

Minimization of Frequency-Weighted l_2 -Sensitivity Subject to l_2 -Scaling Constraints for Two-Dimensional State-Space Digital Filters

Takao Hinamoto, Toru Oumi, Osemekhian I. Omoifo and Wu-Sheng Lu

Abstract—This paper investigates the problem of frequency-weighted l_2 -sensitivity minimization subject to l_2 -scaling constraints for two-dimensional (2-D) state-space digital filters described by the Roesser model. It is shown that the Fornasini-Marchesini second model can be imbedded in the Roesser model. Two iterative methods are developed to solve the constrained optimization problem encountered. The first iterative method introduces a Lagrange function and optimizes it using some matrix-theoretic techniques and an efficient bisection method. The second iterative method converts the problem into an unconstrained optimization formulation by using linear-algebraic techniques and solves it by applying an efficient quasi-Newton algorithm. The optimal filter structure with minimum frequency-weighted l_2 -sensitivity and no overflow is then synthesized by an appropriate coordinate transformation. Case studies are presented to demonstrate the validity and effectiveness of the proposed techniques.

Index Terms—2-D digital filters, Roesser's model, Fornasini-Marchesini's second model, frequency-weighted l_2 -sensitivity minimization, l_2 -scaling constraints, no overflow, Lagrange function, bisection method, quasi-Newton method

I. INTRODUCTION

It is well known that there exist infinitely many minimal state-space realizations for a given transfer function, and some inherent properties such as controllability, observability, stability, etc. are invariant within these realizations. However, performance measures such as coefficient sensitivity, output roundoff noise, overflow oscillations, etc. may be significantly varying among the realizations. Consider a transfer function with coefficients of infinite accuracy, which meets certain design specifications including stability. When the transfer function is implemented by a state-space model with a finite binary representation, truncation or rounding of the state-space model is required to satisfy the finite word length (FWL) constraints. As a result, the characteristics of the stable filter might be so altered that the filter may become unstable. This motivates the study of the coefficient sensitivity minimization problem. To date, several techniques have been reported for synthesizing the state-space descriptions with minimum coefficient sensitivity. The techniques can be divided into two

main classes: those for l_1/l_2 -mixed sensitivity minimization [1]-[5] and those for l_2 -sensitivity minimization [6]-[11]. In [6]-[10], it has been argued that the sensitivity measure based on a sole l_2 -norm 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 [12]-[17] and l_2 -sensitivity minimization problem [10],[17]-[20] have also been investigated. It has been realized that solutions for *frequency-weighted* sensitivity minimization would be of practical use as these solutions allow to emphasize or de-emphasize the filter's sensitivity in certain frequency regions of interest. Synthesis procedures of the optimal FWL 2-D filter structures that minimize the frequency-weighted sensitivity measure have been considered [15]-[18]. However, the minimization methods proposed in the above work do not impose constraints on the scaling of the design variables. As a result, elimination of overflow cannot be ensured. More recently, the minimization problem of l_2 -sensitivity subject to l_2 -scaling constraints has been explored for 1-D and a class of 2-D state-space digital filters [21]-[23]. It is well known that the use of scaling constraints can be beneficial for suppressing overflow [24],[25]. However, frequency-weighted sensitivity measure has not yet been considered in [21]-[23].

In this paper, we investigate the problem of minimizing a frequency-weighted l_2 -sensitivity measure subject to l_2 -scaling constraints for 2-D state-space digital filters described by the Roesser local state-space (LSS) model [26]. We then proceed by introducing an expression for evaluating the frequency-weighted l_2 -sensitivity, and formulating the minimization problem for the frequency-weighted l_2 -sensitivity measure subject to l_2 -scaling constraints. Next, two iterative methods are developed for solving the constrained optimization problem. The first iterative method introduces a Lagrange function, and makes use of some matrix-theoretic techniques and an efficient bisection method. The second iterative method relies on a technique that converts the constrained optimization problem into an unconstrained optimization formulation and utilizes an efficient quasi-Newton method with closed-form formula for gradient evaluation. Finally, case studies are presented to demonstrate the validity and effectiveness of the proposed techniques.

One of the contributions made in this paper is to show that **either the Fornasini-Marchesini (FM) second LSS model [27] or its transposed-structure model [22],[28] can be imbedded in the Roesser model as a special case.** This justifies the use of the Roesser model in our studies. Another

Manuscript received mm dd, 2007; revised mm dd, 2007.

T. Hinamoto, T. Oumi and O. I. Omoifo are with the Graduate School of Engineering, Hiroshima University, Higashi-Hiroshima 739-8527, Japan. (e-mail: {hinamoto,oumi,osei}@hiroshima-u.ac.jp, Phone:+81-82-424-7672, Fax:+81-82-422-7195)

W.-S. Lu is with the Department of Electrical and Computer Engineering, University of Victoria, Victoria, B.C, Canada, V8W 3P6. (e-mail: wslu@ece.uvic.ca, Phone:+1-250-721-8692, Fax:+1-250-721-6052)

contribution is that a bisection method is applied to obtain the Lagrange multipliers, which makes it possible to attain considerably faster convergence than the algorithm reported in [22]. Moreover, unlike [21] and [22], the present paper investigates a *frequency-weighted* l_2 -sensitivity measure under l_2 -scaling constraints. Although this extension is technically manageable, to the best of our knowledge, this is the first time a *frequency-weighted* l_2 -sensitivity measure under l_2 -scaling constraints is addressed in the l_2/l_2 framework for state-space digital filters.

Throughout the paper, \mathbf{I}_n stands for the identity matrix of dimension $n \times n$, \oplus is used to denote the direct sum of matrices, the transpose (conjugate transpose) of a matrix \mathbf{A} is indicated by \mathbf{A}^T (\mathbf{A}^*), and the trace and i th diagonal element of a square matrix \mathbf{A} are denoted by $\text{tr}[\mathbf{A}]$ and $(\mathbf{A})_{ii}$, respectively.

II. PROBLEM FORMULATION

2.1 System Models

Consider a stable, separately locally controllable and separately locally observable LSS model for 2-D recursive digital filters

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^h(i+1, j) \\ \mathbf{x}^v(i, j+1) \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix} \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} u(i, j) \\ y(i, j) &= [\mathbf{c}_1 \quad \mathbf{c}_2] \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + d u(i, j) \end{aligned} \quad (1)$$

which was originally proposed by Roesser [26],[29], where $\mathbf{x}^h(i, j)$ is an $m \times 1$ horizontal state vector, $\mathbf{x}^v(i, j)$ is an $n \times 1$ vertical state vector, $u(i, j)$ is a scalar input, $y(i, j)$ is a scalar output, and $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3, \mathbf{A}_4, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}_1, \mathbf{c}_2$, and d are $m \times m, m \times n, n \times m, n \times n, m \times 1, n \times 1, 1 \times m, 1 \times n$, and 1×1 real constant matrices, respectively. A block diagram of the LSS model in (39) is shown in Fig. 1. The

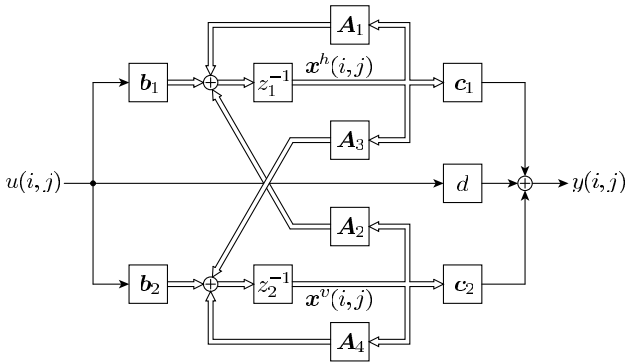


Fig. 1. The block diagram of the Roesser LSS model.

transfer function of the LSS model in (39) is given by

$$H(z_1, z_2) = \mathbf{c}(\mathbf{Z} - \mathbf{A})^{-1} \mathbf{b} + d \quad (2)$$

where \mathbf{A}, \mathbf{b} , and \mathbf{c} are $(m+n) \times (m+n), (m+n) \times 1$, and $1 \times (m+n)$ real constant matrices defined by

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, \quad \mathbf{c} = [\mathbf{c}_1 \quad \mathbf{c}_2],$$

respectively, and $\mathbf{Z} = z_1 \mathbf{I}_m \oplus z_2 \mathbf{I}_n$. For the sake of simplicity, the LSS model in (39) is denoted hereafter by $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$.

Alternatively, an LSS model for a class of 2-D recursive digital filters can be described by [22],[28]

$$\begin{aligned} \begin{bmatrix} \mathbf{x}(i+1, j+1) \\ y(i, j) \end{bmatrix} &= \begin{bmatrix} \mathbf{A}'_1 & \mathbf{A}'_2 \\ \mathbf{c}'_1 & \mathbf{c}'_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}(i, j+1) \\ \mathbf{x}(i+1, j) \end{bmatrix} \\ &+ \begin{bmatrix} \mathbf{b}' \\ d \end{bmatrix} u(i, j) \end{aligned} \quad (3)$$

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 $N \times N, N \times N, N \times 1, 1 \times N, 1 \times N$, and 1×1 real constant matrices, respectively. The transfer function of the LSS model in (3) is given by

$$\begin{aligned} D(z_1, z_2) &= (z_1^{-1} \mathbf{c}'_1 + z_2^{-1} \mathbf{c}'_2) \\ &\cdot (\mathbf{I}_n - z_1^{-1} \mathbf{A}'_1 - z_2^{-1} \mathbf{A}'_2)^{-1} \mathbf{b}' + d. \end{aligned} \quad (4)$$

If we define

$$\mathbf{x}(i, j+1) = \mathbf{x}^h(i, j), \quad \mathbf{x}(i+1, j) = \mathbf{x}^v(i, j) \quad (5)$$

then the LSS model in (3) can then be imbedded in that of (39) as a special case as follows:

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^h(i+1, j) \\ \mathbf{x}^v(i, j+1) \end{bmatrix} &= \begin{bmatrix} \mathbf{A}'_1 & \mathbf{A}'_2 \\ \mathbf{A}'_1 & \mathbf{A}'_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + \begin{bmatrix} \mathbf{b}' \\ \mathbf{b}' \end{bmatrix} u(i, j) \\ y(i, j) &= [\mathbf{c}'_1 \quad \mathbf{c}'_2] \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + d u(i, j) \end{aligned} \quad (6)$$

where $m = n = N$. It is noted that $D(z_1, z_2)^T$ can be viewed as a transfer function of the FM second LSS model [27], which reveals that the LSS model of $D(z_1, z_2)^T$ can be realized by a transposed structure of that in (6). Therefore, the LSS model in (39) is more general than either the LSS model in (3) or the FM second LSS model [27] (**and vice versa with the same dimension, but increased number of coefficients [27]**). In addition, we note that the technique reported in [22] has merely treated the l_2 -sensitivity minimization problem for the LSS model in (3) subject to l_2 -scaling constraints. Recall that the total numbers of the coefficients in (39) and (3) are $(m+n)^2 + 2(m+n) + 1$ and $2N^2 + 3N + 1$, respectively. This means that the LSS model in (39) has less number of the coefficients than that of (3) when their local state vectors possess the same dimension, i.e., $m+n = N$. Under these circumstances, it is worthwhile to consider the more general problem of minimizing the *frequency-weighted* l_2 -sensitivity for the LSS model in (39) subject to l_2 -scaling constraints, because its solutions will make it possible to emphasize or de-emphasize the filter's sensitivity in certain frequency regions of interest, and because the LSS model in (39) owns less number of the coefficients in case $m+n = N$.

2.2 A Frequency-Weighted l_2 -Sensitivity Measure

Suppose that the LSS model in (39) is implemented by FWL fixed-point arithmetic with a B -bit fractional representation, and is realized with coefficient matrices

$$\begin{aligned}\tilde{\mathbf{A}} &= \mathbf{A} + \Delta\mathbf{A}, & \tilde{\mathbf{b}} &= \mathbf{b} + \Delta\mathbf{b} \\ \tilde{\mathbf{c}} &= \mathbf{c} + \Delta\mathbf{c}, & \tilde{d} &= d + \Delta d\end{aligned}\quad (7)$$

where $\Delta\mathbf{A}$, $\Delta\mathbf{b}$, $\Delta\mathbf{c}$, and Δd stand for the quantization errors of the coefficient matrices. The transfer function of the FWL realization is then expressed as

$$\tilde{H}(z_1, z_2) = \tilde{\mathbf{c}}(\mathbf{Z} - \tilde{\mathbf{A}})^{-1}\tilde{\mathbf{b}} + \tilde{d}. \quad (8)$$

Let $\{p_i, i = 1, 2, \dots, M\}$ be the set of the ideal parameters of a realization and let $\{\tilde{p}_i, i = 1, 2, \dots, M\}$ be its FWL version where $\tilde{p}_i = p_i + \Delta p_i$ with Δp_i indicating the corresponding parameter perturbation. If all Δp_i are sufficiently small in magnitude, then the first-order approximation of the Taylor series expansion yields

$$\begin{aligned}\Delta H(z_1, z_2) &= \tilde{H}(z_1, z_2) - H(z_1, z_2) \\ &\simeq \sum_{i=1}^M \frac{\partial H(z_1, z_2)}{\partial p_i} \Delta p_i.\end{aligned}\quad (9)$$

Obviously, smaller $\partial H(z_1, z_2)/\partial p_i$ for $i = 1, 2, \dots, M$ yield smaller transfer function error $\Delta H(z_1, z_2)$. For a fixed-point implementation of B bits, the parameter perturbations can be considered to be independent random-variables uniformly distributed within the range $[-2^{-B-1}, 2^{-B-1}]$. Under these circumstances, a measure of the transfer function error can statistically be defined as

$$\sigma_{\Delta H}^2 = \frac{1}{(2\pi j)^2} \oint_{|z_1|=1} \oint_{|z_2|=1} E[|\Delta H(z_1, z_2)|^2] \frac{dz_1 dz_2}{z_1 z_2} \quad (10)$$

where $E(\cdot)$ denotes the ensemble-average operation. Since Δp_i 's are uniformly-distributed independent random variables, it follows that

$$E[|\Delta H(z_1, z_2)|^2] = \sum_{i=1}^M \left| \frac{\partial H(z_1, z_2)}{\partial p_i} \right|^2 \sigma^2 \quad (11)$$

where

$$\sigma^2 = E[(\Delta p_i)^2] = \frac{1}{12} 2^{-2B}.$$

Thus, if the measure

$$S_o = \frac{1}{(2\pi j)^2} \oint_{|z_1|=1} \oint_{|z_2|=1} \sum_{i=1}^M \left| \frac{\partial H(z_1, z_2)}{\partial p_i} \right|^2 \frac{dz_1 dz_2}{z_1 z_2} \quad (12)$$

is minimized then the minimum variance $\sigma_{\Delta H}^2$ can be attained because of the relation $\sigma_{\Delta H}^2 = S_o \sigma^2$. The measure in (12) is referred to as an l_2 -sensitivity measure.

The frequency-weighted l_2 -sensitivity of the LSS model in (39) 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} , that is differentiable with respect to all the entries of \mathbf{X} . The sensitivity function of $f(\mathbf{X})$ with respect to \mathbf{X} is then defined as [5]

$$\mathbf{S}_{\mathbf{X}} = \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}, \quad (\mathbf{S}_{\mathbf{X}})_{ij} = \frac{\partial f(\mathbf{X})}{\partial x_{ij}} \quad (13)$$

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

Definition 2: In order to take into account the sensitivity of the transfer function in a specified frequency band, or even at some discrete frequency points, the weighted sensitivity functions are defined as [5],[17]

$$\begin{aligned}\frac{\delta H(z_1, z_2)}{\delta \mathbf{A}} &= W_A(z_1, z_2) \frac{\partial H(z_1, z_2)}{\partial \mathbf{A}} \\ \frac{\delta H(z_1, z_2)}{\delta \mathbf{b}} &= W_B(z_1, z_2) \frac{\partial H(z_1, z_2)}{\partial \mathbf{b}} \\ \frac{\delta H(z_1, z_2)}{\delta \mathbf{c}^T} &= W_C(z_1, z_2) \frac{\partial H(z_1, z_2)}{\partial \mathbf{c}^T}\end{aligned}\quad (14)$$

where $W_A(z_1, z_2)$, $W_B(z_1, z_2)$, and $W_C(z_1, z_2)$ are scalar, stable, causal functions of the complex variables z_1 and z_2 .

Notice that δ in (14) is not meant to be a derivative operator, but rather a notation for defining the weighted parameter sensitivity.

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 by

$$\begin{aligned}\|\mathbf{X}(z_1, z_2)\|_2 &= \left(\text{tr} \left[\frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{X}(z_1, z_2) \mathbf{X}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \right] \right)^{\frac{1}{2}}\end{aligned}\quad (15)$$

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

From (2) and Definitions 1-3, the overall frequency-weighted l_2 -sensitivity measure for the LSS model in (39) can be defined as

$$\begin{aligned}S &= \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{A}} \right\|_2^2 + \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{b}} \right\|_2^2 + \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{c}^T} \right\|_2^2 \\ &= \left\| W_A(z_1, z_2) [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T \right\|_2^2 \\ &\quad + \left\| W_B(z_1, z_2) \mathbf{G}^T(z_1, z_2) \right\|_2^2 + \left\| W_C(z_1, z_2) \mathbf{F}(z_1, z_2) \right\|_2^2\end{aligned}\quad (16)$$

where

$$\mathbf{F}(z_1, z_2) = (\mathbf{Z} - \mathbf{A})^{-1} \mathbf{b}, \quad \mathbf{G}(z_1, z_2) = \mathbf{c}(\mathbf{Z} - \mathbf{A})^{-1}.$$

It follows that the frequency-weighted l_2 -sensitivity measure in (16) can be written as

$$S = \text{tr}[\mathbf{M}_A] + \text{tr}[\mathbf{W}_B] + \text{tr}[\mathbf{K}_C] \quad (17)$$

where \mathbf{M}_A , \mathbf{W}_B , and \mathbf{K}_C are obtained by the following general expression:

$$\mathbf{X} = \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{Y}(z_1, z_2) \mathbf{Y}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2}$$

with $\mathbf{Y}(z_1, z_2) = W_A(z_1, z_2) [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T$ for $\mathbf{X} = \mathbf{M}_A$, $\mathbf{Y}(z_1, z_2) = W_B(z_1, z_2) \mathbf{G}^T(z_1, z_2)$ for $\mathbf{X} = \mathbf{W}_B$, and $\mathbf{Y}(z_1, z_2) = W_C(z_1, z_2) \mathbf{F}(z_1, z_2)$ for $\mathbf{X} = \mathbf{K}_C$.

2.3 Problem Formulation

Define a state-space coordinate transformation by [26],[29]

$$\begin{bmatrix} \bar{\mathbf{x}}^h(i, j) \\ \bar{\mathbf{x}}^v(i, j) \end{bmatrix} = \begin{bmatrix} \mathbf{T}_1^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_4^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} \quad (18)$$

where \mathbf{T}_1 and \mathbf{T}_4 are $m \times m$ and $n \times n$ nonsingular matrices, respectively. New realizations can then be characterized as $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m,n}$ with

$$\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T} \quad (19)$$

where $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$. For a new realization, the frequency-weighted l_2 -sensitivity measure in (17) is changed to

$$\begin{aligned} S(\mathbf{P}) &= \text{tr}[\mathbf{M}_A(\mathbf{P})\mathbf{P}] + \text{tr}[\mathbf{W}_B\mathbf{P}] + \text{tr}[\mathbf{K}_C\mathbf{P}^{-1}] \\ &= \text{tr}[\mathbf{N}_A(\mathbf{P})\mathbf{P}^{-1}] + \text{tr}[\mathbf{W}_B\mathbf{P}] + \text{tr}[\mathbf{K}_C\mathbf{P}^{-1}] \end{aligned} \quad (20)$$

with $\mathbf{P} = \mathbf{T}\mathbf{T}^T = \mathbf{P}_1 \oplus \mathbf{P}_4$ and

$$\begin{aligned} \mathbf{M}_A(\mathbf{P}) &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{Y}(z_1, z_2) \mathbf{P}^{-1} \mathbf{Y}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \\ \mathbf{N}_A(\mathbf{P}) &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{Y}^*(z_1, z_2) \mathbf{P} \mathbf{Y}(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \end{aligned}$$

where $\mathbf{Y}(z_1, z_2) = \mathbf{W}_A(z_1, z_2) [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T$.

If l_2 -scaling constraints are imposed on the horizontal and vertical state vectors $\bar{\mathbf{x}}^h(i, j)$ and $\bar{\mathbf{x}}^v(i, j)$, we require that [30]

$$\begin{aligned} (\bar{\mathbf{K}}_1)_{\xi\xi} &= (\mathbf{T}_1^{-1} \mathbf{K}_1 \mathbf{T}_1^{-T})_{\xi\xi} = 1 \quad \text{for } \xi = 1, 2, \dots, m \\ (\bar{\mathbf{K}}_4)_{\zeta\zeta} &= (\mathbf{T}_4^{-1} \mathbf{K}_4 \mathbf{T}_4^{-T})_{\zeta\zeta} = 1 \quad \text{for } \zeta = 1, 2, \dots, n \end{aligned} \quad (21)$$

where

$$\begin{aligned} \mathbf{K} &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{F}(z_1, z_2) \mathbf{F}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \\ &= \begin{bmatrix} \mathbf{K}_1 & \mathbf{K}_2 \\ \mathbf{K}_3 & \mathbf{K}_4 \end{bmatrix} \end{aligned}$$

is the local controllability Gramian for the LSS model in (39) with an $m \times m$ submatrix \mathbf{K}_1 and an $n \times n$ submatrix \mathbf{K}_4 along its diagonal [29].

Thus, the l_2 -scaling constrained frequency-weighted l_2 -sensitivity minimization problem can be formulated as follows: *Given matrices \mathbf{A} , \mathbf{b} , and \mathbf{c} , obtain a block-diagonal nonsingular matrix $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ which minimizes $S(\mathbf{P})$ in (20) subject to l_2 -scaling constraints in (21).*

III. PROBLEM SOLUTION

3.1 A Constrained Optimization Method

Solving the optimization problem formulated above consists of several steps. First, we relax the problem of minimizing $S(\mathbf{P})$ in (20) subject to l_2 -scaling constraints in (21) into the problem

$$\begin{aligned} &\text{minimize } S(\mathbf{P}) \text{ in (20)} \\ &\text{subject to } \text{tr}[\mathbf{K}_1 \mathbf{P}_1^{-1}] = m \text{ and } \text{tr}[\mathbf{K}_4 \mathbf{P}_4^{-1}] = n. \end{aligned} \quad (22)$$

If $\text{tr}[\mathbf{K}_1 \mathbf{P}_1^{-1}] = m$ ($\text{tr}[\mathbf{K}_4 \mathbf{P}_4^{-1}] = n$) is satisfied, then an $m \times m$ ($n \times n$) orthogonal matrix \mathbf{U}_1 (\mathbf{U}_4) matrix can always be constructed so that $\mathbf{T}_1 = \mathbf{P}_1^{1/2} \mathbf{U}_1$ ($\mathbf{T}_4 = \mathbf{P}_4^{1/2} \mathbf{U}_4$) satisfies l_2 -scaling constraints in (21) [22]. This justifies the relaxation made in (22).

In order to solve problem (22) for $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$, we define the Lagrange function of the problem as

$$\begin{aligned} J(\mathbf{P}, \lambda_1, \lambda_4) &= \text{tr}[\mathbf{M}_1(\mathbf{P})\mathbf{P}_1] + \text{tr}[\mathbf{M}_4(\mathbf{P})\mathbf{P}_4] + \text{tr}[\mathbf{W}_{1B}\mathbf{P}_1] \\ &\quad + \text{tr}[\mathbf{W}_{4B}\mathbf{P}_4] + \text{tr}[\mathbf{K}_{1C}\mathbf{P}_1^{-1}] + \text{tr}[\mathbf{K}_{4C}\mathbf{P}_4^{-1}] \\ &\quad + \lambda_1(\text{tr}[\mathbf{K}_1\mathbf{P}_1^{-1}] - m) + \lambda_4(\text{tr}[\mathbf{K}_4\mathbf{P}_4^{-1}] - n) \end{aligned} \quad (23)$$

where λ_1 and λ_4 are the Lagrange multipliers, and

$$\mathbf{M}_A(\mathbf{P}) = \begin{bmatrix} \mathbf{M}_1(\mathbf{P}) & \mathbf{M}_2(\mathbf{P}) \\ \mathbf{M}_3(\mathbf{P}) & \mathbf{M}_4(\mathbf{P}) \end{bmatrix}$$

$$\mathbf{W}_B = \begin{bmatrix} \mathbf{W}_{1B} & \mathbf{W}_{2B} \\ \mathbf{W}_{3B} & \mathbf{W}_{4B} \end{bmatrix}, \quad \mathbf{K}_C = \begin{bmatrix} \mathbf{K}_{1C} & \mathbf{K}_{2C} \\ \mathbf{K}_{3C} & \mathbf{K}_{4C} \end{bmatrix}$$

with an $m \times m$ submatrix and an $n \times n$ submatrix along its diagonal for each matrix. Using the formula for evaluating matrix gradient [31, p.275]

$$\frac{\partial [\text{tr}(\mathbf{M}\mathbf{X})]}{\partial \mathbf{X}} = \mathbf{M}^T \quad (24)$$

$$\frac{\partial [\text{tr}(\mathbf{M}\mathbf{X}^{-1})]}{\partial \mathbf{X}} = -(\mathbf{X}^{-1} \mathbf{M} \mathbf{X}^{-1})^T$$

we compute

$$\begin{aligned} \frac{\partial J(\mathbf{P}, \lambda_1, \lambda_4)}{\partial \mathbf{P}_1} &= \mathbf{M}_1(\mathbf{P}) - \mathbf{P}_1^{-1} \mathbf{N}_1(\mathbf{P}) \mathbf{P}_1^{-1} + \mathbf{W}_{1B} \\ &\quad - \mathbf{P}_1^{-1} \mathbf{K}_{1C} \mathbf{P}_1^{-1} - \lambda_1 \mathbf{P}_1^{-1} \mathbf{K}_1 \mathbf{P}_1^{-1} \\ \frac{\partial J(\mathbf{P}, \lambda_1, \lambda_4)}{\partial \mathbf{P}_4} &= \mathbf{M}_4(\mathbf{P}) - \mathbf{P}_4^{-1} \mathbf{N}_4(\mathbf{P}) \mathbf{P}_4^{-1} + \mathbf{W}_{4B} \\ &\quad - \mathbf{P}_4^{-1} \mathbf{K}_{4C} \mathbf{P}_4^{-1} - \lambda_4 \mathbf{P}_4^{-1} \mathbf{K}_4 \mathbf{P}_4^{-1} \end{aligned} \quad (25)$$

where

$$\mathbf{N}_A(\mathbf{P}) = \begin{bmatrix} \mathbf{N}_1(\mathbf{P}) & \mathbf{N}_2(\mathbf{P}) \\ \mathbf{N}_3(\mathbf{P}) & \mathbf{N}_4(\mathbf{P}) \end{bmatrix}$$

with an $m \times m$ submatrix $\mathbf{N}_1(\mathbf{P})$ and an $n \times n$ submatrix $\mathbf{N}_4(\mathbf{P})$ along its diagonal. Setting $\partial J(\mathbf{P}, \lambda_1, \lambda_4) / \partial \mathbf{P}_1 = \mathbf{0}$ and $\partial J(\mathbf{P}, \lambda_1, \lambda_4) / \partial \mathbf{P}_4 = \mathbf{0}$, it follows that

$$\begin{aligned} \mathbf{P}_1 \mathbf{F}_1(\mathbf{P}) \mathbf{P}_1 &= \mathbf{G}_1(\mathbf{P}, \lambda_1) \\ \mathbf{P}_4 \mathbf{F}_4(\mathbf{P}) \mathbf{P}_4 &= \mathbf{G}_4(\mathbf{P}, \lambda_4) \end{aligned} \quad (26)$$

where

$$\mathbf{F}_1(\mathbf{P}) = \mathbf{M}_1(\mathbf{P}) + \mathbf{W}_{1B}, \quad \mathbf{F}_4(\mathbf{P}) = \mathbf{M}_4(\mathbf{P}) + \mathbf{W}_{4B}$$

$$\mathbf{G}_1(\mathbf{P}, \lambda_1) = \mathbf{N}_1(\mathbf{P}) + \mathbf{K}_{1C} + \lambda_1 \mathbf{K}_1$$

$$\mathbf{G}_4(\mathbf{P}, \lambda_4) = \mathbf{N}_4(\mathbf{P}) + \mathbf{K}_{4C} + \lambda_4 \mathbf{K}_4.$$

The equations in (26) are highly nonlinear with respect to \mathbf{P}_1 and \mathbf{P}_4 . Namely, $\mathbf{P}_i \mathbf{F}_i(\mathbf{P}) \mathbf{P}_i$ for $i = 1, 4$ has a rational type R/D of nonlinearity, where the degree of the nonlinearity of R is $m+n+1$ and the degree of the nonlinearity of D is $m+n$, while $\mathbf{G}_i(\mathbf{P}, \lambda_i)$ for $i = 1, 4$ depends on \mathbf{P} linearly. An effective approach for solving these equations is to *relax* them into the following recursive second-order matrix equations:

$$\begin{aligned} \mathbf{P}_1^{(k+1)} \mathbf{F}_1(\mathbf{P}^{(k)}) \mathbf{P}_1^{(k+1)} &= \mathbf{G}_1(\mathbf{P}^{(k)}, \lambda_1^{(k+1)}) \\ \mathbf{P}_4^{(k+1)} \mathbf{F}_4(\mathbf{P}^{(k)}) \mathbf{P}_4^{(k+1)} &= \mathbf{G}_4(\mathbf{P}^{(k)}, \lambda_4^{(k+1)}) \end{aligned} \quad (27)$$

with initial condition $\mathbf{P}^{(0)} = \mathbf{P}_1^{(0)} \oplus \mathbf{P}_4^{(0)} = \mathbf{I}_{m+n}$. Noting that $\mathbf{P}\mathbf{W}\mathbf{P} = \mathbf{M}$ has the unique solution [5]

$$\mathbf{P} = \mathbf{W}^{-\frac{1}{2}}[\mathbf{W}^{\frac{1}{2}}\mathbf{M}\mathbf{W}^{\frac{1}{2}}]^{\frac{1}{2}}\mathbf{W}^{-\frac{1}{2}} \quad (28)$$

where $\mathbf{W} > 0$ and $\mathbf{M} \geq 0$ are symmetric, the unique solutions $\mathbf{P}_1^{(k+1)}$ and $\mathbf{P}_4^{(k+1)}$ of (27) are found to be

$$\begin{aligned} \mathbf{P}_1^{(k+1)} &= \mathbf{F}_1(\mathbf{P}^{(k)})^{-\frac{1}{2}}[\mathbf{F}_1(\mathbf{P}^{(k)})^{\frac{1}{2}} \\ &\quad \cdot \mathbf{G}_1(\mathbf{P}^{(k)}, \lambda_1^{(k+1)})\mathbf{F}_1(\mathbf{P}^{(k)})^{\frac{1}{2}}]^{\frac{1}{2}}\mathbf{F}_1(\mathbf{P}^{(k)})^{-\frac{1}{2}} \\ \mathbf{P}_4^{(k+1)} &= \mathbf{F}_4(\mathbf{P}^{(k)})^{-\frac{1}{2}}[\mathbf{F}_4(\mathbf{P}^{(k)})^{\frac{1}{2}} \\ &\quad \cdot \mathbf{G}_4(\mathbf{P}^{(k)}, \lambda_4^{(k+1)})\mathbf{F}_4(\mathbf{P}^{(k)})^{\frac{1}{2}}]^{\frac{1}{2}}\mathbf{F}_4(\mathbf{P}^{(k)})^{-\frac{1}{2}}. \end{aligned} \quad (29)$$

Here, the Lagrange multipliers $\lambda_1^{(k+1)}$ and $\lambda_4^{(k+1)}$ can be efficiently obtained using a bisection method [32] so that

$$\begin{aligned} f_1(\lambda_1^{(k+1)}) &= m - \text{tr}[\tilde{\mathbf{K}}_1^{(k)} \tilde{\mathbf{G}}_1^{(k)}(\lambda_1^{(k+1)})] = 0 \\ f_4(\lambda_4^{(k+1)}) &= n - \text{tr}[\tilde{\mathbf{K}}_4^{(k)} \tilde{\mathbf{G}}_4^{(k)}(\lambda_4^{(k+1)})] = 0 \end{aligned} \quad (30)$$

are satisfied where

$$\begin{aligned} \tilde{\mathbf{K}}_1^{(k)} &= \mathbf{F}_1(\mathbf{P}^{(k)})^{\frac{1}{2}}\mathbf{K}_1\mathbf{F}_1(\mathbf{P}^{(k)})^{\frac{1}{2}} \\ \tilde{\mathbf{K}}_4^{(k)} &= \mathbf{F}_4(\mathbf{P}^{(k)})^{\frac{1}{2}}\mathbf{K}_4\mathbf{F}_4(\mathbf{P}^{(k)})^{\frac{1}{2}} \\ \tilde{\mathbf{G}}_1^{(k)}(\lambda_1^{(k+1)}) &= [\mathbf{F}_1(\mathbf{P}^{(k)})^{\frac{1}{2}}\mathbf{G}_1(\mathbf{P}^{(k)}, \lambda_1^{(k+1)})\mathbf{F}_1(\mathbf{P}^{(k)})^{\frac{1}{2}}]^{-\frac{1}{2}} \\ \tilde{\mathbf{G}}_4^{(k)}(\lambda_4^{(k+1)}) &= [\mathbf{F}_4(\mathbf{P}^{(k)})^{\frac{1}{2}}\mathbf{G}_4(\mathbf{P}^{(k)}, \lambda_4^{(k+1)})\mathbf{F}_4(\mathbf{P}^{(k)})^{\frac{1}{2}}]^{-\frac{1}{2}}. \end{aligned}$$

The iteration process continues until

$$|J(\mathbf{P}^{(k)}, \lambda_1^{(k+1)}, \lambda_4^{(k+1)}) - J(\mathbf{P}^{(k-1)}, \lambda_1^{(k)}, \lambda_4^{(k)})| < \varepsilon \quad (31)$$

for a prescribed tolerance $\varepsilon > 0$. If the iteration is terminated at step k , then $\mathbf{P}^{(k)}$ is claimed to be a solution point.

It is noted that a straightforward extension of the iterative algorithm reported in [22] for updating λ_1 and λ_4 does not work well in this case, where we do not require $\text{tr}[\mathbf{K}\mathbf{P}^{-1}] = m+n$, but both $\text{tr}[\mathbf{K}_1\mathbf{P}_1^{-1}] = m$ and $\text{tr}[\mathbf{K}_4\mathbf{P}_4^{-1}] = n$. In addition, the bisection method offers an exponential convergence rate of $(1/2)^L$ where L is the number of iterations used. As such, accurate solutions of the Lagrange multipliers $\lambda_1^{(k+1)}$ and $\lambda_4^{(k+1)}$ in (30) can be identified with just a few iterations. In our simulation studies, the bisection method was found considerably faster than the iterative algorithm proposed in [22].

All the Gramians can be evaluated by truncating the corresponding infinite summations, see Appendix I for details.

3.2 An Unconstrained Optimization Method

By defining

$$\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4 = (\mathbf{T}_1 \oplus \mathbf{T}_4)^T (\mathbf{K}_1 \oplus \mathbf{K}_4)^{-\frac{1}{2}}, \quad (32)$$

it follows that

$$\mathbf{T}_i^{-1} \mathbf{K}_i \mathbf{T}_i^{-T} = \hat{\mathbf{T}}_i^{-T} \hat{\mathbf{T}}_i^{-1} \quad \text{for } i = 1, 4. \quad (33)$$

Thus, a convenient way to eliminate the l_2 -scaling constraints in (21) is to choose $\hat{\mathbf{T}}_1^{-1}$ and $\hat{\mathbf{T}}_4^{-1}$ as

$$\begin{aligned} \hat{\mathbf{T}}_1^{-1} &= \begin{bmatrix} \mathbf{t}_{11} & & \\ \|\mathbf{t}_{11}\| & & \\ & \mathbf{t}_{12} & \\ & \|\mathbf{t}_{12}\| & \\ & & \ddots & \\ & & & \mathbf{t}_{1m} \\ & & & \|\mathbf{t}_{1m}\| \end{bmatrix} \\ \hat{\mathbf{T}}_4^{-1} &= \begin{bmatrix} \mathbf{t}_{41} & & \\ \|\mathbf{t}_{41}\| & & \\ & \mathbf{t}_{42} & \\ & \|\mathbf{t}_{42}\| & \\ & & \ddots & \\ & & & \mathbf{t}_{4n} \\ & & & \|\mathbf{t}_{4n}\| \end{bmatrix} \end{aligned} \quad (34)$$

where \mathbf{t}_{1i} is an $m \times 1$ vector for $i = 1, 2, \dots, m$ and \mathbf{t}_{4j} is an $n \times 1$ vector for $j = 1, 2, \dots, n$. With (34), all the diagonal elements of $\hat{\mathbf{T}}_1^{-T} \hat{\mathbf{T}}_1^{-1}$ and $\hat{\mathbf{T}}_4^{-T} \hat{\mathbf{T}}_4^{-1}$ are found to be unity. Taking (32) and (34) into account, we conclude that the coordinate transformation matrix \mathbf{T} of the form

$$\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4 = (\mathbf{K}_1 \oplus \mathbf{K}_4)^{\frac{1}{2}} (\hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4)^T \quad (35)$$

automatically satisfies the l_2 -scaling constraints in (21). By substituting (35) into $S(\mathbf{P})$ in (20), the frequency-weighted l_2 -sensitivity measure can be expressed as

$$J_o(\mathbf{x}) = \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{M}}_A(\hat{\mathbf{P}}) \hat{\mathbf{T}}^T] + \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{W}}_B \hat{\mathbf{T}}^T] + \text{tr}[\hat{\mathbf{T}}^{-T} \hat{\mathbf{K}}_C \hat{\mathbf{T}}^{-1}] \quad (36)$$

with

$$\mathbf{x} = (\mathbf{t}_{11}^T, \mathbf{t}_{12}^T, \dots, \mathbf{t}_{1m}^T, \mathbf{t}_{41}^T, \mathbf{t}_{42}^T, \dots, \mathbf{t}_{4n}^T)^T$$

$$\hat{\mathbf{M}}_A(\hat{\mathbf{P}}) = \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \hat{\mathbf{Y}}(z_1, z_2) \hat{\mathbf{P}}^{-1} \hat{\mathbf{Y}}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2}$$

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

$$\begin{aligned} \hat{\mathbf{Y}}(z_1, z_2) &= (\mathbf{K}_1 \oplus \mathbf{K}_4)^{\frac{1}{2}} \mathbf{Y}(z_1, z_2) (\mathbf{K}_1 \oplus \mathbf{K}_4)^{-\frac{1}{2}} \\ \mathbf{Y}(z_1, z_2) &= \mathbf{W}_A(z_1, z_2) [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T \\ \hat{\mathbf{W}}_B &= (\mathbf{K}_1 \oplus \mathbf{K}_4)^{\frac{1}{2}} \mathbf{W}_B (\mathbf{K}_1 \oplus \mathbf{K}_4)^{\frac{1}{2}} \\ \hat{\mathbf{K}}_C &= (\mathbf{K}_1 \oplus \mathbf{K}_4)^{-\frac{1}{2}} \mathbf{K}_C (\mathbf{K}_1 \oplus \mathbf{K}_4)^{-\frac{1}{2}}. \end{aligned}$$

This shows that the problem of obtaining the block-diagonal nonsingular matrix $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ which minimizes $S(\mathbf{P})$ in (20) subject to the l_2 -scaling constraints in (21) can be converted into an unconstrained optimization problem of obtaining an $(m^2 + n^2) \times 1$ vector \mathbf{x} which minimizes $J_o(\mathbf{x})$ in (36).

Applying a quasi-Newton algorithm to minimize $J_o(\mathbf{x})$ in (36), in the k th iteration the most recent point \mathbf{x}_k is updated to point \mathbf{x}_{k+1} as [33]

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad (37)$$

where

$$\begin{aligned} \mathbf{d}_k &= -\mathbf{S}_k \nabla J_o(\mathbf{x}_k), \quad \alpha_k = \arg \min_{\alpha} J_o(\mathbf{x}_k + \alpha \mathbf{d}_k) \\ \mathbf{S}_{k+1} &= \mathbf{S}_k + \left(1 + \frac{\gamma_k^T \mathbf{S}_k \gamma_k}{\gamma_k^T \delta_k} \right) \frac{\delta_k \delta_k^T}{\gamma_k^T \delta_k} - \frac{\delta_k \gamma_k^T \mathbf{S}_k + \mathbf{S}_k \gamma_k \delta_k^T}{\gamma_k^T \delta_k} \\ \mathbf{S}_0 &= \mathbf{I}_{m^2+n^2}, \quad \delta_k = \mathbf{x}_{k+1} - \mathbf{x}_k \\ \gamma_k &= \nabla J_o(\mathbf{x}_{k+1}) - \nabla J_o(\mathbf{x}_k). \end{aligned}$$

Here, $\nabla J_o(\mathbf{x})$ is the gradient of $J_o(\mathbf{x})$ with respect to \mathbf{x} , and \mathbf{S}_k is a positive-definite approximation of the inverse Hessian matrix of $J_o(\mathbf{x})$. The algorithm starts with a trivial initial point

\mathbf{x}_0 obtained from an initial assignment $\hat{\mathbf{T}} = \mathbf{I}_{m+n}$, and this iteration process continues until

$$|J_o(\mathbf{x}_{k+1}) - J_o(\mathbf{x}_k)| < \varepsilon \quad (38)$$

where $\varepsilon > 0$ is a prescribed tolerance.

The implementation of (37) requires the gradient of $J_o(\mathbf{x})$, which can be efficiently evaluated using closed-form expressions, see Appendix II for details.

IV. CASE STUDIES

In this section, we examine the proposed algorithms by applying them to a recursive 2-D state-space digital filter. In addition, the algorithms are applied to a 2-D filter utilized in [22] and the results are compared with those obtained by the method of [22].

Example 1: Consider a 2-D stable recursive digital filter realization $(\mathbf{A}^o, \mathbf{b}^o, \mathbf{c}^o, d)_{2,2}$ where [30, p. 971]

$$\mathbf{A}^o = \begin{bmatrix} 1.88899 & -0.91219 & -1.00000 & 0.00000 \\ 1.00000 & 0.00000 & 0.00000 & 0.00000 \\ 0.02771 & -0.02580 & 1.88899 & 1.00000 \\ -0.02580 & 0.02431 & -0.91219 & 0.00000 \end{bmatrix}$$

$$\mathbf{b}^o = [0.219089 \quad 0.000000 \quad -0.028889 \quad 0.091219]^T$$

$$\mathbf{c}^o = [0.028889 \quad -0.091219 \quad -0.219089 \quad 0.000000]$$

$$d = 0.08900.$$

After carrying out the l_2 -scaling for the above realization with a diagonal coordinate matrix

$$\mathbf{T}^o = \text{diag}\{9.336421, 9.336414, 1.065102, 0.986642\},$$

we obtain the 2-D state-space digital filter $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{2,2}$ characterized by

$$\mathbf{A} = \begin{bmatrix} 1.888990 & -0.912189 & -0.114080 & 0.000000 \\ 1.000001 & 0.000000 & 0.000000 & 0.000000 \\ 0.242899 & -0.226156 & 1.888990 & 0.926336 \\ -0.244141 & 0.230041 & -0.984729 & 0.000000 \end{bmatrix}$$

$$\mathbf{b} = [0.023466 \quad 0.000000 \quad -0.027123 \quad 0.092454]^T$$

$$\mathbf{c} = [0.269720 \quad -0.851658 \quad -0.233352 \quad 0.000000]$$

$$d = 0.08900$$

where $\mathbf{A} = \mathbf{T}^{o-1} \mathbf{A}^o \mathbf{T}^o$, $\mathbf{b} = \mathbf{T}^{o-1} \mathbf{b}^o$, and $\mathbf{c} = \mathbf{c}^o \mathbf{T}^o$. The 2-D state-space digital filter $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{2,2}$ now satisfies the l_2 -scaling constraints. The frequency-weighting functions used in this example were given by a 2-D nonrecursive low-pass digital filter with the following unit-sample response [34, p. 895]:

$$w_A(i, j) = w_B(i, j) = w_C(i, j)$$

$$= 0.256322 \exp[-0.103203\{(i-4)^2 + (j-4)^2\}]$$

for $(0, 0) \leq (i, j) \leq (20, 20)$, and zero elsewhere.

Using (A.1) and (A.2) with truncation $(0, 0) \leq (i, j) \leq (150, 150)$ to evaluate the Gramians \mathbf{K}_C , \mathbf{W}_B , \mathbf{M}_A and \mathbf{K} ,

it was found that

$$\mathbf{K}_C = \begin{bmatrix} 32.944701 & 32.414885 & 2.177594 & -2.390760 \\ 32.414885 & 32.944694 & 2.732580 & -2.852159 \\ 2.177594 & 2.732580 & 4.347875 & -4.136567 \\ -2.390760 & -2.852159 & -4.136567 & 4.056383 \end{bmatrix}$$

$$\mathbf{W}_B = 10^3 \begin{bmatrix} 0.429952 & -0.378923 & 0.215338 & 0.250313 \\ -0.378923 & 0.344207 & -0.219002 & -0.242021 \\ 0.215338 & -0.219002 & 3.257832 & 2.969311 \\ 0.250313 & -0.242021 & 2.969311 & 2.795534 \end{bmatrix}$$

$$\mathbf{M}_A = 10^5 \begin{bmatrix} 0.599216 & -0.523300 & 0.711632 & 0.788400 \\ -0.523300 & 0.466607 & -0.639356 & -0.707249 \\ 0.711632 & -0.639356 & 6.203669 & 5.638379 \\ 0.788400 & -0.707249 & 5.638379 & 5.322915 \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} 1.000000 & 0.978030 & 0.164886 & -0.167063 \\ 0.978030 & 1.000000 & 0.132847 & -0.133855 \\ 0.164886 & 0.132847 & 1.000000 & -0.985382 \\ -0.167063 & -0.133855 & -0.985382 & 1.000000 \end{bmatrix}$$

In what follows, we evaluate the frequency-weighted l_2 -sensitivity of the 2-D state-space digital filter $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{2,2}$. Later on, this sensitivity measure will be used as a benchmark in the examination of the performance of the proposed algorithms. Using (17), the frequency-weighted l_2 -sensitivity of the LSS model $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{2,2}$ was found to be

$$S = 126.614237 \times 10^4.$$

As will be seen next, the proposed algorithms are able to deduce an equivalent state-space realization with much reduced frequency-weighted l_2 -sensitivity.

4.1 Application of the Lagrange method

Choosing $\mathbf{P}^{(0)} = \mathbf{P}_1^{(0)} \oplus \mathbf{P}_4^{(0)} = \mathbf{I}_4$ in (29) as an initial estimate and a tolerance $\varepsilon = 10^{-8}$ in (31) as well as in the bisection method, it took the Lagrange-based algorithm 10 iterations to converge to the solution

$$\mathbf{P}^{opt} = \begin{bmatrix} 1.639466 & 1.715106 \\ 1.715106 & 1.828836 \end{bmatrix} \oplus \begin{bmatrix} 0.901558 & -0.915221 \\ -0.915221 & 0.962667 \end{bmatrix}$$

or equivalently,

$$\mathbf{T}^{opt} = \begin{bmatrix} 1.142108 & 0.578841 \\ 1.110711 & 0.771464 \end{bmatrix} \oplus \begin{bmatrix} 0.266022 & -0.911477 \\ -0.094157 & 0.976628 \end{bmatrix}.$$

The minimized frequency-weighted l_2 -sensitivity measure in (23) corresponding to the above solution was found to be

$$J(\mathbf{P}^{opt}, \lambda_1, \lambda_4) = 4.076278 \times 10^4$$

with $\lambda_1 = -16880.585503$ and $\lambda_4 = 16962.901711$, and the optimal state-space filter structure $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{2,2}$ (that

minimizes (20) subject to the l_2 -scaling constraints in (21)) was synthesized by substituting matrix \mathbf{T}^{opt} into (19) as

$$\bar{\mathbf{A}} = \begin{bmatrix} 0.930636 & -0.144497 & -0.098301 & 0.336809 \\ 0.140565 & 0.958354 & 0.141528 & -0.484919 \\ 0.024994 & -0.000749 & 0.958805 & 0.115614 \\ -0.021475 & 0.036943 & -0.175791 & 0.930185 \end{bmatrix}$$

$$\bar{\mathbf{b}} = \begin{bmatrix} 0.076010 & -0.109434 & 0.332103 & 0.126685 \end{bmatrix}^T$$

$$\bar{\mathbf{c}} = \begin{bmatrix} -0.637897 & -0.500899 & -0.062077 & 0.212695 \end{bmatrix}.$$

The profile of the frequency-weighted l_2 -sensitivity measure $J(\mathbf{P}, \lambda_1, \lambda_4)$ and the profiles of the Lagrange multipliers λ_1 and λ_4 for the first 10 iterations are shown in Figs. 2 and 3, respectively.

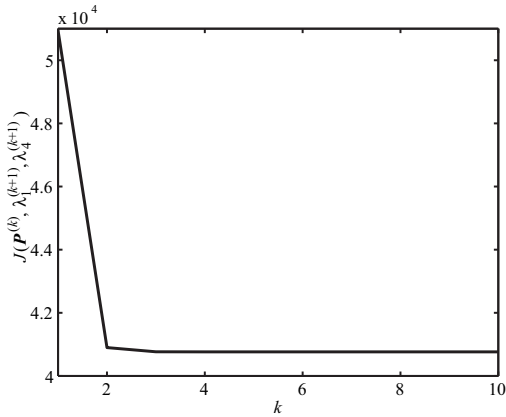


Fig. 2. Profile of $J(\mathbf{P}, \lambda_1, \lambda_4)$ during the first 10 iterations.

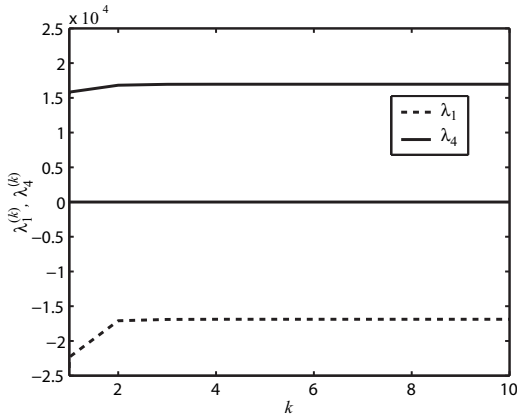


Fig. 3. Profiles of λ_1 and λ_4 during the first 10 iterations.

4.2 Application of the quasi-Newton method

By choosing $\hat{\mathbf{T}} = \mathbf{I}_2 \oplus \mathbf{I}_2$ (therefore $\mathbf{T} = (\mathbf{K}_1 \oplus \mathbf{K}_4)^{1/2}$ in (32)) as an initial estimate and a tolerance $\varepsilon = 10^{-8}$ in (38), the quasi-Newton algorithm took 16 iterations to converge to the solution

$$\hat{\mathbf{T}}^{opt} = \begin{bmatrix} 0.906814 & 0.694991 \\ -0.169732 & 1.129830 \end{bmatrix} \oplus \begin{bmatrix} 0.838552 & 0.582931 \\ -0.400572 & 0.939425 \end{bmatrix}$$

or equivalently,

$$\mathbf{T}^{opt} = \begin{bmatrix} 1.142108 & 0.578841 \\ 1.110711 & 0.771464 \end{bmatrix} \oplus \begin{bmatrix} 0.266022 & -0.911477 \\ -0.094157 & 0.976628 \end{bmatrix}$$

and the minimized frequency-weighted l_2 -sensitivity was found to be

$$J_o(\hat{\mathbf{T}}^{opt}) = 4.076278 \times 10^4.$$

We see that the results obtained by this method are identical to those obtained by the Lagrange-based method. The profile of the l_2 -sensitivity measure $J_o(\hat{\mathbf{T}})$ during the first 16 iterations is shown in Fig. 4.

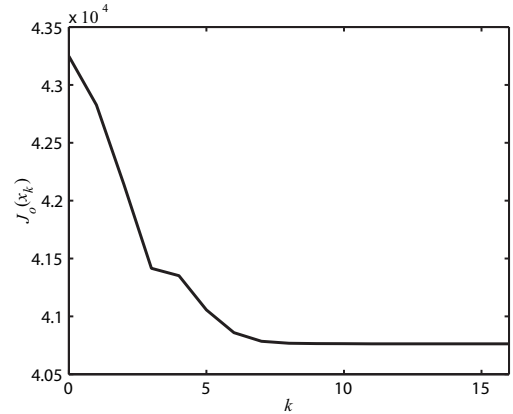


Fig. 4. Profile of $J_o(\hat{\mathbf{T}})$ during the first 16 iterations.

We now explain the effectiveness of the optimal realization that minimizes the frequency-weighted l_2 -sensitivity subject to l_2 -scaling constraints. The magnitude response of the original realization $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{2,2}$ is shown in Fig. 5, where the maximum value and the l_2 -norm (which corresponds to a square root of the summation of squared values at 201×201 sampling points) were 11.975265 and 133.212501, respectively. When all coefficients in the original realization were rounded to power-of-two representation with 8 bits after binary point, the magnitude-response deviation between the original realization and that with rounded coefficients is shown in Fig. 6, where the maximum deviation and the l_2 -norm were 4.141411 and 22.563131, respectively. Alternatively, when all coefficients in the optimal realization were rounded in the same manner, the magnitude-response deviation between the original realization and the optimal one with rounded coefficients is shown in Fig. 7, where the maximum deviation and the l_2 -norm were 0.617833 and 2.323280, respectively. From these figures and data, it is observed that the coefficient sensitivity of the optimal realization is considerably lower than that of the original realization.

Example 2: Let a 2-D stable recursive digital filter realization be specified by $(\mathbf{A}^o, \mathbf{b}^o, \mathbf{c}^o, d)_{4,4}$ where

$$\mathbf{A}^o = \begin{bmatrix} \mathbf{A}_1^o & \mathbf{A}_2^o \\ \mathbf{A}_1^o & \mathbf{A}_2^o \end{bmatrix}, \quad \mathbf{b}^o = \begin{bmatrix} \mathbf{b}_1^o \\ \mathbf{b}_2^o \end{bmatrix}, \quad \mathbf{c}^o = \begin{bmatrix} \mathbf{c}_1^o & \mathbf{c}_2^o \end{bmatrix}$$

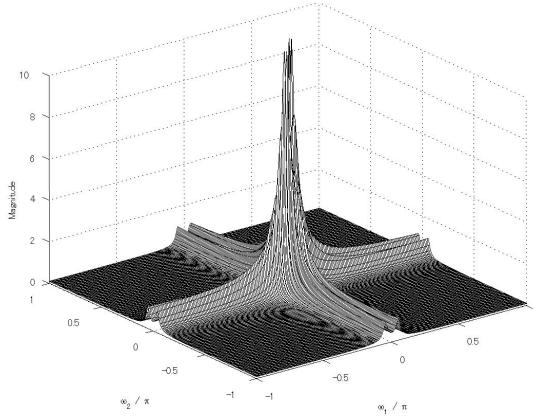


Fig. 5. Magnitude response of the original realization.

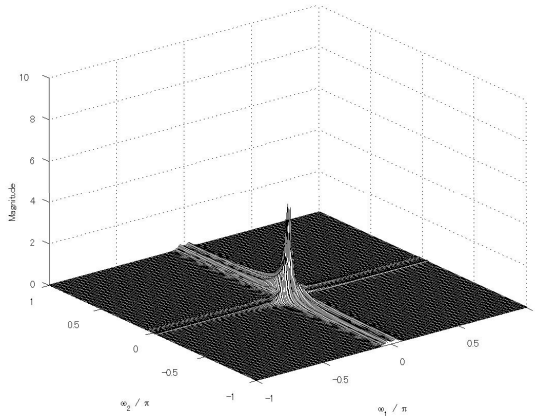


Fig. 6. Magnitude-response deviation between the original realization and that with rounded coefficients.

with

$$\mathbf{A}_1^o = \begin{bmatrix} 0.0 & 0.481228 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.510378 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.525287 \\ -0.031857 & 0.298663 & -0.808282 & 1.044600 \end{bmatrix}$$

$$\mathbf{A}_2^o = \begin{bmatrix} -0.226080 & 0.776837 & 0.024693 & -0.000933 \\ -0.843550 & 1.610400 & -0.309366 & 0.065898 \\ -1.260339 & 2.005100 & -0.453220 & 0.203118 \\ -1.121498 & 1.636435 & -0.590516 & 0.562890 \end{bmatrix}$$

$$\mathbf{b}_1^o = \mathbf{b}_2^o = [0.0 \quad 0.0 \quad 0.0 \quad 0.198473]^T$$

$$\mathbf{c}_1^o = [-0.567054 \quad 0.231913 \quad 0.197016 \quad 0.239932]$$

$$\mathbf{c}_2^o = [0.464344 \quad 0.441837 \quad -0.061100 \quad 0.105505]$$

$$d = 0.009430.$$

This 2-D filter was obtained by imbedding the LSS model of Example 2 in [22] into the Roesser LSS model. By performing the l_2 -scaling for the above realization with a diagonal coordinate matrix

$$\mathbf{T}^o = \text{diag}\{1.000001, 1.000002, 1.000003, 1.000003, \\ 1.000001, 1.000002, 1.000003, 1.000003\}$$

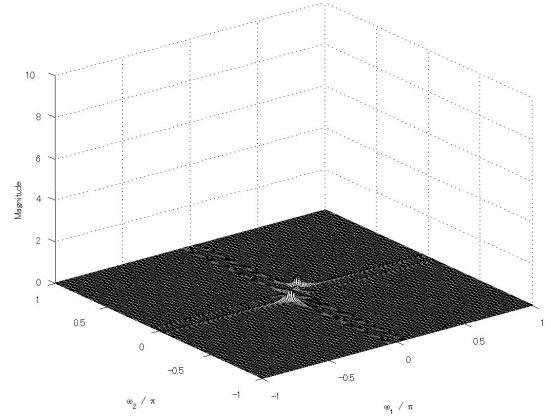


Fig. 7. Magnitude-response deviation between the original realization and the optimal one with rounded coefficients.

and then applying the same frequency-weighted functions as in Example 1 to the resulting realization $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{4,4}$, the frequency-weighted l_2 -sensitivity in (17) was found to be

$$S = 394.423680 \times 10^3$$

with truncation $(0, 0) \leq (i, j) \leq (100, 100)$ in (A.1) and (A.2).

4.3 Application of the Lagrange method

Choosing $\mathbf{P}^{(0)} = \mathbf{P}_1^{(0)} \oplus \mathbf{P}_4^{(0)} = \mathbf{I}_8$ in (29) as an initial estimate and a tolerance $\varepsilon = 10^{-8}$ in (31) as well as in the bisection method, it took the Lagrange-based algorithm 104 iterations to converge to the solution

$$\mathbf{P}^{opt} = \begin{bmatrix} 3.947574 & 3.057718 & 2.498923 & 2.123649 \\ 3.057718 & 2.599605 & 2.183384 & 1.861281 \\ 2.498923 & 2.183384 & 1.894674 & 1.650080 \\ 2.123649 & 1.861281 & 1.650080 & 1.483679 \end{bmatrix} \\ \oplus \begin{bmatrix} 2.234666 & 1.993263 & 1.781601 & 1.653502 \\ 1.993263 & 1.814492 & 1.653747 & 1.534665 \\ 1.781601 & 1.653747 & 1.556000 & 1.464299 \\ 1.653502 & 1.534665 & 1.464299 & 1.432256 \end{bmatrix}$$

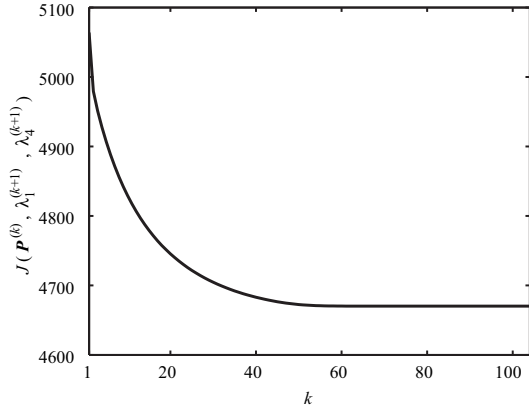
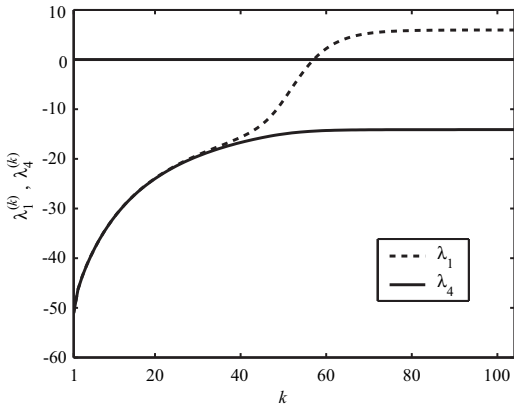
or equivalently,

$$\mathbf{T}^{opt} = \begin{bmatrix} 0.614712 & 1.232311 & 1.040088 & 0.984546 \\ 0.899544 & 0.797151 & 0.663356 & 0.845538 \\ 0.874957 & 0.504513 & 0.679077 & 0.642997 \\ 0.727174 & 0.269835 & 0.705983 & 0.619414 \end{bmatrix} \\ \oplus \begin{bmatrix} 0.531664 & 0.880708 & 0.800816 & 0.731468 \\ 0.620875 & 0.737566 & 0.617178 & 0.709996 \\ 0.753633 & 0.564844 & 0.551426 & 0.604084 \\ 0.723352 & 0.347355 & 0.658078 & 0.596067 \end{bmatrix}.$$

The minimized frequency-weighted l_2 -sensitivity measure in (23) corresponding to the above solution was found to be

$$J(\mathbf{P}^{opt}, \lambda_1, \lambda_4) = 4.670177 \times 10^3$$

with $\lambda_1 = 5.955188$ and $\lambda_4 = -14.111003$. The profile of the frequency-weighted l_2 -sensitivity measure $J(\mathbf{P}, \lambda_1, \lambda_4)$ and the profiles of the Lagrange multipliers λ_1 and λ_4 for the first 104 iterations are shown in Figs. 8 and 9, respectively.


 Fig. 8. Profile of $J(\mathbf{P}, \lambda_1, \lambda_4)$ during the first 104 iterations.

 Fig. 9. Profiles of λ_1 and λ_4 during the first 104 iterations.

4.4 Application of the quasi-Newton method

By choosing $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4 = \mathbf{I}_8$ in (37) as an initial estimate and a tolerance $\varepsilon = 10^{-8}$ in (38), the quasi-Newton algorithm took 54 iterations to converge to

$$\hat{\mathbf{T}}^{opt} = \begin{bmatrix} 3.056671 & -2.673365 & 0.575882 & -0.429287 \\ -0.331629 & 2.142411 & -0.401503 & -0.192081 \\ -2.530651 & 0.932586 & 0.553002 & -0.136935 \\ 1.754363 & -0.312582 & 0.624509 & 0.515370 \end{bmatrix} \oplus \begin{bmatrix} 1.307170 & -0.419919 & 0.045538 & -0.194118 \\ 0.762443 & 0.830435 & -0.297531 & 0.062104 \\ -0.405202 & 0.189220 & 0.976564 & -0.250656 \\ 1.071478 & -0.069804 & 0.315533 & 0.828727 \end{bmatrix}$$

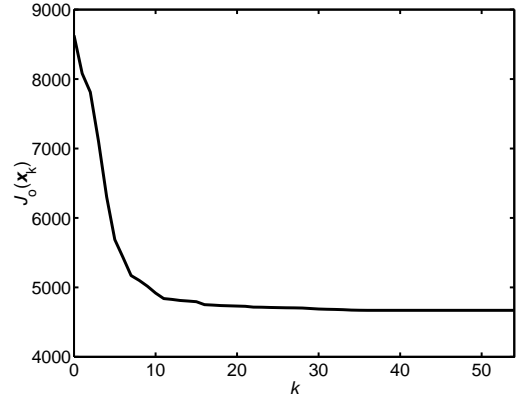
or equivalently,

$$\mathbf{T}^{opt} = \begin{bmatrix} 0.690639 & 0.697890 & -0.986414 & 1.417930 \\ 0.205523 & 0.802226 & -0.589043 & 1.251725 \\ 0.091157 & 0.584768 & -0.335877 & 1.196490 \\ -0.023116 & 0.340604 & -0.305900 & 1.128518 \end{bmatrix} \oplus \begin{bmatrix} 0.595272 & 0.843931 & 0.159783 & 1.068904 \\ 0.406418 & 0.767684 & 0.282059 & 0.990157 \\ 0.270942 & 0.580437 & 0.379363 & 1.000879 \\ 0.148987 & 0.449337 & 0.230545 & 1.074708 \end{bmatrix}$$

The minimized frequency-weighted l_2 -sensitivity in (36) was found to be

$$J_o(\hat{\mathbf{T}}^{opt}) = 4.670177 \times 10^3$$

which is identical to the minimum value of the frequency-weighted l_2 -sensitivity measure, obtained by the Lagrange-based method. The profile of the l_2 -sensitivity measure $J_o(\hat{\mathbf{T}})$ during the first 54 iterations is shown in Fig. 10.


 Fig. 10. Profile of $J_o(\hat{\mathbf{T}})$ during the first 54 iterations.

To compare the technique reported in [22] with the proposed ones, the optimal realization derived in Example 2 of [22] was imbedded in the Roesser LSS model by using (6). Then the frequency-weighted l_2 -sensitivity of the resulting imbedded model was computed from (17) as

$$S = 5.564659 \times 10^3.$$

It is observed that this value is 1.192 times greater than the minimized frequency-weighted l_2 -sensitivity obtained by the proposed techniques. Moreover, the proposed Lagrange-based method attains considerably faster convergence than that reported in [22], which required 2000 iterations for its convergence.

V. CONCLUSION

We have investigated the problem of minimizing the frequency-weighted l_2 -sensitivity subject to l_2 -scaling constraints for 2-D state-space digital filters described by the Roesser LSS model. It has been shown that the FM second LSS model can be imbedded in the Roesser LSS model as a special case. Two iterative methods have been developed to solve the problem at hand. The first iterative method is based on the introduction of a Lagrange function and makes use of an efficient bisection method. In our simulation studies, the bisection method was found to be considerably faster than the iterative method proposed in [22]. The second iterative method relies on the conversion of the constrained optimization problem into an unconstrained optimization formulation and utilizes an efficient quasi-Newton algorithm. The optimal state-space realization with minimum frequency-weighted l_2 -sensitivity and no overflow has then been constructed by applying an appropriate coordinate transformation. Our computer simulation results have demonstrated the validity and effectiveness of the proposed techniques.

$$\begin{aligned}
 \begin{bmatrix} \mathbf{x}^h(i+1, j+1) \\ \mathbf{x}^v(i+1, j+1) \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}^h(i, j+1) \\ \mathbf{x}^v(i, j+1) \end{bmatrix} \\
 &+ \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix} \begin{bmatrix} \mathbf{x}^h(i+1, j) \\ \mathbf{x}^v(i+1, j) \end{bmatrix} \\
 &+ \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{0} \end{bmatrix} u(i, j+1) + \begin{bmatrix} \mathbf{0} \\ \mathbf{b}_2 \end{bmatrix} u(i+1, j) \\
 y(i, j) &= [\mathbf{c}_1 \quad \mathbf{c}_2] \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + du(i, j)
 \end{aligned} \tag{39}$$

APPENDIX I

COMPUTATIONS OF GRAMIANS

The matrices \mathbf{K}_C , \mathbf{W}_B , and \mathbf{M}_A can be computed using

$$\begin{aligned}
 \mathbf{K}_C &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}_C(i, j) \mathbf{f}_C^T(i, j) \\
 \mathbf{W}_B &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{g}_B^T(i, j) \mathbf{g}_B(i, j) \\
 \mathbf{M}_A &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{H}_A^T(i, j) \mathbf{H}_A(i, j)
 \end{aligned} \tag{A.1}$$

where

$$\begin{aligned}
 \mathbf{A}^{(1,0)} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{A}^{(0,1)} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix} \\
 \mathbf{A}^{(0,0)} &= \mathbf{I}_{m+n}, \quad \mathbf{A}^{(-i,j)} = \mathbf{0} \ (i \geq 1), \quad \mathbf{A}^{(i,-j)} = \mathbf{0} \ (j \geq 1) \\
 \mathbf{A}^{(i,j)} &= \mathbf{A}^{(1,0)} \mathbf{A}^{(i-1,j)} + \mathbf{A}^{(0,1)} \mathbf{A}^{(i,j-1)} \\
 &= \mathbf{A}^{(i-1,j)} \mathbf{A}^{(1,0)} + \mathbf{A}^{(i,j-1)} \mathbf{A}^{(0,1)}, \quad (i, j) > (0, 0) \\
 \mathbf{f}(i, j) &= \mathbf{A}^{(i-1,j)} \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{0} \end{bmatrix} + \mathbf{A}^{(i,j-1)} \begin{bmatrix} \mathbf{0} \\ \mathbf{b}_2 \end{bmatrix} \\
 \mathbf{g}(i, j) &= \mathbf{c} \mathbf{A}^{(i-1,j)} \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \mathbf{c} \mathbf{A}^{(i,j-1)} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \\
 \mathbf{H}(i, j) &= \sum_{(0,0) \leq (k,r) < (i,j)} \mathbf{f}(k, r) \mathbf{g}(i-k, j-r) \\
 \mathbf{f}_C(i, j) &= \sum_{(0,0) \leq (k,r) < (i,j)} w_C(k, r) \mathbf{f}(i-k, j-r) \\
 \mathbf{g}_B(i, j) &= \sum_{(0,0) \leq (k,r) < (i,j)} w_B(k, r) \mathbf{g}(i-k, j-r) \\
 \mathbf{H}_A(i, j) &= \sum_{(0,0) \leq (k,r) < (i,j)} w_A(k, r) \mathbf{H}(i-k, j-r)
 \end{aligned}$$

with partial ordering for integer pairs (i, j) [26, p. 2], and $w_A(k, r)$, $w_B(k, r)$, and $w_C(k, r)$ denoting the unit-sample responses of frequency-weighting functions $W_A(z_1, z_2)$, $W_B(z_1, z_2)$, and $W_C(z_1, z_2)$, respectively.

The local controllability Gramian \mathbf{K} and the other Gramians

$\mathbf{M}_A(\mathbf{P})$, $\mathbf{N}_A(\mathbf{P})$ and $\hat{\mathbf{M}}_A(\hat{\mathbf{P}})$ can be computed using

$$\begin{aligned}
 \mathbf{K} &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}(i, j) \mathbf{f}^T(i, j) \\
 \mathbf{M}_A(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{H}_A^T(i, j) \mathbf{P}^{-1} \mathbf{H}_A(i, j) \\
 \mathbf{N}_A(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{H}_A(i, j) \mathbf{P} \mathbf{H}_A^T(i, j) \\
 \hat{\mathbf{M}}_A(\hat{\mathbf{P}}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \hat{\mathbf{H}}_A^T(i, j) \hat{\mathbf{P}}^{-1} \hat{\mathbf{H}}_A(i, j)
 \end{aligned} \tag{A.2}$$

where

$$\hat{\mathbf{H}}_A(i, j) = (\mathbf{K}_1 \oplus \mathbf{K}_4)^{-\frac{1}{2}} \mathbf{H}_A(i, j) (\mathbf{K}_1 \oplus \mathbf{K}_4)^{-\frac{1}{2}}.$$

APPENDIX II

 GRADIENT EVALUATION OF $J_o(\mathbf{x})$

$$\begin{aligned}
 \frac{\partial J_o(\hat{\mathbf{T}})}{\partial t_{ij}} &= \lim_{\Delta \rightarrow 0} \frac{J_o(\hat{\mathbf{T}}_{ij}) - J_o(\hat{\mathbf{T}})}{\Delta} \\
 &= 2\beta_1 - 2\beta_2 + 2\beta_3 - 2\beta_4
 \end{aligned} \tag{A.3}$$

where $\hat{\mathbf{T}}_{ij}$ is the matrix obtained from $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4$ with a perturbed (i, j) th component, which is given by [35, p. 655]

$$\begin{aligned}
 \hat{\mathbf{T}}_{ij} &= \hat{\mathbf{T}} + \frac{\Delta \hat{\mathbf{T}} \mathbf{g}_{ij} \mathbf{e}_j^T \hat{\mathbf{T}}}{1 - \Delta \mathbf{e}_j^T \hat{\mathbf{T}} \mathbf{g}_{ij}}, \quad \hat{\mathbf{T}}_{ij}^{-1} = \hat{\mathbf{T}}^{-1} - \Delta \mathbf{g}_{ij} \mathbf{e}_j^T \\
 \mathbf{g}_{ij} &= \partial \left\{ \frac{t_j}{\|t_j\|} \right\} / \partial t_{ij} = \frac{1}{\|t_j\|^3} (t_{ij} t_j - \|t_j\|^2 \mathbf{e}_i) \\
 \beta_1 &= \mathbf{e}_j^T \hat{\mathbf{M}}_A(\hat{\mathbf{T}}) \hat{\mathbf{T}} \mathbf{g}_{ij} \\
 \beta_2 &= \mathbf{e}_j^T \hat{\mathbf{T}}^{-T} \left[\sum_{p=0}^{\infty} \sum_{q=0}^{\infty} \hat{\mathbf{H}}_A(p, q) \hat{\mathbf{T}}^T \hat{\mathbf{T}} \hat{\mathbf{H}}_A^T(p, q) \right] \mathbf{g}_{ij} \\
 \beta_3 &= \mathbf{e}_j^T \hat{\mathbf{T}} \hat{\mathbf{W}}_B \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij}, \quad \beta_4 = \mathbf{e}_j^T \hat{\mathbf{T}}^{-T} \hat{\mathbf{K}}_c \mathbf{g}_{ij}
 \end{aligned}$$

with

$$\{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_{m+n}\} = \{\mathbf{t}_{11}, \mathbf{t}_{12}, \dots, \mathbf{t}_{1m}\} \cup \{\mathbf{t}_{41}, \mathbf{t}_{42}, \dots, \mathbf{t}_{4n}\}.$$

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, pp. 884-888, Oct. 1984.
- [3] L. Thiele, "On the sensitivity of linear state-space systems," *IEEE Trans. Circuits Syst.*, vol. CAS-33, 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, 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, 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, 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. Acoust., Speech, Signal Processing*, vol. 4, pp. 429-432.
- [8] M. Gevers and G. Li, *Parameterizations in Control, Estimation and Filtering Problems: Accuracy Aspects*, Springer-Verlag, 1993.
- [9] C. Xiao, "Improved L_2 -sensitivity for state-space digital system," *IEEE Trans. Signal Processing*, vol. 45, 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*, vol. 49, pp. 1279-1289, Sept. 2002.
- [11] S. Yamaki, M. Abe and M. Kawamata, "A closed form solution to L_2 -sensitivity minimization of second-order state-space digital filters," in *Proc. 2006 IEEE Int. Symp. Circuits Syst.*, pp. 5223-5226.
- [12] 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 Syst.*, pp. 710-713.
- [13] 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 Processing*, vol. ASSP-38, pp. 1587-1594, Sept. 1990.
- [14] T. Hinamoto, T. Takao and M. Muneyasu, "Synthesis of 2-D separable-denominator digital filters with low sensitivity," *J. Franklin Institute*, vol. 329, pp. 1063-1080, 1992.
- [15] T. Hinamoto and T. Takao, "Synthesis of 2-D state-space filter structures with low frequency-weighted sensitivity," *IEEE Trans. Circuits Syst. II*, vol. 39, pp. 646-651, Sept. 1992.
- [16] T. Hinamoto and T. Takao, "Minimization of frequency-weighting sensitivity in 2-D systems based on the Fornasini-Marchesini second model," in *1992 IEEE Int. Conf. Acoust., Speech, Signal Processing*, pp. 401-404.
- [17] 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*, vol. 46, pp. 1172-1183, Oct. 1999.
- [18] G. Li, "On frequency weighted minimal L_2 sensitivity of 2-D systems using Fornasini-Marchesini LSS model", *IEEE Trans. Circuits Syst. I*, vol. 44, pp. 642-646, July 1997.
- [19] G. Li, "Two-dimensional system optimal realizations with L_2 -sensitivity minimization," *IEEE Trans. Signal Processing*, vol. 46, pp. 809-813, Mar. 1998.
- [20] T. Hinamoto and Y. Sugie, " L_2 -sensitivity analysis and minimization of 2-D separable-denominator state-space digital filters," *IEEE Trans. Signal Processing*, vol. 50, pp. 3107-3114, Dec. 2002.
- [21] T. Hinamoto, H. Ohnishi and W.-S. Lu, "Minimization of L_2 -sensitivity for state-space digital filters subject to L_2 -dynamic-range scaling constraints," *IEEE Trans. Circuits Syst.-II*, vol. 52, pp. 641-645, Oct. 2005.
- [22] T. Hinamoto, K. Iwata and W.-S. Lu, " L_2 -sensitivity Minimization of one- and two-dimensional state-space digital filters subject to L_2 -scaling constraints," *IEEE Trans. Signal Processing*, vol. 54, pp. 1804-1812, May 2006.
- [23] S. Yamaki, M. Abe and M. Kawamata, "A novel approach to L_2 -sensitivity minimization of digital filters subject to L_2 -scaling constraints," in *Proc. 2006 IEEE Int. Symp. Circuits Syst.*, pp. 5219-5222.
- [24] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed-point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, pp. 551-562, Sept. 1976.
- [25] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-25, pp. 273-281, Aug. 1977.
- [26] R. P. Roesser, "A discrete state-space model for linear image processing," *IEEE Trans. Automat. Contr.*, vol. AC-20, pp. 1-10, Feb. 1975.
- [27] E. Fornasini and G. Marchesini, "Doubly-indexed dynamical systems: State-space models and structural properties," *Math Syst. Theory*, vol. 12, pp. 59-72, 1978.
- [28] T. Hinamoto, "A novel local state-space model for 2-D digital filters and its properties," in *Proc. 2001 IEEE Int. Symp. Circuits Syst.*, vol.2, pp. 545-548.
- [29] S. Kung, B. C. Levy, M. Morf and T. Kailath, "New results in 2-D systems theory, Part II: 2-D state-space model—Realization and the notions of controllability, observability, and minimality," *Proc. IEEE*, vol. 65, pp. 945-961, June 1977.
- [30] W.-S. Lu and A. Antoniou, "Synthesis of 2-D state-space fixed-point digital filter structures with minimum roundoff noise," *IEEE Trans. Circuits Syst.*, vol. CAS-33, pp. 965-973, Oct. 1986.
- [31] L. L. Scharf, *Statistical Signal Processing*, Reading, MA: Addison-Wesley, 1991.
- [32] H. Togawa, *Handbook of Numerical Methods*, Saiensu-sha, Tokyo, 1992.
- [33] R. Fletcher, *Practical Methods of Optimization*, 2nd ed., Wiley, New York, 1987.
- [34] S. A. H. Aly and M. M. Fahmy, "Spatial-domain design of two-dimensional recursive digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-27, pp. 892-901, Oct. 1980.
- [35] T. Kailath, *Linear System*, Englewood Cliffs, N.J.: Prentice Hall, 1980.