

# $L_2$ -Sensitivity Minimization of One- and Two-Dimensional State-Space Digital Filters Subject to $L_2$ -Scaling Constraints

Takao Hinamoto, *Fellow, IEEE*, Ken-ichi Iwata, and Wu-Sheng Lu, *Fellow, IEEE*

**Abstract**—This paper investigates the problem of minimizing an  $L_2$ -sensitivity measure subject to  $L_2$ -norm dynamic-range scaling constraints for one-dimensional (1-D) as well as a class of two-dimensional (2-D) state-space digital filters, where the 2-D digital filters are described by a transposed structure of the Fornasini–Marchesini second local state-space model. In each case, a novel iterative technique is developed to solve the constraint optimization problem directly. The proposed solution methods are largely based on the use of a Lagrange function and some matrix-theoretic techniques. Numerical examples are presented to demonstrate the effectiveness of the proposed techniques.

**Index Terms**—A class of 2-D state-space digital filters, Lagrange's function,  $L_2$ -norm dynamic-range scaling constraints,  $L_2$ -sensitivity minimization, optimal realization, state-space digital filters.

## I. INTRODUCTION

FOR over three decades, the issue on finite word length (FWL) effects has been a significant research topic in the implementation of fixed-point state-space digital filters. A great deal of effort has been made to synthesize the optimal FWL state-space filter structures that minimize the FWL effects on the efficiency and performance of the filter actually implemented. This paper is concerned with the problem of minimizing coefficient sensitivity in state-space digital filters. Given a transfer function with infinite accuracy coefficients that meets specification requirements, we often need to implement its state-space model using a finite binary representation. This requires the truncation or rounding of coefficients in the state-space model to fit the FWL constraints. As a result, the characteristics of a stable filter might be so altered that the filter may become unstable. This motivates the study of the coefficient sensitivity minimization problem. In the literature, techniques for synthesizing the state-space descriptions that minimize the coefficient sensitivity can be divided into two main classes:  $L_1/L_2$ -mixed sensitivity minimization [1]–[5] and  $L_2$ -sensitivity minimization [6]–[10]. In [6]–[10], it has been argued that the sensitivity measure

based on the  $L_2$ -norm only is more natural and reasonable relative to the  $L_1/L_2$ -mixed sensitivity minimization. For 2-D state-space digital filters, the  $L_1/L_2$ -mixed sensitivity minimization problem [11]–[15] and  $L_2$ -sensitivity minimization problem [10], [16]–[19] have also been investigated. However, to our best knowledge, little has been done for the minimization of  $L_2$ -sensitivity subject to the  $L_2$ -norm dynamic-range scaling constraints for state-space digital filters [20], although it has been known that the use of scaling constraints can be beneficial for suppressing overflow oscillations [21], [22].

This paper investigates the problem of minimizing an  $L_2$ -sensitivity measure subject to  $L_2$ -norm dynamic-range scaling constraints for state-space digital filters as well as a class of 2-D state-space digital filters. To this end, an expression for evaluating the  $L_2$ -sensitivity is introduced, and the  $L_2$ -sensitivity minimization problem subject to the scaling constraints is formulated in each case. An iterative algorithm is then developed in each case to solve the constraint optimization problem directly. Our solution methods are largely based on the use of a Lagrange function and some matrix-theoretic techniques. Next, the coordinate transformation matrix is adjusted without altering the  $L_2$ -sensitivity to satisfy the scaling constraints, and is used to construct the optimal state-space filter structure that minimizes the  $L_2$ -sensitivity measure subject to the scaling constraints. Two numerical examples are presented to demonstrate the utility of the proposed algorithms.

Unlike the work reported in [20], the proposed iterative technique relies on neither converting the problem into an unconstrained optimization formulation nor using a quasi-Newton algorithm. From computer simulation results, it has turned out that the proposed iterative technique requires less than half amount of computations to attain practically the same convergence accuracy as compared with the technique reported in [20].

Throughout this paper,  $\mathbf{I}_n$  denotes the identity matrix of dimension  $n \times n$ . The transpose (conjugate transpose) of a matrix  $\mathbf{A}$  and trace of a square matrix  $\mathbf{A}$  are denoted by  $\mathbf{A}^T$  ( $\mathbf{A}^*$ ) and  $\text{tr}[\mathbf{A}]$ , respectively. The  $i$ th diagonal element of a square matrix  $\mathbf{A}$  is denoted by  $(\mathbf{A})_{ii}$ .

## II. $L_2$ -SENSITIVITY MINIMIZATION FOR 1-D FILTERS

### A. $L_2$ -Sensitivity Analysis

Consider a stable, controllable, and observable, state-space digital filter  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_n$  described by

$$\begin{aligned} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{b}u(k) \\ y(k) &= \mathbf{c}\mathbf{x}(k) + du(k) \end{aligned} \quad (1)$$

Manuscript received February 9, 2005; revised June 6, 2005. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Sheila S. Hemami.

T. Hinamoto is with the Graduate School of Engineering, Hiroshima University, Hiroshima 739-8527, Japan (e-mail: hinamoto@hiroshima-u.ac.jp).

K. Iwata was with the Graduate School of Engineering, Hiroshima University, Hiroshima 739-8527, Japan. He is now with Mitsubishi Electric Corporation, Himeji, Japan (e-mail: iwata@hiroshima-u.ac.jp).

W.-S. Lu is with the Department of Electrical and Computer Engineering, University of Victoria, Victoria, BC V8W 3P6, Canada (e-mail: wslu@ece.uvic.ca).

Digital Object Identifier 10.1109/TSP.2006.870584

where  $\mathbf{x}(k)$  is an  $n \times 1$  state-variable vector,  $u(k)$  is a scalar input,  $y(k)$  is a scalar output, and  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$  and  $d$  are real constant matrices of appropriate dimensions. The transfer function of the filter in (1) is given by

$$H(z) = \mathbf{c}(z\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{b} + d. \quad (2)$$

The  $L_2$ -sensitivity of the filter in (1) is defined as follows.

*Definition 1:* Let  $\mathbf{X}$  be an  $m \times n$  real matrix and let  $f(\mathbf{X})$  be a scalar complex function of  $\mathbf{X}$ , differentiable with respect to all the entries of  $\mathbf{X}$ . The sensitivity function of  $f$  with respect to  $\mathbf{X}$  is then defined as

$$\mathbf{S}_X = \frac{\partial f}{\partial \mathbf{X}}, \quad (\mathbf{S}_X)_{ij} = \frac{\partial f}{\partial x_{ij}} \quad (3)$$

where  $x_{ij}$  denotes the  $(i, j)$ th entry of matrix  $\mathbf{X}$ .

*Definition 2:* Let  $\mathbf{X}(z)$  be an  $m \times n$  complex matrix-valued function of a complex variable  $z$  and let  $x_{pq}(z)$  be the  $(p, q)$ th entry of  $\mathbf{X}(z)$ . The  $L_2$ -norm of  $\mathbf{X}(z)$  is then defined as

$$\begin{aligned} \|\mathbf{X}(z)\|_2 &= \left[ \frac{1}{2\pi} \int_0^{2\pi} \sum_{p=1}^m \sum_{q=1}^n |x_{pq}(e^{j\omega})|^2 d\omega \right]^{\frac{1}{2}} \\ &= \left( \text{tr} \left[ \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{X}(z)\mathbf{X}^*(z) \frac{dz}{z} \right] \right)^{\frac{1}{2}}. \end{aligned} \quad (4)$$

From (2) and Definitions 1 and 2, the overall  $L_2$ -sensitivity measure for the filter in (1) is defined as

$$\begin{aligned} S_1 &= \left\| \frac{\partial H(z)}{\partial \mathbf{A}} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial \mathbf{b}} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial \mathbf{c}^T} \right\|_2^2 \\ &= \left\| [\mathbf{F}(z)\mathbf{G}(z)]^T \right\|_2^2 + \left\| \mathbf{G}^T(z) \right\|_2^2 + \|\mathbf{F}(z)\|_2^2 \end{aligned} \quad (5)$$

where

$$\mathbf{F}(z) = (z\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{b}, \quad \mathbf{G}(z) = \mathbf{c}(z\mathbf{I}_n - \mathbf{A})^{-1}.$$

The term  $d$  in (2) and the sensitivity with respect to it are coordinate independent, and therefore they are neglected here.

It is easy to show that the  $L_2$ -sensitivity measure in (5) can be expressed as

$$S_1 = \text{tr} [\mathbf{M}(\mathbf{I}_n)] + \text{tr}[\mathbf{W}_o] + \text{tr}[\mathbf{K}_c] \quad (6)$$

where

$$\begin{aligned} \mathbf{K}_c &= \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{F}(z)\mathbf{F}^T(z^{-1}) \frac{dz}{z} \\ \mathbf{W}_o &= \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{G}^T(z)\mathbf{G}(z^{-1}) \frac{dz}{z} \\ \mathbf{M}(\mathbf{P}) &= \frac{1}{2\pi j} \oint_{|z|=1} [\mathbf{F}(z)\mathbf{G}(z)]^T \mathbf{P}^{-1} \mathbf{F}(z^{-1})\mathbf{G}(z^{-1}) \frac{dz}{z}. \end{aligned}$$

The matrices  $\mathbf{K}_c$  and  $\mathbf{W}_o$  in (6) are called the controllability and observability Gramians, respectively, and can be obtained by solving the following Lyapunov equations [23]:

$$\begin{aligned} \mathbf{K}_c &= \mathbf{A}\mathbf{K}_c\mathbf{A}^T + \mathbf{b}\mathbf{b}^T \\ \mathbf{W}_o &= \mathbf{A}^T\mathbf{W}_o\mathbf{A} + \mathbf{c}^T\mathbf{c}. \end{aligned} \quad (7)$$

In order to evaluate the term  $\text{tr}[\mathbf{M}(\mathbf{I}_n)]$  in (6), one only needs to compute  $\mathbf{M}(\mathbf{P})$  with  $\mathbf{P} = \mathbf{I}_n$ . However, the general solution method for computing  $\mathbf{M}(\mathbf{P})$  with an arbitrary positive-definite matrix  $\mathbf{P}$  will be required in our subsequent development. We shall address this issue shortly.

If a coordinate transformation defined by

$$\bar{\mathbf{x}}(k) = \mathbf{T}^{-1}\mathbf{x}(k) \quad (8)$$

is applied to the filter in (1), then the new realization  $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_n$  can be characterized by

$$\begin{aligned} \bar{\mathbf{A}} &= \mathbf{T}^{-1}\mathbf{A}\mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1}\mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c}\mathbf{T} \\ \bar{\mathbf{K}}_c &= \mathbf{T}^{-1}\mathbf{K}_c\mathbf{T}^{-T}, \quad \bar{\mathbf{W}}_o = \mathbf{T}^T\mathbf{W}_o\mathbf{T}. \end{aligned} \quad (9)$$

From (2) and (9), it is clear that the transfer function  $H(z)$  is invariant under the coordinate transformation in (8). In addition, under the coordinate transformation in (8), the Gramian  $\mathbf{M}(\mathbf{I}_n)$  becomes  $\mathbf{T}^T\mathbf{M}(\mathbf{P})\mathbf{T}$  and the  $L_2$ -sensitivity measure in (6) is changed to

$$S_1(\mathbf{P}) = \text{tr} [\mathbf{M}(\mathbf{P})\mathbf{P}] + \text{tr}[\mathbf{W}_o\mathbf{P}] + \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] \quad (10)$$

where  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ . Noting that

$$\begin{aligned} \bar{\mathbf{F}}(z)\bar{\mathbf{G}}(z) &= \mathbf{T}^{-1}\mathbf{F}(z)\mathbf{G}(z)\mathbf{T} \\ &= [\mathbf{T}^{-1} \mathbf{0}] \begin{bmatrix} z\mathbf{I}_n - \mathbf{A} & -\mathbf{bc} \\ \mathbf{0} & z\mathbf{I}_n - \mathbf{A} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{0} \\ \mathbf{T} \end{bmatrix} \end{aligned} \quad (11)$$

where

$$\bar{\mathbf{F}}(z) = (z\mathbf{I}_n - \bar{\mathbf{A}})^{-1}\bar{\mathbf{b}}, \quad \bar{\mathbf{G}}(z) = \bar{\mathbf{c}}(z\mathbf{I}_n - \bar{\mathbf{A}})^{-1}$$

and denoting the observability Gramian of the system in (11) by  $\mathbf{Y}$ , it can be shown for an arbitrary  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ , the matrix  $\mathbf{M}(\mathbf{P})$  can be obtained by solving the Lyapunov equation

$$\mathbf{Y} = \begin{bmatrix} \mathbf{A} & \mathbf{bc} \\ \mathbf{0} & \mathbf{A} \end{bmatrix}^T \mathbf{Y} \begin{bmatrix} \mathbf{A} & \mathbf{bc} \\ \mathbf{0} & \mathbf{A} \end{bmatrix} + \begin{bmatrix} \mathbf{P}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (12)$$

and then taking the lower-right  $n \times n$  block of  $\mathbf{Y}$  as  $\mathbf{M}(\mathbf{P})$ , i.e.,

$$\mathbf{M}(\mathbf{P}) = [\mathbf{0} \quad \mathbf{I}_n] \mathbf{Y} \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_n \end{bmatrix}. \quad (13)$$

Moreover, if the  $L_2$ -norm dynamic-range scaling constraints are imposed on the new state-variable vector  $\bar{\mathbf{x}}(k)$ , it is required that

$$(\bar{\mathbf{K}}_c)_{ii} = (\mathbf{T}^{-1}\mathbf{K}_c\mathbf{T}^{-T})_{ii} = 1, \quad \text{for } i = 1, 2, \dots, n. \quad (14)$$

The problem of  $L_2$ -sensitivity minimization subject to  $L_2$ -norm dynamic-range scaling constraints can now be formulated as follows: *For given  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$ , obtain an  $n \times n$*

nonsingular matrix  $\mathbf{T}$  that minimizes the sensitivity measure  $S_1(\mathbf{P})$  in (10) subject to the scaling constraints in (14).

### B. $L_2$ -Sensitivity Minimization

The problem of minimizing  $S_1(\mathbf{P})$  in (10) subject to the constraints in (14) is a constrained nonlinear optimization problem where the variable matrix is  $\mathbf{P}$ . If we sum the  $n$  constraints in (14) up, then we have

$$\text{tr}[\mathbf{T}^{-1}\mathbf{K}_c\mathbf{T}^{-T}] = \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] = n. \quad (15)$$

Consequently, the problem of minimizing (10) subject to the constraints in (14) can be “relaxed” into the following problem:

$$\begin{aligned} &\text{minimize } S_1(\mathbf{P}) \text{ in (10)} \\ &\text{subject to } \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] = n. \end{aligned} \quad (16)$$

Although clearly a solution of problem (16) is not necessarily a solution of the problem of minimizing (10) subject to the constraints in (14), it is important to emphasize that the ultimate solution we seek is not matrix  $\mathbf{P}$  but a nonsingular matrix  $\mathbf{T}$  that is related to the solution of the problem of minimizing (10) subject to the constraints in (14) as  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ . If matrix  $\mathbf{P}$  is a solution of problem (16) and  $\mathbf{P}^{1/2}$  denotes a matrix square root of  $\mathbf{P}$ , i.e.,  $\mathbf{P} = \mathbf{P}^{1/2}\mathbf{P}^{1/2}$ , then it is easy to see that any matrix  $\mathbf{T}$  of the form  $\mathbf{T} = \mathbf{P}^{1/2}\mathbf{U}$ , where  $\mathbf{U}$  is an arbitrary orthogonal matrix, still holds the relation  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ . As will be shown shortly, under the constraint  $\text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] = n$  in (16), there exists an orthogonal matrix  $\mathbf{U}$  such that matrix  $\mathbf{T} = \mathbf{P}^{1/2}\mathbf{U}$  satisfies the constraints in (14), where  $\mathbf{P}^{1/2}$  is a square root of the solution matrix  $\mathbf{P}$  for problem (16).

It is for these reasons, we now address problem (16) as the first step of our solution strategy. To solve (16), we define the Lagrange function of the problem as

$$J_1(\mathbf{P}, \lambda) = \text{tr}[\mathbf{M}(\mathbf{P})\mathbf{P}] \text{tr}[\mathbf{W}_o\mathbf{P}] + \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] + \lambda \left( \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] - n \right) \quad (17)$$

where  $\lambda$  is a Lagrange multiplier. It is well known that the solution of problem (16) must satisfy the Karush–Kuhn–Tucker (KKT) conditions  $\partial J_1(\mathbf{P}, \lambda)/\partial \mathbf{P} = \mathbf{0}$  and  $\partial J_1(\mathbf{P}, \lambda)/\partial \lambda = 0$ , where the gradients are found to be [24]

$$\begin{aligned} \frac{\partial J_1(\mathbf{P}, \lambda)}{\partial \mathbf{P}} &= \mathbf{M}(\mathbf{P}) - \mathbf{P}^{-1}\mathbf{N}(\mathbf{P})\mathbf{P}^{-1} \\ &\quad + \mathbf{W}_o - (\lambda + 1)\mathbf{P}^{-1}\mathbf{K}_c\mathbf{P}^{-1} \\ \frac{\partial J_1(\mathbf{P}, \lambda)}{\partial \lambda} &= \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] - n \end{aligned} \quad (18)$$

where  $\mathbf{N}(\mathbf{P})$  can be obtained by solving the following Lyapunov equation:

$$\mathbf{Z} = \begin{bmatrix} \mathbf{A} & \mathbf{bc} \\ \mathbf{0} & \mathbf{A} \end{bmatrix} \mathbf{Z} \begin{bmatrix} \mathbf{A} & \mathbf{bc} \\ \mathbf{0} & \mathbf{A} \end{bmatrix}^T + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{P} \end{bmatrix}$$

and then taking the upper-left  $n \times n$  block of  $\mathbf{Z}$ , i.e.,

$$\mathbf{N}(\mathbf{P}) = [\mathbf{I}_n \quad \mathbf{0}]\mathbf{Z} \begin{bmatrix} \mathbf{I}_n \\ \mathbf{0} \end{bmatrix}.$$

Note that matrix  $\mathbf{Z}$  corresponds to the controllability Gramian of the system in (11).

Hence, the KKT conditions become

$$\begin{aligned} \mathbf{P}\mathbf{F}(\mathbf{P})\mathbf{P} &= \mathbf{G}(\mathbf{P}, \lambda) \\ \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] &= n \end{aligned} \quad (19)$$

where

$$\begin{aligned} \mathbf{F}(\mathbf{P}) &= \mathbf{M}(\mathbf{P}) + \mathbf{W}_o \\ \mathbf{G}(\mathbf{P}, \lambda) &= \mathbf{N}(\mathbf{P}) + (\lambda + 1)\mathbf{K}_c. \end{aligned}$$

The first equation in (19) is highly nonlinear with respect to  $\mathbf{P}$ . An effective approach to solving the first equation in (19) is to “relax” it into the following recursive second-order matrix equation:

$$\mathbf{P}_{i+1}\mathbf{F}(\mathbf{P}_i)\mathbf{P}_{i+1} = \mathbf{G}(\mathbf{P}_i, \lambda_i) \quad (20)$$

where  $\mathbf{P}_i$  is assumed to be known from the previous recursion. Noting that  $\mathbf{P}\mathbf{W}\mathbf{P} = \mathbf{M}$  has the unique solution [5]

$$\mathbf{P} = \mathbf{W}^{-\frac{1}{2}} \left[ \mathbf{W}^{\frac{1}{2}} \mathbf{M} \mathbf{W}^{\frac{1}{2}} \right]^{\frac{1}{2}} \mathbf{W}^{-\frac{1}{2}} \quad (21)$$

where  $\mathbf{W} > 0$  and  $\mathbf{M} \geq 0$  are symmetric, the solution  $\mathbf{P}_{i+1}$  of (20) is given by

$$\mathbf{P}_{i+1} = \mathbf{F}(\mathbf{P}_i)^{-\frac{1}{2}} \left[ \mathbf{F}(\mathbf{P}_i)^{\frac{1}{2}} \mathbf{G}(\mathbf{P}_i, \lambda_i) \mathbf{F}(\mathbf{P}_i)^{\frac{1}{2}} \right]^{\frac{1}{2}} \mathbf{F}(\mathbf{P}_i)^{-\frac{1}{2}}. \quad (22)$$

To derive a recursive formula for the Lagrange multiplier  $\lambda$ , we employ (19) to write

$$\text{tr}[\mathbf{P}\mathbf{F}(\mathbf{P})] = \text{tr}[\mathbf{N}(\mathbf{P})\mathbf{P}^{-1}] + n(\lambda + 1) \quad (23)$$

which naturally suggests the following recursion for  $\lambda$ :

$$\lambda_{i+1} = \frac{\text{tr}[\mathbf{P}_i\mathbf{F}(\mathbf{P}_i)] - \text{tr}[\mathbf{N}(\mathbf{P}_i)\mathbf{P}_i^{-1}]}{n} - 1. \quad (24)$$

In the above algorithm,  $\lambda_i$  is obtained from the previous iteration. The iteration process starts with  $\mathbf{P}_1 = \mathbf{I}_n$  and any value of  $\lambda_1 > 0$  and continues until (19) is satisfied within a prescribed numerical tolerance.

As the second step of the solution strategy, we now turn our attention to the construction of the optimal coordinate transformation matrix  $\mathbf{T}$  that solves the problem of minimizing (10) subject to the constraints in (14). As analyzed earlier, the optimal  $\mathbf{T}$  assumes the form

$$\mathbf{T} = \mathbf{P}^{\frac{1}{2}}\mathbf{U} \quad (25)$$

where  $\mathbf{P}^{1/2}$  is the square root of the matrix  $\mathbf{P}$  obtained above, and  $\mathbf{U}$  is an  $n \times n$  orthogonal matrix to be determined as follows. From (9) and (25), it follows that

$$\begin{aligned} \bar{\mathbf{K}}_c &= \mathbf{T}^{-1}\mathbf{K}_c\mathbf{T}^{-T} \\ &= \mathbf{U}^T \mathbf{P}^{-\frac{1}{2}} \mathbf{K}_c \mathbf{P}^{-\frac{1}{2}} \mathbf{U}. \end{aligned} \quad (26)$$

In order to find an orthogonal matrix  $\mathbf{U}$  such that the matrix  $\bar{\mathbf{K}}_c$  in (26) satisfies the constraints in (14), we perform the eigenvalue-eigenvector decomposition for the positive definite matrix  $\mathbf{P}^{-1/2}\mathbf{K}_c\mathbf{P}^{-1/2}$  as

$$\mathbf{P}^{-\frac{1}{2}}\mathbf{K}_c\mathbf{P}^{-\frac{1}{2}} = \mathbf{R}\mathbf{O}\mathbf{R}^T \quad (27)$$

where  $\Theta = \text{diag}\{\theta_1, \theta_2, \dots, \theta_n\}$  with  $\theta_i > 0$  and  $\mathbf{R}$  is an orthogonal matrix. Next, an orthogonal matrix  $\mathbf{S}$  such that

$$\mathbf{S}\Theta\mathbf{S}^T = \begin{bmatrix} 1 & * & \cdots & * \\ * & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & * \\ * & \cdots & * & 1 \end{bmatrix} \quad (28)$$

can be obtained by numerical manipulations [22, p. 278]. Using (26), (27), and (28), it can be readily verified that the orthogonal matrix  $\mathbf{U} = \mathbf{R}\mathbf{S}^T$  leads to a  $\bar{\mathbf{K}}_c$  in (26) whose diagonal elements are equal to unity; hence, the constraints in (14) are now satisfied. This matrix  $\mathbf{T}$  together with (25) gives the solution of the problem of minimizing (10) subject to the constraints in (14) as

$$\mathbf{T} = \mathbf{P}^{\frac{1}{2}}\mathbf{R}\mathbf{S}^T. \quad (29)$$

### III. $L_2$ -SENSITIVITY MINIMIZATION FOR 2-D DIGITAL FILTERS

#### A. $L_2$ -Sensitivity Analysis

Consider a local state–space model  $(\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}, \mathbf{c}_1, \mathbf{c}_2, d)_n$  [25] for a class of 2-D state–space digital filters that is described by

$$\begin{bmatrix} \mathbf{x}(i+1, j+1) \\ y(i, j) \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{c}_1 & \mathbf{c}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}(i, j+1) \\ \mathbf{x}(i+1, j) \end{bmatrix} + \begin{bmatrix} \mathbf{b} \\ d \end{bmatrix} u(i, j) \quad (30)$$

where  $\mathbf{x}(i, j)$  is an  $n \times 1$  local state vector,  $u(i, j)$  is a scalar input,  $y(i, j)$  is a scalar output, and  $\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}, \mathbf{c}_1, \mathbf{c}_2$  and  $d$  are real constant matrices of appropriate dimensions. The 2-D filter in (30) is assumed to be stable, locally controllable, and locally observable. The transfer function of the 2-D filter in (30) is given by

$$H(z_1, z_2) = (z_1^{-1}\mathbf{c}_1 + z_2^{-1}\mathbf{c}_2) \cdot (\mathbf{I}_n - z_1^{-1}\mathbf{A}_1 - z_2^{-1}\mathbf{A}_2)^{-1} \mathbf{b} + d. \quad (31)$$

A block diagram of the local state–space (LSS) model in (30) is shown in Fig. 1. It is interesting to note that

$$H^T(z_1, z_2) = \mathbf{b}^T \left( \mathbf{I}_n - z_1^{-1}\mathbf{A}_1^T - z_2^{-1}\mathbf{A}_2^T \right)^{-1} \cdot (z_1^{-1}\mathbf{c}_1^T + z_2^{-1}\mathbf{c}_2^T) + d \quad (32)$$

can be viewed as a transfer function of the Fornasini–Marchesini second LSS model [26]. Since  $H(z_1, z_2) = H^T(z_1, z_2)$ , the LSS model in (30) corresponds to a transposed structure of the Fornasini–Marchesini second LSS model.

From (31) and Definition 1, it can easily be verified that

$$\begin{aligned} \frac{\partial H(z_1, z_2)}{\partial \mathbf{A}_k} &= z_k^{-1} [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T \\ \frac{\partial H(z_1, z_2)}{\partial \mathbf{b}} &= \mathbf{G}^T(z_1, z_2) \\ \frac{\partial H(z_1, z_2)}{\partial \mathbf{c}_k^T} &= z_k^{-1} \mathbf{F}(z_1, z_2), \quad k = 1, 2 \end{aligned} \quad (33)$$

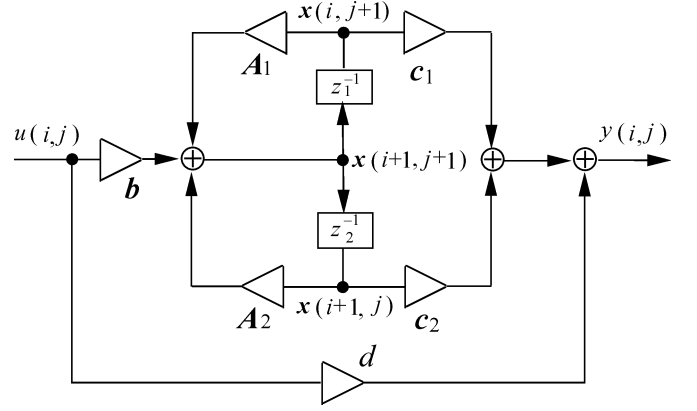


Fig. 1. LSS model for 2-D filters.

where

$$\begin{aligned} \mathbf{F}(z_1, z_2) &= (\mathbf{I}_n - z_1^{-1}\mathbf{A}_1 - z_2^{-1}\mathbf{A}_2)^{-1} \mathbf{b} \\ \mathbf{G}(z_1, z_2) &= (z_1^{-1}\mathbf{c}_1 + z_2^{-1}\mathbf{c}_2) \cdot (\mathbf{I}_n - z_1^{-1}\mathbf{A}_1 - z_2^{-1}\mathbf{A}_2)^{-1}. \end{aligned}$$

The term  $d$  in (31) and its sensitivity are independent on the LSS coordinate, and therefore they are neglected here.

*Definition 3:* Let  $\mathbf{X}(z_1, z_2)$  be an  $m \times n$  complex matrix valued function of the complex variables  $z_1$  and  $z_2$ . The  $L_2$  norm of  $\mathbf{X}(z_1, z_2)$  is then defined as

$$\begin{aligned} \|\mathbf{X}(z_1, z_2)\|_2 &= \left[ \frac{1}{(2\pi)^2} \int_0^{2\pi} \int_0^{2\pi} \sum_{p=1}^m \sum_{q=1}^n |x_{pq}(e^{j\omega_1}, e^{j\omega_2})|^2 d\omega_1 d\omega_2 \right]^{\frac{1}{2}} \\ &= \left( \text{tr} \left[ \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{X}(z_1, z_2) \mathbf{X}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \right] \right)^{\frac{1}{2}} \end{aligned} \quad (34)$$

where  $\Gamma_i = \{z_i : |z_i| = 1\}$  for  $i = 1, 2$ .

From (33) and Definition 3, the overall  $L_2$ -sensitivity measure for the 2-D filter in (30) is evaluated by

$$\begin{aligned} S_2 &= \sum_{k=1}^2 \left\| \frac{\partial H(z_1, z_2)}{\partial \mathbf{A}_k} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2)}{\partial \mathbf{b}} \right\|_2^2 \\ &\quad + \sum_{k=1}^2 \left\| \frac{\partial H(z_1, z_2)}{\partial \mathbf{c}_k^T} \right\|_2^2 \\ &= 2 \left\| [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T \right\|_2^2 + \left\| \mathbf{G}^T(z_1, z_2) \right\|_2^2 \\ &\quad + 2 \left\| \mathbf{F}(z_1, z_2) \right\|_2^2. \end{aligned} \quad (35)$$

The  $L_2$ -sensitivity measure in (35) can be written as

$$S_2 = 2\text{tr}[\mathbf{M}(\mathbf{I}_n)] + \text{tr}[\mathbf{W}_o] + 2\text{tr}[\mathbf{K}_c] \quad (36)$$

where

$$\begin{aligned} \mathbf{K}_c &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{F}(z_1, z_2) \mathbf{F}^T(z_1^{-1}, z_2^{-1}) \frac{dz_1 dz_2}{z_1 z_2} \\ \mathbf{W}_o &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{G}^T(z_1, z_2) \mathbf{G}(z_1^{-1}, z_2^{-1}) \frac{dz_1 dz_2}{z_1 z_2} \end{aligned}$$

$$\mathbf{M}(\mathbf{P}) = \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T \mathbf{P}^{-1} \cdot \mathbf{F}(z_1^{-1}, z_2^{-1}) \mathbf{G}(z_1^{-1}, z_2^{-1}) \frac{dz_1 dz_2}{z_1 z_2}.$$

Matrices  $\mathbf{K}_c$ ,  $\mathbf{W}_o$ , and  $\mathbf{M}(\mathbf{P})$  are called 2-D Gramians and can be derived from

$$\begin{aligned} \mathbf{K}_c &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}(i, j) \mathbf{f}^T(i, j) \\ \mathbf{W}_o &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{g}^T(i, j) \mathbf{g}(i, j) \\ \mathbf{M}(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{H}^T(i, j) \mathbf{P}^{-1} \mathbf{H}(i, j) \end{aligned} \quad (37)$$

where

$$\begin{aligned} \mathbf{f}(i, j) &= \mathbf{A}^{(i, j)} \mathbf{b} \\ \mathbf{g}(i, j) &= \mathbf{c}_1 \mathbf{A}^{(i-1, j)} + \mathbf{c}_2 \mathbf{A}^{(i, j-1)} \\ \mathbf{A}^{(0, 0)} &= \mathbf{I}_n, \quad \mathbf{A}^{(i, j)} = \mathbf{0}, \quad i < 0 \text{ or } j < 0 \\ \mathbf{A}^{(i, j)} &= \mathbf{A}_1 \mathbf{A}^{(i-1, j)} + \mathbf{A}_2 \mathbf{A}^{(i, j-1)} \\ &= \mathbf{A}^{(i-1, j)} \mathbf{A}_1 + \mathbf{A}^{(i, j-1)} \mathbf{A}_2, \quad (i, j) > (0, 0) \\ \mathbf{H}(i, j) &= \sum_{(0, 0) \leq (k, r) \leq (i, j)} \mathbf{f}(k, r) \mathbf{g}(i-k, j-r) \end{aligned}$$

with the partial ordering for integer pairs  $(i, j)$  used in [27, p. 2].

If a coordinate transformation defined by

$$\bar{\mathbf{x}}(i, j) = \mathbf{T}^{-1} \mathbf{x}(i, j) \quad (38)$$

is applied to the 2-D filter in (30), we obtain a new realization  $(\bar{\mathbf{A}}_1, \bar{\mathbf{A}}_2, \bar{\mathbf{b}}, \bar{\mathbf{c}}_1, \bar{\mathbf{c}}_2, d)_n$  characterized by

$$\begin{aligned} \bar{\mathbf{A}}_1 &= \mathbf{T}^{-1} \mathbf{A}_1 \mathbf{T}, \quad \bar{\mathbf{A}}_2 = \mathbf{T}^{-1} \mathbf{A}_2 \mathbf{T} \\ \bar{\mathbf{b}} &= \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}}_1 = \mathbf{c}_1 \mathbf{T}, \quad \bar{\mathbf{c}}_2 = \mathbf{c}_2 \mathbf{T} \\ \bar{\mathbf{K}}_c &= \mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T}, \quad \bar{\mathbf{W}}_o = \mathbf{T}^T \mathbf{W}_o \mathbf{T}. \end{aligned} \quad (39)$$

Noting that the coordinate transformation in (38) transforms the Gramian  $\mathbf{M}(\mathbf{I}_n)$  into  $\mathbf{T}^T \mathbf{M}(\mathbf{P}) \mathbf{T}$ , (36) becomes

$$S_2(\mathbf{P}) = 2\text{tr}[\mathbf{M}(\mathbf{P})\mathbf{P}] + \text{tr}[\mathbf{W}_o\mathbf{P}] + 2\text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] \quad (40)$$

where  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ . Moreover, if the  $L_2$ -norm dynamic-range scaling constraints are imposed on the LSS vector  $\bar{\mathbf{x}}(i, j)$ , then

$$(\bar{\mathbf{K}}_c)_{ii} = (\mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T})_{ii} = 1, \quad \text{for } i = 1, 2, \dots, n \quad (41)$$

are required.

The problem considered here is as follows: *For given  $\mathbf{A}_1$ ,  $\mathbf{A}_2$ ,  $\mathbf{b}$ ,  $\mathbf{c}_1$ , and  $\mathbf{c}_2$ , obtain an  $n \times n$  nonsingular matrix  $\mathbf{T}$  which minimizes (40) subject to the scaling constraints in (41).*

### B. $L_2$ -Sensitivity Minimization

In order to minimize (40) over an  $n \times n$  symmetric positive-definite matrix  $\mathbf{P}$  subject to the ‘‘relaxed’’ constraint

$\text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] = n$  derived from (41), we define the Lagrange function

$$J_2(\mathbf{P}, \lambda) = 2\text{tr}[\mathbf{M}(\mathbf{P})\mathbf{P}] + \text{tr}[\mathbf{W}_o\mathbf{P}] + 2\text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] + \lambda (\text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] - n) \quad (42)$$

where  $\lambda$  is a Lagrange multiplier. As is well known, the solution of the minimization problem of (42) must satisfy the KKT conditions  $\partial J_2(\mathbf{P}, \lambda)/\partial \mathbf{P} = \mathbf{0}$  and  $\partial J_2(\mathbf{P}, \lambda)/\partial \lambda = 0$ , where the gradients are found to be [24]

$$\begin{aligned} \frac{\partial J_2(\mathbf{P}, \lambda)}{\partial \mathbf{P}} &= 2\mathbf{M}(\mathbf{P}) - 2\mathbf{P}^{-1} \mathbf{N}(\mathbf{P}) \mathbf{P}^{-1} \\ &\quad + \mathbf{W}_o - (\lambda + 2) \mathbf{P}^{-1} \mathbf{K}_c \mathbf{P}^{-1} \\ \frac{\partial J_2(\mathbf{P}, \lambda)}{\partial \lambda} &= \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] - n \end{aligned} \quad (43)$$

where  $\mathbf{N}(\mathbf{P})$  is derived from

$$\mathbf{N}(\mathbf{P}) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{H}(i, j) \mathbf{P} \mathbf{H}^T(i, j).$$

Hence the KKT conditions become

$$\begin{aligned} \mathbf{P} \mathbf{F}(\mathbf{P}) \mathbf{P} &= \mathbf{G}(\mathbf{P}, \lambda) \\ \text{tr}[\mathbf{K}_c\mathbf{P}^{-1}] &= n \end{aligned} \quad (44)$$

where

$$\begin{aligned} \mathbf{F}(\mathbf{P}) &= 2\mathbf{M}(\mathbf{P}) + \mathbf{W}_o \\ \mathbf{G}(\mathbf{P}, \lambda) &= 2\mathbf{N}(\mathbf{P}) + (\lambda + 2) \mathbf{K}_c. \end{aligned}$$

The first equation in (44) is highly nonlinear with respect to  $\mathbf{P}$ . An effective approach to solving the first equation in (44) is to *relax* it into the following recursive second-order matrix equation:

$$\mathbf{P}_{i+1} \mathbf{F}(\mathbf{P}_i) \mathbf{P}_{i+1} = \mathbf{G}(\mathbf{P}_i, \lambda_i) \quad (45)$$

where  $\mathbf{P}_i$  is assumed to be known from the previous recursion and then the solution  $\mathbf{P}_{i+1}$  is given by

$$\mathbf{P}_{i+1} = \mathbf{F}(\mathbf{P}_i)^{-\frac{1}{2}} \left[ \mathbf{F}(\mathbf{P}_i)^{\frac{1}{2}} \mathbf{G}(\mathbf{P}_i, \lambda_i) \mathbf{F}(\mathbf{P}_i)^{\frac{1}{2}} \right]^{\frac{1}{2}} \mathbf{F}(\mathbf{P}_i)^{-\frac{1}{2}}. \quad (46)$$

To derive a recursive formula for the Lagrange multiplier  $\lambda$ , we use (44) to write

$$\text{tr}[\mathbf{P} \mathbf{F}(\mathbf{P})] = 2\text{tr}[\mathbf{N}(\mathbf{P}) \mathbf{P}^{-1}] + n(\lambda + 2) \quad (47)$$

that naturally reveals the following recursion for  $\lambda$ :

$$\lambda_{i+1} = \frac{\text{tr}[\mathbf{P}_i \mathbf{F}(\mathbf{P}_i)] - 2\text{tr}[\mathbf{N}(\mathbf{P}_i) \mathbf{P}_i^{-1}]}{n} - 2. \quad (48)$$

The iteration process starts with  $\mathbf{P}_1 = \mathbf{I}_n$  and any value of  $\lambda_1 > 0$  and continues until (44) is satisfied within a prescribed numerical tolerance.

Finally, the process in (25)–(29) is applied to the resulting optimal matrix  $\mathbf{P}$  in order to construct the optimal transformation matrix  $\mathbf{T} = \mathbf{P}^{1/2} \mathbf{R} \mathbf{S}^T$  minimizing (40) subject to the constraints in (41).

## IV. NUMERICAL EXAMPLES

*Example 1:* Let a state-space digital filter in (1) be specified by

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0.453770 & -1.556160 & 1.974860 \end{bmatrix} \\ \mathbf{b} &= [0 \ 0 \ 0.242096]^T \\ \mathbf{c} &= [0.095706 \ 0.095086 \ 0.327556] \\ d &= 0.015940. \end{aligned}$$

Performing the computation of (7), (12) and (13), the Gramians  $\mathbf{K}_c$ ,  $\mathbf{W}_o$ , and  $\mathbf{M}(\mathbf{I}_3)$  are calculated as

$$\begin{aligned} \mathbf{K}_c &= \begin{bmatrix} 1.000000 & 0.872501 & 0.562821 \\ 0.872501 & 1.000000 & 0.872501 \\ 0.562821 & 0.872501 & 1.000000 \end{bmatrix} \\ \mathbf{W}_o &= \begin{bmatrix} 0.820741 & -2.035328 & 1.628161 \\ -2.035328 & 5.307273 & -4.264903 \\ 1.628161 & -4.264903 & 3.941491 \end{bmatrix} \\ \mathbf{M}(\mathbf{I}_3) &= \begin{bmatrix} 8.921380 & -22.046457 & 17.916285 \\ -22.046457 & 55.671710 & -46.052011 \\ 17.916285 & -46.052011 & 42.522082 \end{bmatrix}. \end{aligned}$$

The  $L_2$ -sensitivity measure in (6) is then computed as

$$S_1 = 120.184677.$$

Choosing  $\mathbf{P}_1 = \mathbf{I}_3$  and  $\lambda_1 = 100$  in (22) and (24) as the initial estimates, it took the proposed iterative algorithm 500 iterations to converge to

$$\mathbf{P} = \begin{bmatrix} 2.307529 & 1.375667 & 0.514400 \\ 1.375667 & 1.103115 & 0.678193 \\ 0.514400 & 0.678193 & 0.666912 \end{bmatrix}$$

which yields

$$\mathbf{T} = \begin{bmatrix} 0.906372 & 0.756223 & 0.956110 \\ 0.196978 & 0.857123 & 0.574155 \\ -0.369823 & 0.597630 & 0.415910 \end{bmatrix}.$$

In this case, the Gramian  $\mathbf{M}(\mathbf{P})$  is computed from (12) and (13) as

$$\mathbf{M}(\mathbf{P}) = \begin{bmatrix} 1.908677 & -0.301984 & -1.313686 \\ -0.301984 & 1.701052 & 0.430349 \\ -1.313686 & 0.430349 & 1.395025 \end{bmatrix}$$

and the  $L_2$ -sensitivity measure in (10) is minimized subject to the scaling constraints in (14) to

$$S_1(\mathbf{P}) = 8.672129.$$

Profiles of the  $L_2$ -sensitivity, parameter  $\lambda$ , as well as  $\text{tr}[\mathbf{K}_c \mathbf{P}^{-1}]$  during the first 500 iterations of the algorithm are shown in Figs. 2 and 3, respectively. Together, these figures clearly reveal a two-stage convergence behavior of the algorithm in that the first stage (which consists of just one iteration) of the algorithm reduces the  $L_2$ -sensitivity drastically without maintaining

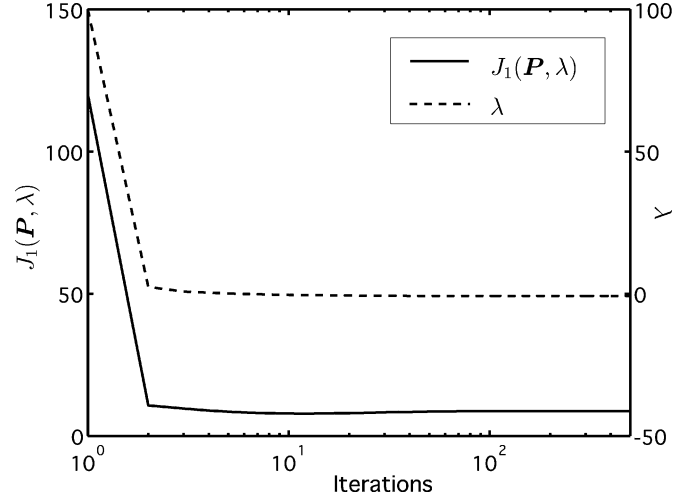


Fig. 2.  $L_2$ -sensitivity and  $\lambda$  performances.

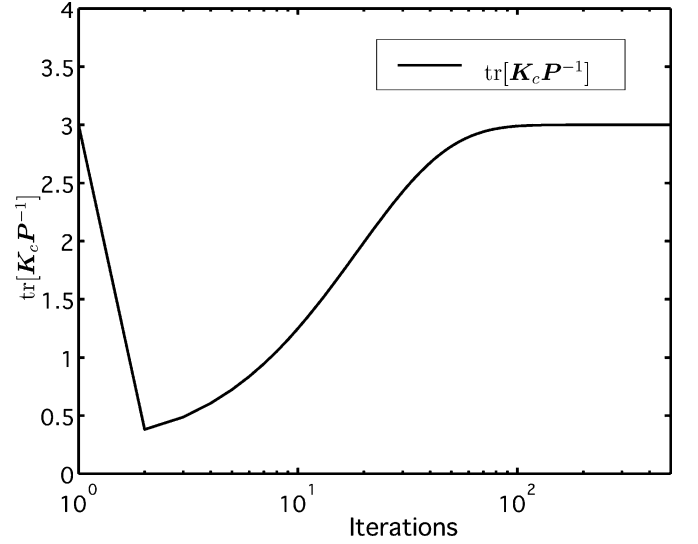


Fig. 3.  $\text{tr}[\mathbf{K}_c \mathbf{P}^{-1}]$  performance.

the constraint  $\text{tr}[\mathbf{K}_c \mathbf{P}^{-1}] = n$ , and the second stage of the algorithm is able to restore the constraint  $\text{tr}[\mathbf{K}_c \mathbf{P}^{-1}] = n$  while further reducing the  $L_2$ -sensitivity slightly.

For comparison purposes, only the iterative algorithm in (22) is applied by letting  $\lambda_i = 0$  for any  $i$  and setting  $\mathbf{P}_1 = \mathbf{I}_3$  in order to minimize the  $L_2$ -sensitivity measure in (10) (without considering the scaling constraints in (14)), and after 500 iterations it converges to

$$\mathbf{P} = \begin{bmatrix} 4.774934 & 2.835816 & 1.053819 \\ 2.835816 & 2.287705 & 1.415049 \\ 1.053819 & 1.415049 & 1.403809 \end{bmatrix}$$

which yields

$$\mathbf{T} = \begin{bmatrix} 2.185162 & 0.0 & 0.0 \\ 1.297760 & 0.776868 & 0.0 \\ 0.482261 & 1.015861 & 0.373174 \end{bmatrix}$$

and  $S_1(\mathbf{P}) = 7.832680$ . Note that this method is essentially the same as in [10]. The above coordinate transformation matrix  $\mathbf{T}$  is then scaled by an appropriate nonsingular diagonal matrix, so that the scaling constraints in (14) are satisfied. Then, the result is

$$S_1(\mathbf{P}) = 9.822372$$

where  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$  and

$$\mathbf{T} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.593\,896 & 0.562\,461 & 0.0 \\ 0.220\,698 & 0.735\,495 & 0.307\,225 \end{bmatrix}.$$

Applying the technique reported in [20] yields

$$S_1(\mathbf{P}) = 8.683\,274$$

where

$$\mathbf{T} = \begin{bmatrix} 0.913\,655 & -0.857\,313 & 0.877\,296 \\ 0.905\,773 & -0.121\,938 & 0.493\,844 \\ 0.576\,905 & 0.415\,235 & 0.377\,361 \end{bmatrix}.$$

From these results, it is observed that the proposed technique offers a smaller value of the  $L_2$ -sensitivity measure subject to the scaling constraints relative to a method (which is essentially the same as in [10]) for performing the scaling so as to satisfy the constraints in (14) after minimizing the  $L_2$ -sensitivity measure in (10) as well as the existing method in [20]. It has also turned out that the proposed iterative technique requires less than half amount of computations to attain practically the same convergence accuracy as compared with the method reported in [20].

*Example 2:* Consider a class of 2-D digital filters in (30) specified by

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 0.481\,228 & 0 & 0 \\ 0 & 0 & 0.510\,378 & 0 \\ 0 & 0 & 0 & 0.525\,287 \\ -0.031\,857 & 0.298\,663 & -0.808\,282 & 1.044\,600 \end{bmatrix}$$

$$\mathbf{A}_2 = \begin{bmatrix} -0.226\,080 & 0.776\,837 & 0.024\,693 & -0.000\,933 \\ -0.843\,550 & 1.610\,400 & -0.309\,366 & 0.065\,898 \\ -1.260\,339 & 2.005\,100 & -0.453\,220 & 0.203\,118 \\ -1.121\,498 & 1.636\,435 & -0.590\,516 & 0.562\,890 \end{bmatrix}$$

$$\mathbf{b} = [0 \ 0 \ 0 \ 0.198\,473]^T$$

$$\mathbf{c}_1 = [-0.567\,054 \ 0.231\,913 \ 0.197\,016 \ 0.239\,932]$$

$$\mathbf{c}_2 = [0.464\,344 \ 0.441\,837 \ -0.061\,100 \ 0.105\,505]$$

$$d = 0.009\,43.$$

Using (37), the Gramians  $\mathbf{K}_c$ ,  $\mathbf{W}_o$ , and  $\mathbf{M}(\mathbf{I}_4)$  are calculated as

$$\mathbf{K}_c = \begin{bmatrix} 1.000\,000 & 0.987\,279 & 0.940\,868 & 0.844\,274 \\ 0.987\,279 & 1.000\,000 & 0.976\,755 & 0.888\,478 \\ 0.940\,868 & 0.976\,755 & 1.000\,000 & 0.952\,963 \\ 0.844\,274 & 0.888\,478 & 0.952\,963 & 1.000\,000 \end{bmatrix}$$

$$\mathbf{W}_o = \begin{bmatrix} 1.337\,108 & -1.304\,050 & 0.189\,462 & -0.556\,646 \\ -1.304\,050 & 1.637\,345 & -0.429\,399 & 0.576\,183 \\ 0.189\,462 & -0.429\,399 & 2.122\,604 & -2.191\,942 \\ -0.556\,646 & 0.576\,183 & -2.191\,942 & 2.672\,474 \end{bmatrix}$$

$$\mathbf{M}(\mathbf{I}_4) = 10^3 \begin{bmatrix} 1.043\,052 & -1.097\,577 & 0.637\,299 & -0.982\,714 \\ -1.097\,577 & 1.238\,937 & -0.830\,495 & 1.153\,558 \\ 0.637\,299 & -0.830\,495 & 2.324\,302 & -2.574\,574 \\ -0.982\,714 & 1.153\,558 & -2.574\,574 & 3.019\,844 \end{bmatrix}$$

where the infinite sums in (37) were truncated with  $(i, j) = (100, 100)$ . The  $L_2$ -sensitivity measure  $S_2$  in (36) is then computed as

$$S_2 = 15\,337.965\,477.$$

Choosing  $\mathbf{P}_1 = \mathbf{I}_4$  and  $\lambda_1 = 100$  in (46) and (48) as the initial estimates, it took the proposed iterative algorithm 2000 iterations to converge to

$$\mathbf{P} = \begin{bmatrix} 1.688\,238 & 1.500\,480 & 1.311\,958 & 1.149\,324 \\ 1.500\,480 & 1.373\,665 & 1.224\,597 & 1.077\,089 \\ 1.311\,958 & 1.224\,597 & 1.148\,785 & 1.035\,997 \\ 1.149\,324 & 1.077\,089 & 1.035\,997 & 0.980\,059 \end{bmatrix}$$

which yields

$$\mathbf{T} = \begin{bmatrix} -1.174\,654 & 0.198\,729 & 0.493\,644 & 0.158\,893 \\ -1.115\,961 & 0.031\,366 & 0.322\,000 & 0.153\,714 \\ -1.034\,448 & -0.100\,879 & 0.255\,344 & -0.057\,664 \\ -0.962\,590 & -0.009\,915 & 0.107\,523 & -0.204\,501 \end{bmatrix}.$$

In this case, the Gramian  $\mathbf{M}(\mathbf{P})$  is computed from (37) as

$$\mathbf{M}(\mathbf{P}) = 10^3 \times \begin{bmatrix} 0.605\,548 & -0.599\,700 & 0.142\,673 & -0.347\,188 \\ -0.599\,700 & 0.672\,622 & -0.215\,513 & 0.360\,191 \\ 0.142\,673 & -0.215\,513 & 1.146\,333 & -1.234\,023 \\ -0.347\,188 & 0.360\,191 & -1.234\,023 & 1.478\,310 \end{bmatrix}$$

over  $(0, 0) \leq (i, j) \leq (100, 100)$ , and the  $L_2$ -sensitivity measure in (40) is minimized subject to the scaling constraints in (41) to

$$S_2(\mathbf{P}) = 372.776\,304.$$

Profiles of the  $L_2$ -sensitivity, parameter  $\lambda$ , as well as  $\text{tr}[\mathbf{K}_c\mathbf{P}^{-1}]$  during the first 2000 iterations of the algorithm are shown in Figs. 4 and 5, respectively. From these figures, a two-stage convergence behavior similar to that of Example 1 can also be observed.

For comparison purposes, only the iterative algorithm in (46) is applied by letting  $\lambda_i = 0$  for any  $i$  and setting  $\mathbf{P}_1 = \mathbf{I}_4$  in order to minimize the  $L_2$ -sensitivity measure in (40) (without considering the scaling constraints in (41)), and after 2000 iterations it converges to

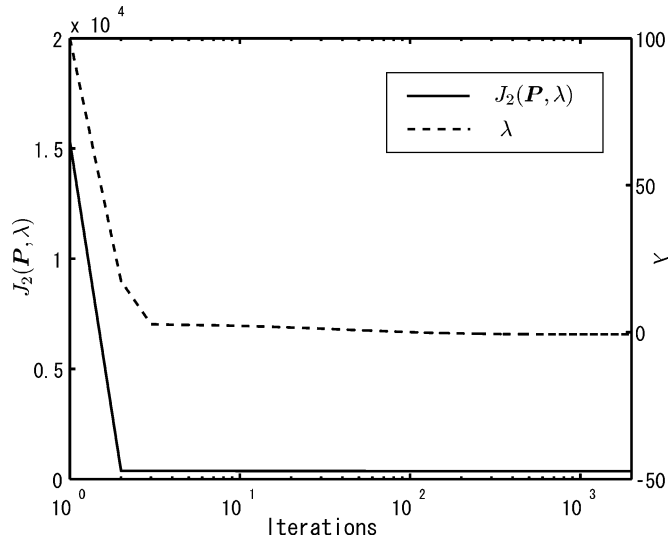
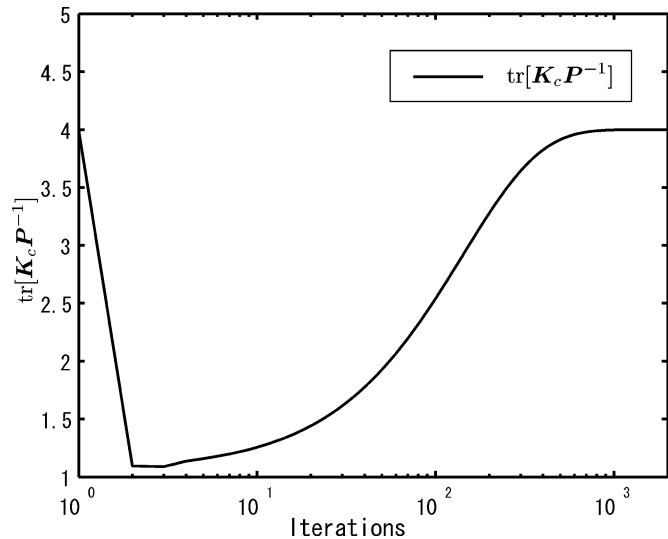
$$\mathbf{P} = \begin{bmatrix} 2.095\,327 & 1.862\,416 & 1.628\,491 & 1.426\,572 \\ 1.862\,416 & 1.705\,267 & 1.520\,239 & 1.337\,047 \\ 1.628\,491 & 1.520\,239 & 1.426\,747 & 1.286\,632 \\ 1.426\,572 & 1.337\,047 & 1.286\,632 & 1.217\,307 \end{bmatrix}$$

which yields

$$\mathbf{T} = \begin{bmatrix} 1.447\,524 & 0.0 & 0.0 & 0.0 \\ 1.286\,621 & 0.223\,322 & 0.0 & 0.0 \\ 1.125\,018 & 0.325\,839 & 0.234\,330 & 0.0 \\ 0.985\,525 & 0.309\,192 & 0.329\,243 & -0.205\,050 \end{bmatrix}$$

and  $S_2(\mathbf{P}) = 372.464\,800$ . The above coordinate transformation matrix  $\mathbf{T}$  is then scaled by an appropriate nonsingular diagonal matrix, so that the scaling constraints in (41) are satisfied. Then, the result is

$$S_2(\mathbf{P}) = 458.063\,668$$

Fig. 4.  $L_2$ -sensitivity and  $\lambda$  performances.Fig. 5.  $\text{tr}[\mathbf{K}_c \mathbf{P}^{-1}]$  performance.

where  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$  and

$$\mathbf{T} = \begin{bmatrix} 1.0 & 0.0 & 0.0 & 0.0 \\ 0.888\ 843 & 0.187\ 005 & 0.0 & 0.0 \\ 0.777\ 201 & 0.272\ 850 & 0.171\ 289 & 0.0 \\ 0.680\ 835 & 0.258\ 910 & 0.240\ 668 & -0.250\ 773 \end{bmatrix}.$$

From the above result, it is observed that the constrained optimization technique proposed here offers smaller  $L_2$ -sensitivity subject to the scaling constraints relative to a method for carrying out the scaling so as to satisfy the constraints in (41) after minimizing the  $L_2$ -sensitivity measure in (40).

## V. CONCLUSION

The problem of minimizing an  $L_2$ -sensitivity measure subject to  $L_2$ -norm dynamic-range scaling constraints for state-space digital filters as well as a class of 2-D state-space digital filters have been investigated. In each case, a novel iterative algorithm has been developed to solve the constraint optimization problem directly. This has been performed by using a Lagrange function and some matrix-theoretic techniques. The optimal state-space

filter structure has been constructed by applying the coordinate-transformation matrix adjusted to satisfy the scaling constraints without altering the  $L_2$ -sensitivity. In addition, it has been shown that the 2-D LSS model used to describe a class of 2-D state-space digital filters corresponds to a transposed structure of the Fornasini–Marchesini second LSS model. Our computer simulation results have demonstrated the effectiveness of the proposed techniques compared with the existing methods.

## REFERENCES

- [1] L. Thiele, "Design of sensitivity and round-off noise optimal state-space discrete systems," *Int. J. Circuit Theory Appl.*, vol. 12, pp. 39–46, Jan. 1984.
- [2] V. Tavsanoğlu and L. Thiele, "Optimal design of state-space digital filters by simultaneous minimization of sensitivity and roundoff noise," *IEEE Trans. Circuits Syst.*, vol. CAS-31, no. 10, pp. 884–888, Oct. 1984.
- [3] L. Thiele, "On the sensitivity of linear state-space systems," *IEEE Trans. Circuits Syst.*, vol. CAS-33, no. 5, pp. 502–510, May 1986.
- [4] M. Iwatsuki, M. Kawamata, and T. Higuchi, "Statistical sensitivity and minimum sensitivity structures with fewer coefficients in discrete time linear systems," *IEEE Trans. Circuits Syst.*, vol. 37, no. 1, pp. 72–80, Jan. 1989.
- [5] G. Li, B. D. O. Anderson, M. Gevers, and J. E. Perkins, "Optimal FWL design of state-space digital systems with weighted sensitivity minimization and sparseness consideration," *IEEE Trans. Circuits Syst. I*, vol. 39, no. 5, pp. 365–377, May 1992.
- [6] W.-Y. Yan and J. B. Moore, "On  $L_2$ -sensitivity minimization of linear state-space systems," *IEEE Trans. Circuits Syst. I*, vol. 39, no. 8, pp. 641–648, Aug. 1992.
- [7] G. Li and M. Gevers, "Optimal synthetic FWL design of state-space digital filters," in *Proc. 1992 IEEE Int. Conf. Acoustics, Speech, Signal Processing (ICASSP)*, vol. 4, pp. 429–432.
- [8] M. Gevers and G. Li, *Parameterizations in Control, Estimation and Filtering Problems: Accuracy Aspects*. New York: Springer-Verlag, 1993.
- [9] C. Xiao, "Improved  $L_2$ -sensitivity for state-space digital system," *IEEE Trans. Signal Process.*, vol. 45, no. 4, pp. 837–840, Apr. 1997.
- [10] T. Hinamoto, S. Yokoyama, T. Inoue, W. Zeng, and W.-S. Lu, "Analysis and minimization of  $L_2$ -sensitivity for linear systems and two-dimensional state-space filters using general controllability and observability Gramians," *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. 49, no. 9, pp. 1279–1289, Sep. 2002.
- [11] M. Kawamata, T. Lin, and T. Higuchi, "Minimization of sensitivity of 2-D state-space digital filters and its relation to 2-D balanced realizations," in *Proc. 1987 IEEE Int. Symp. Circuits Systems*, pp. 710–713.
- [12] T. Hinamoto, T. Hamanaka, and S. Maekawa, "Synthesis of 2-D state-space digital filters with low sensitivity based on the Fornasini–Marchesini model," *IEEE Trans. Acoust. Speech, Signal Process.*, vol. ASSP-38, no. 9, pp. 1587–1594, Sep. 1990.
- [13] T. Hinamoto, T. Takao, and M. Muneyasu, "Synthesis of 2-D separable-denominator digital filters with low sensitivity," *J. Franklin Inst.*, vol. 329, pp. 1063–1080, 1992.
- [14] T. Hinamoto and T. Takao, "Synthesis of 2-D state-space filter structures with low frequency-weighted sensitivity," *IEEE Trans. Circuits Syst. II, Analog Digit. Signal Process.*, vol. 39, no. 9, pp. 646–651, Sep. 1992.
- [15] —, "Minimization of frequency-weighting sensitivity in 2-D systems based on the Fornasini–Marchesini second model," in *Proc. 1992 IEEE Int. Conf. Acoustics, Speech, Signal Processing (ICASSP)*, pp. 401–404.
- [16] G. Li, "On frequency weighted minimal  $L_2$  sensitivity of 2-D systems using Fornasini–Marchesini LSS model," *IEEE Trans. Circuits Syst. I, Fundam. Theory Appl.*, vol. 44, no. 7, pp. 642–646, Jul. 1997.
- [17] —, "Two-dimensional system optimal realizations with  $L_2$ -sensitivity minimization," *IEEE Trans. Signal Process.*, vol. 46, no. 3, pp. 809–813, Mar. 1998.
- [18] T. Hinamoto, Y. Zempo, Y. Nishino, and W.-S. Lu, "An analytical approach for the synthesis of two-dimensional state-space filter structures with minimum weighted sensitivity," *IEEE Trans. Circuits Syst. I, Fundam. Theory Appl.*, vol. 46, no. 10, pp. 1172–1183, Oct. 1999.
- [19] T. Hinamoto and Y. Sugie, " $L_2$ -sensitivity analysis and minimization of 2-D separable-denominator state-space digital filters," *IEEE Trans. Signal Process.*, vol. 50, no. 12, pp. 3107–3114, Dec. 2002.
- [20] T. Hinamoto, H. Ohnishi, and W.-S. Lu, "Minimization of  $L_2$ -sensitivity for state-space digital filters subject to  $L_2$ -scaling constraints," in *Proc. 2004 IEEE Int. Symp. Circuits Systems*, vol. III, pp. 137–140.



- [21] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed-point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, no. 9, pp. 551–562, Sep. 1976.
- [22] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust. Speech, Signal Process.*, vol. 25, no. 4, pp. 273–281, Aug. 1977.
- [23] T. Kailath, *Linear System*. Englewood Cliffs, NJ: Prentice-Hall, 1980.
- [24] L. L. Scharf, *Statistical Signal Processing*. Reading, MA: Addison-Wesley, 1991.
- [25] T. Hinamoto, "A novel local state-space model for 2-D digital filters and its properties," in *Proc. 2001 IEEE Int. Symp. Circuits Systems*, vol. 2, pp. 545–548.
- [26] E. Fornasini and G. Marchesini, "Doubly-indexed dynamical systems: State-space models and structural properties," *Math Syst. Theory*, vol. 12, pp. 59–72, 1978.
- [27] R. P. Roesser, "A discrete state-space model for linear image processing," *IEEE Trans. Autom. Control*, vol. AC-20, no. 1, pp. 1–10, Feb. 1975.

**Takao Hinamoto** (M'77–SM'84–F'01) received the B.E. degree from Okayama University, Okayama, Japan, in 1969, the M.E. degree from Kobe University, Kobe, Japan, in 1971, and the Dr.Eng. degree from Osaka University, Osaka, Japan, in 1977, all in electrical engineering.

From 1972 to 1988, he was with the Faculty of Engineering, Kobe University. From 1979 to 1981, he was a Visiting Member of Staff in the Department of Electrical Engineering, Queen's University, Kingston, ON, Canada, on leave from Kobe University. From 1988 to 1991, he was a Professor of electronic circuits in the Faculty of Engineering, Tottori University, Tottori, Japan. Since January 1992, he has been a Professor of electronic control in the Department of Electrical Engineering, Hiroshima University, Hiroshima, Japan. His research interests include digital signal processing, system theory, and control engineering. He has published about 350 papers in these areas and is the coeditor and coauthor of the book *Two-Dimensional Signal and Image Processing* (Tokyo, Japan: SICE, 1996).

Dr. Hinamoto served as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS—II: ANALOG AND DIGITAL SIGNAL PROCESSING from 1993 to 1995 and as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS—I: FUNDAMENTAL THEORY AND APPLICATIONS from 2002 to 2003. He was the Guest Editor of the special sections on Digital Signal Processing (DSP) in the August 1998 issue as well as Adaptive Signal Processing and Its Applications in the March 2005 issue of the *IEICE Transactions on Fundamentals*. He served as the Chair of the 12th DSP Symposium held in Hiroshima, Japan, in November 1997, sponsored by the DSP Technical Committee of IEICE. He also served as the General Chair of the 47th IEEE International Midwest Symposium on Circuits and Systems held in Hiroshima, Japan, in July 2004. Since 1995, he has been a member of the Steering Committee of the IEEE International Midwest Symposium on Circuits and Systems, and since 1998, a member of the Digital Signal Processing Technical Committee in the IEEE Circuits and Systems Society. He was a member of the Technical Program Committee for the International Symposium of Circuits and Systems (ISCAS) 1999 and an International Coordinator of the Organizing Committee for ISCAS 2004. From 1993 to 2000, he served as Senator or Member of the Board of Directors in the Society of Instrument and Control Engineers (SICE), and from 1999 to 2001, he was the Chair of the Chugoku Chapter of SICE. He played a leading role in establishing the Hiroshima Section of IEEE and served as the Interim Chair of the Section. From June 2003 to May 2004, he served as the Chair of the DSP Technical Committee of IEICE and the Chair of the Chugoku Chapter of IEICE. He is a recipient of the IEEE Third Millennium Medal and a Fellow of IEICE.

**Ken-ichi Iwata** received the B.E. and M.E. degrees in electrical engineering from Hiroshima University, Hiroshima, Japan, in 2003 and 2005, respectively.

He was engaged in research on digital signal processing during his graduate studies. Since April 2005, he has been with Mitsubishi Electric Corporation, Himeji, Japan.

**Wu-Sheng Lu** (S'81–M'85–SM'90–F'99) received the B.Sc. degree in mathematics from Fudan University, Shanghai, China, in 1964 and the M.S. degree in electrical engineering and the Ph.D. degree in control science from the University of Minnesota, MN, in 1983 and 1984, respectively.

He was a Postdoctoral Fellow at the University of Victoria, Victoria, BC, Canada, in 1985 and a Visiting Assistant Professor with the University of Minnesota in 1986. Since 1987, he has been with the University of Victoria, where he is currently a Professor. His current teaching and research interests are in the general areas of digital signal processing and application of optimization methods. He is the coauthor with A. Antoniou of the book *Two-Dimensional Digital Filters* (New York: Marcel Dekker, 1992).

Dr. Lu served as an Associate Editor of the *Canadian Journal of Electrical and Computer Engineering* in 1989 and Editor of the same journal from 1990 to 1992. He served as an Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS, PART II from 1993 to 1995 and for Part I of the same journal from 1999 to 2001. Presently, he is serving as Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS, PART I, and the *International Journal of Multidimensional Systems and Signal Processing*.