



## Stability, controllability, and observability criteria for state-space dynamical systems on measure chains with an application to fixed point arithmetic



Yan Wu<sup>a</sup>, Sailaja P<sup>b</sup>, K. N. Murty<sup>c,\*</sup>

<sup>a</sup>Department of Mathematical Sciences, Georgia Southern University, Statesboro, GA 30460, USA.

<sup>b</sup>Department of Mathematics, Geethanjali Engineering College, Hyderabad, Telangana 501301, India.

<sup>c</sup>Department of Applied Mathematics, Andhra University, Waltair, AP 530017, India.

### Abstract

In this paper, our main attempt is to unify results on stability, controllability, and observability criteria on real-time dynamical systems with non-uniform domains. The results of continuous/discrete systems will now become a particular case of our results. As an application a first-order time scale dynamical system on measure chains in one-dimensional state space having both continuous/discrete filters to minimize the effect of a round of noise at the filter outputs is presented. A set of necessary and sufficient conditions for this dynamical system to be stable and completely stable are established.

**Keywords:** Linear Systems, time scale dynamical systems, control systems, concurrency control.

**2010 MSC:** 93B05, 93B07, 93B20, 93B55, 93D99.

©2020 All rights reserved.

### 1. Introduction

Our aim, in this paper, is to develop the foundation for a comprehensive linear system theory, which in-fact coincides with the existing canonical system theories in the continuous as well as digital systems, but also extend those theories to dynamical systems with non-uniform domains. A fascinating fact is that all the widely different disciplines of an application depend on a common core of Mathematical techniques of the modern control systems theory [1, 2]. Our main object in this paper is to unify results on stability, controllability, and observability criteria on real-time dynamical systems and deduce the result on continuous/discrete systems as a particular case. The paper is organized as follows. Section 2 presents some salient features of time scale dynamic systems that are needed for our later discussion. This section deals with time varying system, and then presents the time invariant system on a time scale dynamical system. We present a set of necessary and sufficient conditions in a more restrictive time invariant setting relative to the time varying system on both controllability and observability. Further,

\*Corresponding author

Email address: [nkanuri@hotmail.com](mailto:nkanuri@hotmail.com) (K. N. Murty)

doi: [10.22436/jnsa.013.04.03](https://doi.org/10.22436/jnsa.013.04.03)

Received: 2019-10-20 Revised: 2019-11-27 Accepted: 2019-12-10

results on realizable criteria are also presented on regressive linear system. We introduce the concepts of stability, controllability, and observability criteria on time scale dynamical system. Section 3 is concerned with round off noise minimization for 1-D state space digital filters using joint optimization of error feedback. Results presented in this section generalize the results of [7] and includes them as a particular case for discrete systems. Further, our theory unifies both continuous and discrete systems on noise minimization for 2-D state space digital/continuous filters. We present a set of necessary and sufficient conditions for the first order time scale dynamical system to be stable, completely stable, completely controllable and completely observable. Further, we present more convenient criteria for controllability and observability under smoothness conditions on time scales. Section 4, presents stability analysis of real time dynamical systems on measure chains in one-dimensional state space digital filters/continuous filters to minimize the effects of round of noise at the filter outputs subjected to suitable norms. Norms that suit for controllability, observability and reliability are the norms that are discussed in [8]. In [5], the authors established useful results on  $\Psi$ -bounded solutions of linear first order differential systems on time scales. These results are used as a tool to develop our results in this paper.

## 2. Basic results on time scale dynamical systems

In this section, we outline some of the basic notions on time scales. A time scale  $T$  is a closed subset of  $\mathbb{R}$ , and examples of time scales include  $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{R}$ , Cantor's set etc. The set  $Q$  and  $A = \{t \in \mathbb{R} : 0 \leq t \leq 1\}$  are time scales. Time scales are not necessarily connected. In order to overcome these deficiencies, we introduce the notion of jump operators. The mappings  $\sigma, \rho : T \rightarrow T$  defined by  $\sigma(t) = \inf\{s \in T : s > t\}$ ,  $\rho(t) = \sup\{s \in T : s < t\}$  are called jump operators. A point  $t \in T$  is said to be right dense, right scattered, left dense, and left scattered accordingly as  $\sigma(t) = t$ ;  $\sigma(t) > t$ ,  $\rho(t) = t$ , and  $\rho(t) < t$ , respectively. The graininess  $\mu : T \rightarrow \{0, \infty\}$  is defined by  $\mu(t) = \sigma(t) - t$ .

We say that  $f$  is rd-continuous, if it is continuous at right-dense points and if  $\lim_{s \rightarrow t} f(s)$  as  $s \rightarrow t$  exists for all left dense points  $t \in T$ .

A function  $f : T \rightarrow T$  is said to be differentiable at  $t \in T^k = \{T - \rho(t), \max(T)\}$ , if

$$\lim_{s \rightarrow t} \frac{f(\sigma(t)) - f(s)}{\sigma(t) - s},$$

where  $s \in t - \{\sigma(t)\}$  exists, and is said to be differentiable on  $T$ , provided it is differentiable at each  $t \in T^k$ .

A function  $f : T \rightarrow T$  with  $F(t) = f(t)$  for all  $t \in T^k$  is said to be an anti derivative of  $f$  on  $T$ , and in this case

$$\int_s^t f(\tau) \Delta \tau = F(t) - F(s)$$

for all  $s, t \in T$ . Let  $f : T \rightarrow T$ . We note that, if  $T = \mathbb{R}$  and  $a, b \in T$ , then  $f^\Delta(t) = f'(t)$  and

$$\int_s^b f(t) \Delta t = \int_s^b f(t) dt.$$

On the other hand, if  $T = \mathbb{Z}$ , then  $f^\Delta(t) = \Delta f(t) = f(t+1) - f(t)$  and

$$\int_s^b f(t) \Delta t = \begin{cases} \sum_{k=a}^{b-1} f(k), & \text{if } a < b, \\ 0, & \text{if } a = b, \\ \sum_{k=b}^{a-1} f(k), & \text{if } a > b. \end{cases}$$

If  $f, g : T \rightarrow \mathbb{R}$  and  $t \in T^k$ , then

- 1)  $(f + g)^\Delta(t) = f^\Delta(t) + g^\Delta(t)$ ;
- 2)  $(fg)^\Delta(t) = f^\Delta(t)g(t) + f(\sigma(t))g^\Delta(t)$ ;
- 3)  $f(\sigma(t)) = f(t) + \mu(t)f^\Delta(t)$ ;
- 4)  $(kf)^\Delta(t) = kf^\Delta(t)$ ; (for any scalar  $k$ ).

Note that, if  $f^\Delta(t)$  exists, then  $f$  is continuous at  $t$ , if  $t$  is right-scattered and  $f$  is continuous at  $t$ , then  $f^\Delta(t) = \frac{f(\sigma(t)) - f(t)}{\mu(t)}$ .

Note that  $\mu$  is called the graininess. For those not familiar with the rapidly expanding area of dynamic equations on time scales, excellent survey can be found in [3, 4, 6]. We examine the fundamental notion on controllability, observability and stability commonly dealt with in linear system control theory on time scale dynamical systems. Our theory in fact generalizes the results on continuous and discrete linear systems in a general frame work. Note that our approach to dynamical systems generalizes both continuous/digital filters and include the earlier results as a particular case. More specifically, our focus here is how to generalize these concepts to the non-unified domain setting while at the same time preserving and unifying the well-known bodies of knowledge on these subjects in the continuous as well as discrete cases. This generalized framework has already shown promising application to adaptive control systems [6].

### Controllability, observability, and realizability criteria

In linear systems theory, we say that a system is controllable, if the solution of the relevant dynamical system (continuous/discrete/hybrid) can be obtained to a specified final state in time. For, we have

**Definition 2.1.** Let  $GA(t) \in \mathbb{R}^{n \times n}$ ,  $B(t) \in \mathbb{R}^{n \times n}$ ,  $C(t) \in \mathbb{R}^{n \times n}$  and  $D(t) \in \mathbb{R}^{p \times m}$  all be rd-continuous matrix functions defined on the time scale  $T$  with  $p, m \leq n$ . The system

$$x(t) = A(t)x(t) + B(t)U(t), x(t_0) = x_0, \quad (2.1)$$

$$y(t) = C(t)x(t) + D(t)U(t), \quad (2.2)$$

is controllable on  $[t_0, t_f]$  if, given any initial state  $x(t_0) = x_0$ , there exists a rd-continuous input  $U(t)$  such that the corresponding solution of the system satisfies  $x(t_f) = x_f$ . The time varying system is completely controllable, if it is controllable for all  $t \in [t_0, t_f]$ . When  $T = \mathbb{R}$ , (2.1) is equivalent to

$$x'(t) = A(t)x(t) + B(t)U(t), x(t_0) = x_0$$

and when  $T = \mathbb{Z}$  (discrete), (2.1) is equivalent to

$$x(j+1) = A(j)x(j) + B(j)U(j), x(j_0) = x_0,$$

Our first result establishes a necessary and sufficient condition for the controllability of the linear dynamical system (2.1)-(2.2). Note that for a time varying linear system, the connection of the input signal to the state variables can change with time.

**Theorem 2.2.** Any solution  $x(t)$  of (2.1) satisfying  $x(t_0) = x_0$  is given by

$$x(t) = \Phi(t, t_0)x_0 + \int_{t_0}^t \Phi(t, \sigma(s))B(s)U(s)\Delta s.$$

*Proof.* Any solution of  $x^\Delta(t) = A(t)x(t)$  has the form  $x(t) = \Phi(t)C$ . At  $t = t_0$ , we have  $x(t_0) = \Phi(t_0)C$  or  $C = \Phi^{-1}(t_0)x(t_0) = \Phi^{-1}(t_0)x_0$ . Therefore,  $x(t) = \Phi(t)\Phi^{-1}(t_0)x(t_0) = \Phi(t, t_0)x_0$ .

Further, if  $x(t)$  is any solution of (2.1) and  $\bar{x}(t)$  is a particular solution of (2.1), then  $x(t) - \bar{x}(t)$  is a solution of the dynamic system  $\dot{x}(t) = A(t)x(t)$ . Thus,

$$x(t) = \Phi(t, t_0) x_0 + \int_{t_0}^t \Phi(t, \sigma(s)) B(s) U(s) ds.$$

Note that the second term in the above equation is the particular solution of the Non-homogeneous dynamical system (2.1).  $\square$

**Theorem 2.3.** *The time scale dynamical system (2.1) (The regressive linear system)*

$$\begin{aligned} x^\Delta(t) &= A(t)x(t) + B(t)U(t), \quad x(t_0) = x_0, \\ y(t) &= C(t)x(t) + D(t)U(t) \end{aligned}$$

is complete controllable on  $[t_0, t_f]$  if and only if the  $(n \times n)$  controllability Gramian matrix given by

$$\omega(t_0, t_f) = \int_{t_0}^{t_f} \Phi(t_0, \sigma(s)) B(s) B^*(s) \Phi^*(t_0, \sigma(s)) \Delta s,$$

is non-singular where  $\Phi(t, t_0)$  is the transition matrix for the system  $x^\Delta(t) = A(t)x(t)$ ,  $x(t_0) = I$ .

*Proof.* Suppose  $\omega(t_0, t_f)$  is invertible. Then given  $x_0$  and  $x_f$ , one can choose the input signal  $U(t)$  as  $U(t) = -B^*(t)\Phi^*(t_0, \sigma(t))\omega^{-1}(t_0, t_f)(x_0 - \Phi(t_0, t_f)x_f)$ . Clearly, the input signal  $U(t)$  is continuous on the interval  $[t_0, t_f]$ , and the corresponding solution of the system at  $t = t_f$  is given by

$$\begin{aligned} x(t_f) &= \Phi(t_f, t_0) x_0 + \int_{t_0}^{t_f} \Phi(t_f, \sigma(s)) B(s) U(s) \Delta s \\ &= \Phi(t_f, t_0) x_0 - \int_{t_0}^{t_f} \Phi(t_f, \sigma(s)) B(s) B^*(s) \Phi^*(s) \Phi^*(t_0, t_f) \omega^{-1}(t_0, t_f) (x_0 - \Phi(t_0, t_f) x_f) \Delta s \\ &= \Phi(t_f, t_0) x_0 \\ &\quad - \Phi(t_f, t_0) \int_{t_0}^{t_f} \Phi(t_f, \sigma(s)) B(s) B^*(s) \Phi^*(s) \Phi^*(t_0, \sigma(s)) \Delta s \omega^{-1}(t_0, t_0) (x_0 - \Phi(t_0, t_f) x_f) \\ &= \Phi(t_f, t_0) x_0 - (\Phi(t_f, t_0) x_0 - x_f) \\ &= x_f. \end{aligned}$$

So, the state equation is controllable on  $[t_0, t_f]$ . This is true for all  $t \in [t_0, t_f]$ , it follows that the state equation is completely controllable. Now, to show the reverse implication, suppose that the state equation is completely controllable on  $[t_0, t_f]$  and for the sake of contradiction, assume that the matrix is  $\omega(t_0, t_f)$  is singular. Since  $\omega[t_0, t_f]$  is non-invertible, there exists an  $(n \times 1)$  vector  $\alpha$  such that

$$\alpha^* \omega(t_0, t_f) \alpha = \int_{t_0}^{t_f} \alpha^* \Phi(t_0, \alpha(s) B(s)) B^*(s) \Phi^*(t_0, t_f) \Delta s = 0. \quad (2.3)$$

Thus

$$\int_{t_0}^{t_f} \|\alpha^* \Phi(t_0, \alpha(s) B(s))\|^2 \Delta s = 0. \quad (2.4)$$

This implies  $\alpha^* \Phi(t_0, \alpha(s) B(s)) = 0, t \in [t_0, t_f]$ . However, the state equation is controllable on  $[t_0, t_f]$ , we can chose  $x_0 = \alpha + \Phi(t_0, t_f)x_f$  and an input signal  $U(t)$  such that

$$x_f = \Phi(t_f, t_0) x_0 + \int_{t_0}^{t_f} \Phi(t_f, \sigma(s)) B(s) U_\alpha(s) \Delta s,$$

which is equivalent to the equation

$$U_\alpha = - \int_{t_0}^{t_f} \Phi(t_f, \sigma(s)) B(s) U_\alpha(s) \Delta s.$$

Multiplying both sides with  $U_\alpha^*$  and using (2.3) and (2.4) yields  $U_\alpha^* U_\alpha = 0$ , a contradiction. Thus, the matrix  $\omega(t_0, t_f)$  is non-singular.

Since the controllability Gramian is symmetric and positive definite, the above theorem can be interpreted as saying that (2.1) is controllable on  $[t_0, t_f]$  if and only if the Gramian is positive definite.  $\square$

Next, we turn our attention to the observability criteria of a linear time scale dynamical system. In the linear systems theory when the term observability is used, it refers to the effect that the state equation has the output of the state equation. As such the concept is unchanged by assuming that the response of the system to zero input. As such we define the following.

**Definition 2.4.** The regressive linear system

$$x^\Delta(t) = A(t)x(t), \quad x(t_0) = x_0, \quad y(t) = C(t)x(t), \quad (2.5)$$

is observable on  $[t_0, t_f]$  if any initial state  $x(t_0) = x_0$  is uniquely determined by the corresponding response  $y(t)$  for  $t \in [t_0, t_f]$ .

**Theorem 2.5.** The regressive linear system (2.5) is completely observable on  $[t_0, t_f]$  if and only if the  $(n \times n)$  symmetric observability matrix

$$M(t_0, t_f) = \int_{t_0}^{t_f} \Phi(s, t_0) C^*(s) C(s) \Phi(s, t_0) \Delta s,$$

is non-singular.

*Proof.* First suppose that the system (2.5) is completely observable and suppose that  $M(t_0, t_f)$  is non-invertible. Then there exists a non-zero vector  $x_0$  such that  $M(t_0, t_f)x_0 = 0$ . Then, clearly  $x^* M(t_0, t_f)x_0 = 0$ .

Hence  $C(t)\Phi(t, t_0)x_0 = 0, t \in [t_0, t_f]$ . Thus,  $x(t_0) = x_0 + x_a$  yields the same zero-input response for the system with  $x(t_0) = x_0$ , and the system is not observable on  $[t_0, t_f]$  a contradiction.

Conversely, suppose the Gramian matrix  $M(t_0, t_f)$  is invertible. Then the solution expression  $y(t) = C(t)\Phi(t, t_0)x_0$  is multiplied by  $\Phi^*(t, t_0) x_0 C^*(t)$  and integrating, we get

$$\int_{t_0}^{t_f} \Phi^*(t, t_0) C^*(t) y(t) \Delta t = M(t_0, t_f) x_0.$$

The left hand side of the above expression is determined by  $y(t)$  for  $t \in [t_0, t_f]$  and the equation is a linear algebraic equation in  $x_0$ . Since  $M(t_0, t_f)$  is non-singular, it follows that  $x_0$  is determined uniquely and hence the state equation is observable. This is true for all  $t \in [t_0, t_f]$ , it follows that the state equation is completely observable.  $\square$

It may be noted that the observability Gramian, like the controllability Gramian is positive semi definite symmetric matrix. It is positive definite, if the state equation is observable and in fact completely observable. It may be noted that the Gramian condition is not very practical as it requires explicit knowledge of the transition matrix. Thus, we present a sufficient condition that is much easier to check the criteria of observability.

**Definition 2.6.** If  $T$  is a time scale such that  $\mu$  is sufficiently differentiable with the indicated rd-continuous derivatives, we define a sequence of  $(p \times n)$  matrix functions as

$$L_0(t) = C(t), \quad L_j(t) = L_{j-1}(t)A(t) + L_{j-1}^\Delta(t)(I + \mu(t)A(t)), j = 1, 2, 3, \dots$$

It can easily be verified by induction argument that

$$L_j(t) = \frac{\partial_j}{\Delta t_j} [C(t) \Phi(t, s)]_{s=t}.$$

**Theorem 2.7.** Suppose  $m$  is a positive integer such that for  $t \in [t_0, t_f]$ ,  $C(t)$  is  $m$  times continuously  $\Delta$  differentiable and both  $U(t)$  and  $A(t)$  are  $(m-1)$  times continuously  $\Delta$  differentiable. Then the dynamic equation  $x^\Delta(t) = A(t)x(t)$ ,  $x(t_0) = x_0$  is completely observable on  $[t_0, t_f]$ , if for some  $t_c \in [t_0, t_f]$ ,  $\text{rank}[L_0[t_c], L_1[t_c], \dots, L_m[t_c]]^T = n$ , where

$$L_j(t) = \frac{\partial^j}{\Delta s^j} [C(t)\Phi(t, s)]_{s=t}, j = 0, 1, 2, \dots$$

We now proceed to introduce the concept of reliability criteria on the time scale dynamical system. In linear system theory, the concept of reliability refers to the ability to characterize a known output in terms of a linear system with some input. Our interest here is in the reversal of the computation, and in particular we will establish conditions on some specific  $G(t, \sigma(s))$ . More specifically, we assume zero initial state and the output signal  $y(t)$  corresponding to a given input signal  $U(t)$  is given by

$$y(t) = \int_{t_0}^t G(t, \sigma(s))U(s)\Delta s + D(t)U(t), t > t_0,$$

where

$$G(t, \sigma(s)) = C(t) \Phi(t, \sigma(s)) B(s).$$

**Definition 2.8.** The regressive linear system

$$x^\Delta(t) = A(t)x(t) + B(t)U(t), \quad x(t_0) = 0, \quad y(t) = C(t)x(t)$$

of dimension  $n$  is a realization of the weighting pattern  $G(t, \sigma(s))$ , if  $G(t, \sigma(s)) = C(t) \Phi(t, \sigma(s)) B(s)$  for all  $t, s$ . If a realization of this system exists, then the weighting pattern is realizable. The system is a minimal realization if no realization of  $G(t, \sigma(s))$  with dimension less than  $n$ .

It may be noted that for the system

$$x^\Delta(t) = A(t)x(t) + B(t)U(t), \quad x(t_0) = 0, \quad y(t) = C(t)x(t) + D(t)U(t),$$

the output signal  $y(t)$  corresponding to a given input  $U(t)$  and weighting pattern  $G(t, \sigma(s)) = C(t)\Phi(t, \sigma(s))B(s)$  is given by

$$y(t) = \int_{t_0}^t G(t, \sigma(s))U(s)\Delta s + D(t)U(t), t \geq t_0,$$

then there exists a realization of a particular weighting pattern  $G(t, \sigma(s))$ , there will be in fact many such weighting patterns since a change of state variables will leave the weighting pattern unchanged. Also, there can be many different realizations of the same weighting pattern that all have different dimensions. This is why we are careful to distinguish between realizations and minimal realizations in our discussion. We complete our discussion by considering stability analysis of digital/continuous filter implementation subject to Finite word length (FWL) effects based on eigenvalue sensitivity analysis.

### 3. Stability analysis

Suppose that a local state-space (LSS) model for a 2-D recursive time scale filter is described by

$$x^\Delta(t) = A x(t) + B u(t), \quad y(t) = C x(t) + D u(t),$$

where  $A, B, C$ , and  $D$  are all constant matrices with appropriate dimensions. Note that, if  $A$  is a constant matrix, then its fundamental matrix  $\Phi(t, t_0) = e^{(t-t_0)A}$ . We assume that the initial point  $t_0 = 0$ . The regressive system

$$x^\Delta(t) = A x(t), \quad x(t_0) = 0$$

is said to be stable, if there exists a positive constraint  $M$  such that  $\|\Phi(t, t_0)\| \leq M$  implies  $\|\Phi(t)\| < \varepsilon, \forall t \geq t_0$ .

**Definition 3.1.** The regressive system

$$x^\Delta(t) = A x(t), \quad x(t_0) = x_0,$$

is said to be uniformly stable if there exist positive constants  $\lambda, \mu > 0$  such that for any  $t_0$  and  $x(t_0)$ , the corresponding solution satisfies

$$\|x^\Delta(t)\| \leq \|x(t_0)\| \mu e^{-\lambda(t-t_0)}, \quad t \geq t_0.$$

**Theorem 3.2.** The time invariant system  $x^\Delta(t) = A x(t), x(t_0) = x_0$  is uniformly exponential stable if and only if for some  $\beta > 0$ ,

$$\lim_{t \rightarrow \infty} \int_{t_0}^t \|e^t\| \Delta t < \lambda e^{-\lambda t} \Delta t = \frac{\lambda}{\mu} = \beta. \quad (3.1)$$

*Proof.* Suppose that the system is uniformly exponential stable. It is claimed that there exists a  $\beta > 0$  such that

$$\lim_{t \rightarrow \infty} \int_{t_0}^t \|e^t\| \Delta t < \beta.$$

We have

$$\int_0^\infty \|e^t\| \Delta t < \int_0^\infty \lambda e^{-\lambda t} \Delta t = \frac{\lambda}{\mu} = \beta.$$

Hence the claim holds. Now, to prove the reverse inequality, assume that condition (3.1) is satisfied. To the contrary, suppose the system is not uniformly exponential stable. Then find  $\lambda, \mu > 0$  with  $\lambda \in \mathbb{R}^T$ , we have  $\|e^t\| > \mu e^{-\lambda t}$ . Hence

$$\int_0^\infty \|e^t\| \Delta t < \int_{t_0}^\infty \mu e^{-\lambda t} \Delta t = \frac{\lambda}{\mu}.$$

In partial choose  $\lambda > \beta \mu$ , then

$$\int_0^\infty \|e^t\| \Delta t > \frac{\beta \mu}{\mu} = \beta,$$

a contradiction. □



#### 4. Stability analysis of time scale filters

In this section, we consider a time scale dynamic system in a state space as

$$x^\Delta(t) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t), \quad (4.1)$$

where  $x(t), u(t), y(t)$  are the state, input and output respectively. Under finite word length effects and considering the fixed point arithmetic, we have

$$\tilde{x}^\Delta(t) = \text{fl}(Ax(t)) + \text{fl}(\tilde{B}\tilde{u}(t)), \quad \tilde{y}(t) = \text{fl}(\tilde{C}\tilde{x}(t)) + \text{fl}(\tilde{D}\tilde{u}(t)), \quad (4.2)$$

where  $\text{fl}(\cdot)$  and  $\sim$  denote the fixed point multiplicative operation and the round off/quantization operations on states and parameter matrices, respectively. In fixed-point-format, the word length of a real number can be split into three point: sign bit, bits of integer part and that of fractional part, denoted by  $\omega_s, \omega_i$ , and  $\omega_f$ , respectively. Thus, the total word length of fixed point form is given by

$$\omega = \omega_s + \omega_i + \omega_f = 1 + \omega_i + \omega_f.$$

When a real number is represented by a fixed point form, and we consider that the overflow effect is only limited by the saturation method, the round off error can be defined as  $\text{fl}(a) = a + \varepsilon_{rd}$ , where  $a$  is any real number and  $\varepsilon_{rd}$  denotes the round off error and is bounded by  $|\varepsilon_{rd}| < 2^{-\omega_f}$ , when two real numbers  $a$  and  $b$  are multiplied and assuming that  $a$  and  $b$  are both in the range  $(0, 1)$  and the fixed point error representation is described as

$$\text{fl}(ab) = ab + \hat{\varepsilon}_{op}, \quad (4.3)$$

where  $\hat{\varepsilon}_{op}$  is the operational error and is uniformly distributed in the range  $(-1, 1)$  denoted by  $\hat{\varepsilon}_{op}$ . Note that equation (4.3) can be written as  $\text{fl}(ab) = ab + \Delta \in op$ , where  $\Delta = z^{-\omega_f}$ . Let  $ap, bq \in \mathbb{R}^n$ , then the rounded inner product is give by

$$\begin{aligned} \sum_{i=1}^n \text{fl}(ap_i bq_i) &= ap_1 bq_1 + ap_2 bq_2 + \cdots + ap_n bq_n + \Delta (m(pq)_1 + m(pq)_2 + \cdots + m(pq)_n) \\ &\quad + \underline{\Delta} (ap_1 bq_1 + ap_2 bq_2 + \cdots + ap_n bq_n). \end{aligned}$$

Thus, we have  $\text{Variance}(m((pq)_i)) = \frac{1}{3}$ , where  $\text{var}(\cdot)$  denotes the variance of  $(\cdot)$ . Further for any two matrices  $E \in \mathbb{R}^{n \times n}$  and  $F \in \mathbb{R}^{n \times n}$ , we have

$$\text{fl}[EF] = EF + \Delta \begin{bmatrix} m(11) & m(12) & \cdots & m(1q) \\ m(21) & m(22) & \cdots & m(2q) \\ \cdots & \cdots & \cdots & \cdots \\ m(n1) & m(n2) & \cdots & m(nq) \end{bmatrix} = EF + \Delta M,$$

where  $M \in \mathbb{R}^{n \times q}$  is the stochastic matrix with elements which are uniformly distributed in  $(-1, 1)$  and are mutually independent, and satisfies

$$\|M\|_2 = \sqrt{\lambda_{\max}(MTMJ)}.$$

For obtaining specific quantization error bound and to derive stability criterion in terms of word length  $\omega$  in equation (4.1), we have the following

1. For the quantizational errors, suppose that  $p_{ij}$  is the  $(i, j)$  the element of a filter parameterization set  $P = \{A, B, C, D\}$ . When FWL effects occur and no over flow exists, we can have

$$\tilde{p}_{ij} = \text{Sgn}(p_{ij}) (|p_{ij}| + 2^{-\omega_f}), \quad \text{for } p_{ij} \text{ not an integer}, \quad (4.4)$$

where  $\text{sgn}(p_{ij})$  stands for the sign function of  $p_{ij}$ .



2. For computational error, error bound is given by

$$\|\Delta p\|_2.$$

We now turn our attention to stability analysis for fixed point dynamic filter implementation. Equation (4.2) can be written as

$$\tilde{x}(t) = A\tilde{x}(t) + (A - \tilde{A})\tilde{x}(t) + \Delta M_1 + \tilde{B}\tilde{U}(t) + \Delta M_2$$

and

$$\tilde{y}(t) = \tilde{C}\tilde{x}(t) + \Delta M_3 + \tilde{a}U(t) + \Delta M_4$$

for discrete digital filter implementation, it reduces to

$$\begin{aligned}\tilde{x}(n+1) &= A\tilde{x}(n) + (A - \tilde{A})\tilde{x}(n) + \Delta M_1 + \tilde{B}\tilde{U}(n) + \Delta M_2, \\ \tilde{y}(n) &= \tilde{C}\tilde{x}(n) + \Delta M_3 + \tilde{D}U(n) + \Delta M_4,\end{aligned}\tag{4.5}$$

where  $M_i$  ( $i = 1, 2, 3, 4$ ) are with stochastic elements uniformly distributed in  $(-1, 1)$  and are all mutually independent.

Note that, the implemented FWL digital continuous filter is stable, if and only if

$$\|A\|_2 + \|A - \tilde{A}\|_2 < 1.\tag{4.6}$$

**Theorem 4.1.** Based on equations (4.4) and (4.5), an estimated bit-number for controller implementation subject to stability criteria is given by

$$\omega(\text{est}) = 1 + \text{int} \left[ \log_2 \frac{\|\text{sgn}(A)\|_2}{1 - \|A\|_2} \right]$$

subject to  $0 < \|A\|_2 < 1$ , where  $\omega(\text{est})$  and  $\text{int}[\cdot]$  denote that estimated word-length and the smallest integer equal to or greater than 1, respectively.

*Proof.* By equation (4.4), we have

$$\|A - \tilde{A}\|_2 < 2^{-\omega_f} \|\text{sgn}(A)\|_2.$$

Using (4.6), the filter system is stable if and only if

$$2^{-\omega_f} \|\text{sgn}(A)\|_2 < 1 - \|A\|_2.$$

Since all parameters are normalized, the estimated word-length in terms of stability may be given as

$$\omega(\text{est}) = \text{sign}(\text{bit}) + \text{int}(\omega t).$$

The proof of the theorem is now complete. □

## References

- [1] B. Aulbach, S. Hilger, *A unified Approach to Continuous and Discrete Dynamics*, Qualitative theory of differential equations (Szeged), **1988** (1988), 37–56. 1
- [2] B. Aulbach, S. Hilger, *Linear dynamic processes with inhomogeneous time scale*, Nonlinear dynamics and quantum dynamical systems (Gaussig), **1990** (1990), 9–20. 1
- [3] M. Bohner, A. Peterson, *Dynamic Equations on Time Scales*, Birkhäuser, Boston, (2001). 2
- [4] M. Bohner, A. Peterson, *Advances in Dynamics Equations on Time Scales*, Birkhäuser, Boston, (2003). 2
- [5] K. V. V. Kanuri, R. Suryanarayana, K. N. Murty. *Existence of  $\Psi$ -bounded solutions for linear differential systems on time scales*, Journal of Mathematics and Computer Sciences, **20** (2020), 1–13. 1
- [6] V. Lakshmikantham, S. Sivasundaram, B. Kaymakalan, *Dynamic systems on measure chains*, Kluwer Academic Publishers Group, Dordrecht, (1996). 2
- [7] H. J. Ko, *Stability Analysis of Digital Filters Under Finite Word Length Effects via Normal-Form Transformation*, Asian J. Health Infor. Sci., **1** (2006), 112–121. 1
- [8] K. N. Murty, Y. Wu, V. Kanuri, *Metrics that suit for dichotomy, well conditioning and object oriented design on measure chains*, Nonlinear Stud., **18** (2011), 621–637. 1