t)| \quad \text{if } |y(t)| \geq 3\varepsilon. \quad (4.7)$$

If instead t is so that $|y(t)| < 3\varepsilon$, then since $|y(t) + \xi v(t)| < 4\varepsilon \leq \delta$, it follows that

$$u(t) = y(t) + v(t). \quad (4.8)$$

So

$$w(t) = \frac{v(t)^2 - y(t)^2}{|z(t)| + \sqrt{|z(t)|^2 + v(t)^2 - y(t)^2}}$$

and hence

$$w(t) \leq \frac{v(t)^2 - y(t)^2}{2|z(t)|} \leq \frac{\varepsilon^2}{2|z(t)|}. \quad (4.9)$$

We conclude that, provided $|y(t)| < 3\varepsilon$,

$$|z(t + 1)| \leq |z(t)| + \frac{\varepsilon^2}{2|z(t)|} + |e(t)|. \quad (4.10)$$

In addition,

$$\begin{aligned} |y(t + 1)| &\geq b|x(t)| - |a|(|y(t)| + |u(t)|) - |e_2(t)| \\ &\geq b|x(t)| - (7|a| + 1)\varepsilon, \end{aligned}$$

for $|y(t)| < 3\varepsilon$. If $|x(t)| \geq (1/b)(7|a| + 4)\varepsilon$, then

$$|y(t + 1)| \geq 3\varepsilon, \quad (4.11)$$

and also $|x(t)| \geq 4\varepsilon$ (recall that $b \leq 1$), which implies $|z(t)| \geq 4\varepsilon$. Since $|e(t)| \leq \varepsilon$, from (4.10) it follows that

$$|z(t + 1)| \leq |z(t)| + \frac{\varepsilon^2}{8\varepsilon} + \varepsilon \leq \frac{41}{32}|z(t)|. \quad (4.12)$$

On the other hand, since $|y(t + 1)| \geq 3\varepsilon$, applying (4.7) for $z(t + 2)$, we conclude that

$$|z(t + 2)| \leq |z(t + 1)| - \frac{2\varepsilon^2}{|z(t + 1)|} + |e(t + 1)|. \quad (4.13)$$

Using (4.10) and (4.12) to substitute $|z(t + 1)|$ in the first and second terms of (4.13), we end up with

$$\begin{aligned} |z(t + 2)| &\leq |z(t)| + \frac{\varepsilon^2}{2|z(t)|} - \frac{64\varepsilon^2}{41|z(t)|} + |e(t)| + |e(t + 1)| \\ &< |z(t)| - \frac{\varepsilon^2}{|z(t)|} + |e(t)| + |e(t + 1)|. \end{aligned}$$

Summarizing, we have proved:

Fact I: (i) if $|y(t)| \geq 3\varepsilon$, then

$$|z(t + 1)| \leq |z(t)| - \frac{2\varepsilon^2}{|z(t)|} + |e(t)|; \quad (4.14)$$

(ii) if $|y(t)| < 3\varepsilon$, and $|x(t)| \geq (1/b)(7|a| + 4)\varepsilon$, then

$$|z(t+2)| \leq |z(t)| - \frac{\varepsilon^2}{|z(t)|} + |e(t)| + |e(t+1)|. \quad (4.15)$$

As a consequence of Fact I, we have

Fact II: there exists $t > 0$ such that $z(t)$ is in the region

$$\mathcal{R} = \left\{ x + yi: |x| \leq \frac{1}{b}(7|a| + 4)\varepsilon, |y| \leq 3\varepsilon \right\}.$$

Indeed, if Fact II were not true, then for any $t > 0$ we would have either $|y(t)| \geq 3\varepsilon$ or $|x(t)| \geq (1/b)(7|a| + 4)\varepsilon$. Now we select a sequence (t_0, t_1, t_2, \dots) of integers in the following way:

- $t_0 = 0$,
- for $j \geq 0$, if (4.14) is true for $t = t_j$, then $t_{j+1} = t_j + 1$; otherwise $t_{j+1} = t_j + 2$.

Then we have

$$|z(t_{j+1})| \leq |z(t_j)| - \frac{\varepsilon^2}{|z(t_j)|} + \sum_{k=t_j}^{t_{j+1}-1} |e(k)|. \quad (4.16)$$

Summing (4.16) for $j = 0, 1, 2, \dots, n$, we have

$$|z(t_{n+1})| \leq |z(0)| - \varepsilon^2 \sum_{k=0}^n \frac{1}{|z(t_k)|} + \sum_{k=0}^{t_{n+1}-1} |e(k)|. \quad (4.17)$$

In particular, we have

$$|z(t_{n+1})| \leq |z(0)| + \|e\|_1 = M \quad (4.18)$$

for all $n \geq 0$. So from (4.17) it follows that

$$|z(t_{n+1})| \leq |z(0)| - (n+1)\varepsilon^2/M + \|e\|_1. \quad (4.19)$$

Let $n \rightarrow \infty$. Then $|z(t_{n+1})| \rightarrow -\infty$, which is a contradiction. So Fact II is proved.

To complete the proof of the lemma, it is enough to show the next fact.

Fact III: if $z(T) \in \mathcal{R}$ for some $T \geq 0$, then $|z(t)| \leq r$ for all $t \geq T$.

Note that if $z(t) \in \mathcal{R}$, then

$$|z(t)| \leq \frac{1}{b}(7|a| + 4)\varepsilon + 3\varepsilon. \quad (4.20)$$

If for some T_1 , $z(T_1) \notin \mathcal{R}$, but $z(T_1 - 1) \in \mathcal{R}$, then from Fact II (applied to the trajectory which starts at the state $(x(T_1), y(T_1))$), it follows that there exists $T_2 > T_1$ such that $z(T_2) \in \mathcal{R}$, and $z(t) \notin \mathcal{R}$ for

$T_1 \leq t < T_2$. Now we select $t_0 = T_1, t_1, t_2, \dots, t_n = T_2$ as we did above such that (4.16) is satisfied for $j = 0, 1, 2, \dots, n$. Then

$$|z(t_j)| \leq |z(t_0)| + \sum_{k=t_0}^{t_j-1} |e(k)| \quad (4.21)$$

for $1 \leq j \leq n$. Note that $z(t_0) = e^{i\theta}(z(T_1 - 1) - iu(T_1 - 1)) + e(T_1 - 1)$, and $z(T_1 - 1) \in \mathcal{R}$.

There are two cases to consider now, depending on the sign of $w(T_1 - 1)$. If this quantity is negative, then from (4.5) we know that

$$|z(T_1)| < |z(T_1 - 1)| + \varepsilon.$$

Together with (4.21), we conclude (recall that $t_0 = T_1$) that

$$|z(t_j)| \leq |z(T_1 - 1)| + \varepsilon + \sum_{k=T_1}^{t_j-1} |e(k)|. \quad (4.22)$$

If instead $w(T_1 - 1) > 0$, then from (4.9) it follows that $|y(T_1 - 1)| < \varepsilon$, so we have that also $|u(T_1 - 1)| < 2\varepsilon$. Thus $|z(t_0)| \leq |z(T_1 - 1)| + 2\varepsilon + |e(T_1 - 1)|$. Substituting this into (4.21), we obtain the estimate:

$$|z(t_j)| \leq |z(T_1 - 1)| + 2\varepsilon + \sum_{k=T_1-1}^{t_j-1} |e(k)|, \quad 0 \leq j \leq n. \quad (4.23)$$

For the times of the form t_j , the above bounds will provide the desired conclusions. However, we must take into account as well the cases when $t_j - t_{j-1} = 2$, so that we need to bound the states $x(t_j + 1)$ for such j 's. In that case, from (4.10) and (4.23) we have

$$\begin{aligned} |z(t_j + 1)| &\leq |z(t_j)| + \frac{\varepsilon^2}{2|z(t_j)|} + |e(t_j)| \\ &\leq |z(T_1 - 1)| + \frac{\varepsilon^2}{2|z(t_j)|} + 3\varepsilon. \end{aligned}$$

Since by $z(t_j)$ is not in \mathcal{R} , it follows that $|z(t_j)| \geq \min\{4\varepsilon/b, 3\varepsilon\} > \varepsilon/2$, so we have $|z(t_j + 1)| \leq |z(T_1 - 1)| + 4\varepsilon$. From (4.20) we conclude that

$$|z(t_j + 1)| \leq \frac{1}{b}(7|a| + 4)\varepsilon + 7\varepsilon \quad (4.24)$$

when $t_{j+1} - t_j = 2$. Finally, from inequalities (4.20), together with (4.22) or (4.23) when t is in the sequence of t_j 's, or (4.24) when t is not in this sequence,

imply that $|z(t)| \leq r$ for $T_1 \leq t < T_2$. So Fact III is established. \square

We can summarize the above result, as well as an analogous one-dimensional property, as a general property of certain systems with orthogonal A matrices, as follows.

Corollary 4.2. For $n = 1, 2$, let J be an $n \times n$ matrix, equal to either 1 or -1 if $n = 1$, or of the form

$$\begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix}$$

in the case $n = 2$, with $\alpha^2 + \beta^2 = 1$ and $\beta \neq 0$. Let $b = 1$ if $n = 1$, and $b = (0, 1)'$ if $n = 2$. Then for every $\varepsilon > 0, \delta > 0$ there exists $\theta > 0$ such that for any functions $v: \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}$ and $e: \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}^n$, where $v \leq_{\text{ev}} \theta, e \in l_1$, if $\gamma: \mathbb{Z}_{\geq 0} \rightarrow \mathbb{R}^n$ is any solution of the system

$$\begin{aligned} x(t+1) &= J(x(t) - \sigma(x_n(t) - \xi v(t))b + \eta v(t)b) + e(t), \end{aligned}$$

where $\sigma(s) = \delta \text{sat}(s/\delta), \xi + \eta = 1, \xi, \eta \geq 0$, it follows that

$$\limsup_{t \rightarrow +\infty} |\gamma(t)| < \varepsilon.$$

Proof. Assume first $n = 2$. We pick any $0 < \theta < \min\{(\delta/4), \varepsilon/(7 + (7|a| + 4)/|b|)\}$, and apply Lemma 4.1 with “ ε ” there equal to θ , from which the conclusion follows.

Next, we prove the conclusion for $n = 1$. In this case, the equation of the system becomes

$$x(t+1) = \lambda(x(t) - \sigma(x(t) - \xi v(t)) + \eta v(t)) + e(t),$$

where $\lambda = \pm 1$. Pick any $\theta > 0$. Arguing as earlier (start from a large enough time), we may without loss of generality assume that $\|e\|_1 < \theta$. If $|v(t)| \leq \theta \leq \delta/3$, then for $|x(t)| \geq 3\theta$ we have $|\sigma(x(t) - \xi v(t)) - \eta v(t)| \geq 2\theta$, and $\sigma(x(t) - \xi v(t)) - \eta v(t)$ has the same sign as $x(t)$. So if $|x(t)| \geq 3\theta$, then

$$\begin{aligned} |x(t+1)| &\leq |x(t) - \sigma(x(t) - \xi v(t)) + \eta v(t)| + |e(t)| \\ &\leq |x(t)| - \theta. \end{aligned}$$

Thus there is some t_0 so that $|x(t_0)| \leq 3\theta$. However, the interval $[-3\theta, 3\theta]$ is invariant: it follows from the equation and the fact that $\theta \leq \delta/3$ that $|x(t+1)| \leq 3\theta$

whenever $|x(t)| \leq 3\theta$. So $\limsup_{t \rightarrow +\infty} |x(t)| \leq 3\theta$. Now, to obtain the conclusion of the corollary, it suffices to take $\theta = \min\{\delta/3, \varepsilon/3\}$. \square

5. Proof of Theorem 1

First, we notice that under the conditions of the theorem there exists a linear change of coordinates of the state space that transforms Σ into the block form

$$\Sigma: \begin{cases} x_1(t+1) = A_1 x_1(t) + B_1 u(t), & x_1(t) \in \mathbb{R}^{n_1}, \\ x_2(t+1) = A_2 x_2(t) + B_2 u(t), & x_2(t) \in \mathbb{R}^{n_2}, \end{cases}$$

where (i) $n_1 + n_2 = n$, (ii) all the eigenvalues of A_1 have magnitude 1, (iii) all the eigenvalues of A_2 have magnitude less than 1, and (iv) (A_1, B_1) is a controllable pair. Suppose that we find an IICS-stabilizing feedback $u = k(x_1)$ of Type \mathcal{F} or Type \mathcal{G} for the system $x_1(t+1) = A_1 x_1(t) + B_1 u(t)$ such that the resulting closed-loop system is asymptotically stable. Then this same feedback law will stabilize Σ as well, because the second equation, $x_2(t+1) = A_2 x_2(t) + B_2 k(x_1(t))$, can be seen as an asymptotically stable linear system forced by a function that converges to zero. Thus, in order to stabilize Σ , it is enough to stabilize the “critical subsystem” $x_1(t+1) = A_1 x_1(t) + B_1 u(t)$. Without loss of generality, in our proof of the theorem we will suppose that Σ is already in this form, that is, we assume that all the eigenvalues of A have magnitude 1 and that the pair (A, B) is controllable.

5.1. Single-input case

We start with the single-input case, and prove the theorem by induction on the dimension n of the system.

For dimension zero there is nothing to prove. Now assume that we are given a single-input n -dimensional system, $n \geq 1$, and suppose that Theorem 1 has been established for all single-input systems of dimension less than or equal to $n - 1$. We consider separately the following two possibilities:

- (i) 1 or -1 is an eigenvalue of A ,
- (ii) neither 1 nor -1 is an eigenvalue of A .

Write $N = N(A)$, and pick any $\varepsilon > 0$. We want to prove the existence of IICS-stabilizing feedbacks $u = -k_{\mathcal{F}}(x)$ and $u = -k_{\mathcal{G}}(x)$, where $k_{\mathcal{F}} \in \mathcal{F}_n(\sigma), k_{\mathcal{G}} \in \mathcal{G}_n(\sigma)$, for some finite sequence $\sigma = (\sigma_1, \dots, \sigma_N)$ of functions in \mathcal{L} , with $\|\sigma\| \leq \varepsilon$. (The negative signs are merely for notational convenience; since

saturations are odd functions, the signs can be switched by changing coefficients of linear combinations.)

In Case (i), we apply Part (i) of Lemma 3.1 and rewrite our system in the form

$$\begin{aligned}\bar{y}(t+1) &= A_1 \bar{y}(t) + (y_n(t) + u(t))b_1, \\ y_n(t+1) &= \lambda(y_n(t) + u(t)),\end{aligned}\quad (5.1)$$

where $\bar{y} = (y_1, \dots, y_{n-1})'$. (Note that if $n = 1$, only the second equation appears.) In Case (ii), since $n > 0$, A has a pair of eigenvalues of the form $\alpha + \beta i$, with $\beta \neq 0$. So we apply Part (ii) of Lemma 3.1 and make a linear transformation that puts Σ in the form

$$\begin{aligned}\bar{y}(t+1) &= A_1 \bar{y}(t) + (y_n(t) + u(t))b_1, \\ y_{n-1}(t+1) &= \alpha y_{n-1}(t) - \beta(y_n(t) + u(t)), \\ y_n(t+1) &= \beta y_{n-1}(t) + \alpha(y_n(t) + u(t)),\end{aligned}\quad (5.2)$$

where $\bar{y} = (y_1, y_2, \dots, y_{n-2})'$. (In the special case when $n = 2$, the first equation will be missing.) So, in either case, we can rewrite our system in the form

$$\begin{aligned}\bar{y}(t+1) &= A_1 \bar{y}(t) + (y_n(t) + u(t))b_1, \\ \bar{y}(t+1) &= J(\bar{y}(t) + u(t)b_0),\end{aligned}\quad (5.3)$$

where J is as in Corollary 4.2 and b_0 is like b in that corollary. To consider the problem of IICS-stabilizing feedback, we must study solutions of the following system:

$$\begin{aligned}\bar{y}(t+1) &= A_1 \bar{y}(t) + (y_n(t) + u(t))b_1 + \bar{e}(t), \\ \bar{y}(t+1) &= J(\bar{y}(t) + u(t)b_0) + \bar{e}(t),\end{aligned}\quad (5.4)$$

where \bar{e}, \tilde{e} are arbitrary elements of l_1 .

We will design a feedback of the form

$$u = \sigma_N(-y_n + \xi v) + \eta v = -\sigma_N(y_n - \xi v) + \eta v, \quad (5.5)$$

where ξ and η are constants such that $\xi\eta = 0$, $\xi + \eta = 1$, $\sigma_N(s) = \varepsilon \text{sat}(s/\varepsilon)$, and v is to be chosen later.

From Corollary 4.2 we may pick a $0 < \theta < \varepsilon/2$ such that, if $|v(t)| \leq_{\text{ev}} \theta$, then all trajectories of (5.4) satisfy $|\bar{y}| \leq_{\text{ev}} \varepsilon/2$. Consider one such trajectory. Then, for all t sufficiently large, $u(t) = -y_n(t) + v(t)$, and the first block equation in (5.4) becomes

$$\bar{y}(t+1) = A_1 \bar{y}(t) + v(t)b_1 + \bar{e}(t) \quad (5.6)$$

for all large t . Note that (A_1, b_1) is controllable and all eigenvalues of A_1 have magnitude 1. By the inductive

hypothesis, we conclude that there exist

$$\bar{k}_{\mathcal{F}} \in \mathcal{F}_n(\bar{\sigma}) \quad \text{and} \quad \bar{k}_{\mathcal{G}} \in \mathcal{G}_n(\bar{\sigma}) \quad (5.7)$$

for some $\bar{\sigma} = (\sigma_1, \dots, \sigma_{N-1})$ such that $\|\bar{\sigma}\| \leq \theta$, each of which is IICS-stabilizing for the system $\bar{y}(t+1) = A_1 \bar{y}(t) + u(t)b_1$.

We let

$$k_{\mathcal{F}}(y) = \sigma_N(-y_n + \bar{k}_{\mathcal{F}}(\bar{y}))$$

and

$$k_{\mathcal{G}}(y) = \sigma_N(-y_n) + \bar{k}_{\mathcal{G}}(\bar{y})$$

(cases $\xi = 1, \eta = 0$, and $\xi = 0, \eta = 1$, respectively), and claim that these are IICS-stabilizing for the original system. Locally around the origin, the closed-loop system is linear, so stability is not an issue, and it is enough to prove the attraction property. We must show that, for any \bar{e}, \tilde{e} elements of l_1 , all solutions converge to zero. Pick any such trajectory. As discussed, u is eventually linear in the variables y_n and v , where we are taking $v = \bar{k}_{\mathcal{F}}(\bar{y})$ or $v = \bar{k}_{\mathcal{G}}(\bar{y})$. By the inductive construction, we know that also $\bar{y}(t) \rightarrow 0$ as $t \rightarrow \infty$, which means that, since v is a linear function of \bar{y} when \bar{y} is small, (5.4) will eventually become a linear asymptotically stable system with a converging input, and thus the state indeed converges to zero. The sequence $\sigma = (\sigma_1, \dots, \sigma_{N-1}, \sigma_N)$ clearly satisfies $\|\sigma\| \leq \varepsilon$. The proof for the single-input case is completed.

5.2. The general case

Next, we deal with the general case of $m > 1$ inputs and prove Theorem 1 by induction on m .

First, we know from the proof above that the theorem is true if $m = 1$. Assume that Theorem 1 has been established for all k -input systems, for all $k \leq m - 1$, and let $\Sigma: x(t+1) = Ax(t) + Bu(t)$ be an m -input system.

Assume without loss of generality that the first column b_1 of B is nonzero and consider the Kalman controllability decomposition of the system $\Sigma_1: x(t+1) = Ax(t) + b_1 u_1(t)$ (see e.g. [5, Lemma 3.3.3]). We conclude that, under a change of coordinates $y = T^{-1}x$, Σ_1 has the form

$$y_1(t+1) = A_1 y_1(t) + A_2 y_2(t) + \bar{b}_1 u_1(t),$$

$$y_2(t+1) = A_3 y_2(t),$$

where (A_1, \bar{b}_1) is a controllable pair. In these coordinates Σ has the form

$$\begin{aligned} y_1(t+1) &= A_1 y_1(t) + A_2 y_2(t) + \bar{b}_1 u_1(t) + \bar{B}_1 \bar{u}(t), \\ y_2(t+1) &= A_3 y_2(t) + \bar{B}_2 \bar{u}(t), \end{aligned} \quad (5.8)$$

where $\bar{u} = (u_2, \dots, u_m)'$ and \bar{B}_1, \bar{B}_2 are appropriate matrices. So it suffices to show the conclusion for (5.8). Let n_1, n_2 denote the dimensions of y_1, y_2 , respectively. Recall that $N = N(A)$. For the single-input controllable system

$$y_1(t+1) = A_1 y_1(t) + \bar{b}_1 u_1(t),$$

there is a feedback

$$u_1 = k_1(y_1) \quad (5.9)$$

such that (i) $k_1 \in \mathcal{F}_{n_1}(\sigma_1, \dots, \sigma_{N_1})$ (respectively, $k_1 \in \mathcal{G}_{n_1}(\sigma_1, \dots, \sigma_{N_1})$) where $N_1 = N(A_1)$; (ii) the resulting closed-loop system is IICS; (iii) $\|\sigma_1\| \leq \varepsilon$, where $\sigma_1 = (\sigma_1, \dots, \sigma_{N_1})$. Since (5.8) is controllable, we conclude that the $(m-1)$ -input subsystem $y_2(t+1) = A_3 y_2(t) + \bar{B}_2 \bar{u}(t)$ is controllable as well. By the inductive hypothesis, this subsystem can be stabilized by a feedback

$$\bar{u} = \bar{k}(y_2) = (k_2(y_2), \dots, k_m(y_2)) \quad (5.10)$$

such that (i) $\bar{k} \in \mathcal{F}_{n_2}^{\bar{l}}(\sigma_{N_1+1}, \dots, \sigma_N)$ (respectively, $\bar{k} \in \mathcal{G}_{n_2}^{\bar{l}}$), where $\bar{l} = (N_2, \dots, N_m)$ is an $(m-1)$ -tuple of nonnegative integers and $|\bar{l}| = N - N_1$; (ii) the resulting closed-loop system is IICS; (iii) $\|\sigma_2\| \leq \varepsilon$, where $\sigma_2 = (\sigma_{N_1+1}, \dots, \sigma_N)$. We let $k(y) = (k_1(y_1), \bar{k}(y_2))$. This globally stabilizes (5.8), and the resulting closed-loop system is IICS. Indeed, around the origin the system (5.8) has a block triangular linear form, whose diagonal blocks are asymptotically stable, so stability is automatic. Consider now any $e_1, e_2 \in l_1$ and any solution of (5.8) with e_1, e_2 added to the respective blocks. Then $y_2(t) \rightarrow 0$ as $t \rightarrow \infty$ because \bar{k} is IICS-stabilizing. Moreover, since near the origin

the system is linear, y_2 is an l_1 function itself. Now consider the first block of equations, viewing

$$A_2 y_2(t) + \bar{B}_1 \bar{k}(y_2(t)) + e_1(t)$$

as an l_1 perturbation. Since k_1 is IICS-stabilizing, it follows that $y_1(t) \rightarrow 0$ as $t \rightarrow \infty$ as well. So if we let $l = (N_1, N_2, \dots, N_m)$ and $k = (k_1(y_1), k_2(y_2), \dots, k_m(y_2))$, then $k \in \mathcal{F}_n^l(\sigma)$ (respectively, $k \in \mathcal{G}_n^l(\sigma)$), $\sigma = (\sigma_1, \dots, \sigma_N)$, satisfies all the required properties as desired. \square

References

- [1] Z. Lin, A. Saberi, Semi-global exponential stabilization of linear discrete-time systems subject to "input saturation" via linear feedbacks, *Systems Control Lett.* 24 (1995) 125–132.
- [2] W. Liu, Y. Chitour, E.D. Sontag, On finite gain stabilizability of linear systems subject to input saturation, *SIAM J. Control Optim.* 34 (1996), to appear.
- [3] E.D. Sontag, An algebraic approach to bounded controllability of linear systems, *Internat. J. Control* 39 (1984) 181–188.
- [4] E.D. Sontag, Remarks on stabilization and input-to-state stability, *Proc. IEEE Conf. Decision and Control*, Tampa, 1989, IEEE Publications, New York, 1989, pp. 1376–1378.
- [5] E.D. Sontag, *Mathematical Control Theory: Deterministic Finite Dimensional Systems*, Springer, New York, 1990.
- [6] E.D. Sontag, H.J. Sussmann, Nonlinear output feedback design for linear systems with saturating controls, *Proc. IEEE Conf. Decision and Control*, Honolulu, 1990, IEEE Publications, New York, pp. 3414–3416.
- [7] E.D. Sontag, Y. Yang, Global stabilization of linear systems with bounded feedback, *Technical Report SYCON-91-09*, Rutgers Center for Systems and Control, 1991.
- [8] H.J. Sussmann, E.D. Sontag, Y. Yang, A general result on the stabilization of linear systems using bounded controls, *IEEE Trans. Automat. Control* 39 (1994) 2411–2425.
- [9] A.R. Teel, Global stabilization and restricted tracking for multiple integrators with bounded controls, *Systems Control Lett.* 18 (1992) 165–171.
- [10] A.G. Tsirukis, M. Morari, Controller design with actuator constraints, *Proc. IEEE Conf. Decision and Control*, Tucson, 1992, IEEE Publications, New York, pp. 2623–2628.
- [11] Y. Yang, Global stabilization of linear systems with bounded feedback, Ph.D. Thesis, Mathematics Department, Rutgers University, 1993.