

An Analytical Approach for the Synthesis of Two-Dimensional State-Space Filter Structures with Minimum Weighted Sensitivity

Takao Hinamoto, *Senior Member, IEEE*, Yoshitaka Zempo, Yoshio Nishino, and Wu-Sheng Lu, *Fellow, IEEE*

Abstract—This paper considers the problem of synthesizing the finite-word-length (FWL) two-dimensional (2-D) state-space filter structures with minimum weighted sensitivity. Two kinds of frequency-weighted sensitivity measures, one based on a mixture of L_1/L_2 norms and the other a pure L_2 norm, are defined in place of the usual sensitivity measure and an upper bound expressed in terms of 2-D weighted Gramians is used to evaluate the weighted L_1/L_2 mixed sensitivity. A simple technique is then developed for obtaining a set of filter structures with very low weighted L_1/L_2 -sensitivity. In this connection, the optimal coordinate transformation is characterized in a closed form. Next, an iterative procedure is proposed to obtain the optimal coordinate transformation that minimizes the weighted L_2 -sensitivity measure. Once the initial value is given, the estimate at each iteration can be calculated analytically. Finally, two numerical examples are given to illustrate the utility of the proposed technique.

Index Terms—Finite word length, optimal realization, Roesser model, two-dimensional IIR digital filter, weighted coefficient sensitivity.

I. INTRODUCTION

UNDESIRABLE finite-word-length (FWL) effects arise in the fixed-point implementation of recursive digital filters. One of them is the deviation of the actual transfer function from the ideal transfer function, which is caused by the truncation or rounding of the filter coefficients. As is well known, the state-space approach allows such an effect to be minimized by appropriately choosing a filter structure that minimizes a well-defined FWL effect. Several techniques have been reported to synthesize linear state-space systems that minimize the coefficient sensitivity [1]–[7]. A similar technique for multi-input–multi-output continuous-time systems has also been presented [8]. In addition, the problem of minimizing the coefficient sensitivity of two-dimensional (2-D) state-space digital filters has been studied [9]–[15]. Based on the Roesser local state-space (LSS) model [16], Zilouchian and Carroll have investigated a coefficient sensitivity bound in 2-D state-

space digital filters [9]. Effective methods for synthesizing 2-D filter structures with minimum coefficient sensitivity have been investigated [10], [13]. In [10], all the frequency regions are treated uniformly. The method reported in [13] studies the sensitivity behavior of a transfer function within a specified frequency range. Based on the Fornasini-Marchesini second LSS model [17], similar techniques have been explored [11], [14]. The frequency-weighted sensitivity measures have been introduced in [13], [14], where a constraint on the weights of the various terms of the measure is imposed. More recently, the frequency-weighted L_2 -sensitivity problem has been considered by [15] via a 2-D gradient-flow-based optimization technique that was initiated in [7] for the one-dimensional (1-D) case. It was argued in [7] and [15] that the L_2 sensitivity minimization, although technically more challenging, is more natural and reasonable than the conventional L_1/L_2 mixed sensitivity minimization.

This paper treats the problem of reducing the coefficient sensitivity of 2-D state-space digital filters within a specified frequency range. Here, the Roesser LSS model is employed to describe 2-D state-space digital filters. From a practical viewpoint, we are interested in the sensitivity performance of the transfer function within a specified frequency range. This is achieved by defining a weighted sensitivity function. One contribution of our paper is to address the case of general unconstrained frequency weights for 2-D state-space digital filters. Another, is to solve the corresponding problem of synthesizing the filter structures with minimum weighted sensitivity. First, the sensitivities of a 2-D transfer function with respect to state-space parameters are analyzed in conjunction with frequency weighted functions. The overall frequency-weighted sensitivity measure is then evaluated, using a mixture of L_1/L_2 norms, as well as a pure L_2 norm. Second, a simple technique is developed for synthesizing the 2-D filter structures with very low frequency-weighted L_1/L_2 -sensitivity. A closed-form solution that is optimal in a certain sense is obtained. The 1-D version of this closed-form solution turns out to be more efficient than the one proposed in [5]. Notice that the closed-form solution reported in [5] is restrictive and only exists under a certain constraint. Third, an iterative procedure is presented to find the optimal coordinate transformation that minimizes the weighted L_2 -sensitivity measure. This

Manuscript received October 15, 1997; revised May 26, 1998. This paper was recommended by Associate Editor G. Martinelli.

T. Hinamoto, Y. Zempo and Y. Nishino are with the Faculty of Engineering, Hiroshima University, Higashi-Hiroshima, Japan 739-8527.

W.-S. Lu is with the Department of Electrical and Computer Engineering, University of Victoria, Victoria, BC, V8W 3P6, Canada.

Publisher Item Identifier S 1057-7122(99)08102-7.

procedure is advantageous since the estimate is calculated analytically at each iteration. Finally, two numerical examples are presented to demonstrate the validity of the proposed technique.

Throughout this paper, the n -dimensional identity matrix is denoted by I_n . The transpose (conjugate transpose) of any matrix \mathbf{A} is indicated by \mathbf{A}^t (\mathbf{A}^*) and $\text{tr}\mathbf{A}$ and \oplus are used to denote the trace of a square matrix \mathbf{A} and the direct sum of matrices, respectively.

II. WEIGHTED L_1/L_2 MIXED SENSITIVITY ANALYSIS

Consider the following LSS model $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$ for 2-D digital filters which was originally proposed by Roesser [16]:

$$\begin{bmatrix} \mathbf{x}_{11}(i, j) \\ y(i, j) \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{c} & d \end{bmatrix} \begin{bmatrix} \mathbf{x}(i, j) \\ u(i, j) \end{bmatrix} \quad (1)$$

where

$$\mathbf{x}_{11}(i, j) = \begin{bmatrix} \mathbf{x}^h(i+1, j) \\ \mathbf{x}^v(i, j+1) \end{bmatrix}, \quad \mathbf{x}(i, j) = \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix}$$

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

Here $\mathbf{x}^h(i, j)$ 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 $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3, \mathbf{A}_4, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}_1, \mathbf{c}_2, d$ are real constant matrices of appropriate dimensions. The LSS model (1) is assumed to be BIBO stable, separately locally controllable, and separately locally observable [18]. Define

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

$$H(z_1, z_2) = \mathbf{c}(\mathbf{S} - \mathbf{A})^{-1}\mathbf{b} + d, \quad \mathbf{S} = z_1\mathbf{I}_m \oplus z_2\mathbf{I}_n \quad (2)$$

where $\mathbf{F}(z_1, z_2)$ consists of the transfer functions from the filter input to the filter states, $\mathbf{G}(z_1, z_2)$ is defined as the set of transfer functions from the input of each of delay operators to the output, and $H(z_1, z_2)$ is the transfer function from the filter input to the output.

Let the coordinate transformation be specified as

$$\bar{\mathbf{x}}(i, j) = \mathbf{T}^{-1}\mathbf{x}(i, j) \quad (3)$$

where $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ and \mathbf{T}_1 (\mathbf{T}_4) is an $m \times m$ ($n \times n$) non-singular matrix. Then an algebraically equivalent realization $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m,n}$ given by

$$\bar{\mathbf{A}} = \mathbf{T}^{-1}\mathbf{A}\mathbf{T}, \bar{\mathbf{b}} = \mathbf{T}^{-1}\mathbf{b}, \bar{\mathbf{c}} = \mathbf{c}\mathbf{T} \quad (4)$$

is obtained. From (2) and (4) it is clear that the transfer function $H(z_1, z_2)$ is invariant under such a transformation.

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

$$\mathbf{S}_X = \frac{\partial f}{\partial \mathbf{X}} \quad \text{with} \quad (\mathbf{S}_X)_{ij} = \frac{\partial f}{\partial x_{ij}} \quad (5)$$

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

With these notations, it can easily be shown that

$$\frac{\partial H(z_1, z_2)}{\partial \mathbf{A}} = \mathbf{G}^t(z_1, z_2)\mathbf{F}^t(z_1, z_2)$$

$$\frac{\partial H(z_1, z_2)}{\partial \mathbf{b}} = \mathbf{G}^t(z_1, z_2)$$

$$\frac{\partial H(z_1, z_2)}{\partial \mathbf{c}^t} = \mathbf{F}(z_1, z_2). \quad (6)$$

The term d and the sensitivity with respect to it are coordinate independent and therefore they are neglected here.

To consider the sensitivity behavior of the transfer function in a specified frequency band, or even at some discrete frequency points, the weighted sensitivity functions are defined as

$$\frac{\delta H(z_1, z_2)}{\delta \mathbf{A}} = W_A(z_1, z_2) \frac{\partial H(z_1, z_2)}{\partial \mathbf{A}}$$

$$\frac{\delta H(z_1, z_2)}{\delta \mathbf{b}} = W_B(z_1, z_2) \frac{\partial H(z_1, z_2)}{\partial \mathbf{b}}$$

$$\frac{\delta H(z_1, z_2)}{\delta \mathbf{c}^t} = W_C(z_1, z_2) \frac{\partial H(z_1, z_2)}{\partial \mathbf{c}^t} \quad (7)$$

where $W_A(z_1, z_2)$, $W_B(z_1, z_2)$, and $W_C(z_1, z_2)$ are three stable, causal scalar rational functions of the complex variables z_1 and z_2 . It should be noted that δ in (7) is not meant to be a derivative operator, but rather a notation for defining the weighted parameter sensitivity as seen in (7). Let

$$W_A(z_1, z_2) = W_1(z_1, z_2)W_2(z_1, z_2) \quad (8)$$

be a factorization of $W_A(z_1, z_2)$. Note that, unlike the system considered in [13], there is no assumption such that

$$W_1(z_1, z_2) = W_B(z_1, z_2), \quad W_2(z_1, z_2) = W_C(z_1, z_2)$$

for the system considered here.

Definition 2: 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_p norm of $\mathbf{X}(z_1, z_2)$ is then defined as

$$\|\mathbf{X}\|_p = \left[\frac{1}{(2\pi j)^2} \oint \oint_{\Gamma^2} \|\mathbf{X}(z_1, z_2)\|_F^p \frac{dz_1 dz_2}{z_1 z_2} \right]^{1/p} \quad (9)$$

where $\Gamma^2 = \{(z_1, z_2): |z_1| = 1, |z_2| = 1\}$ and $\|\mathbf{X}(z_1, z_2)\|_F$ is the Frobenius norm of the matrix $\mathbf{X}(z_1, z_2)$ defined by

$$\|\mathbf{X}(z_1, z_2)\|_F = \left[\sum_{p=1}^m \sum_{q=1}^n |x_{pq}(z_1, z_2)|^2 \right]^{1/2}.$$

The overall weighted L_1/L_2 mixed sensitivity measure is now defined as

$$m_{1/2} = \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{A}} \right\|_1^2 + \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{b}} \right\|_2^2$$

$$+ \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{c}^t} \right\|_2^2. \quad (10)$$

From (6)—(8), we can write (10) as

$$m_{1/2} = \|W_1(z_1, z_2)\mathbf{G}^t(z_1, z_2)W_2(z_1, z_2)\mathbf{F}^t(z_1, z_2)\|_1^2$$

$$+ \|W_B(z_1, z_2)\mathbf{G}^t(z_1, z_2)\|_2^2$$

$$+ \|W_C(z_1, z_2)\mathbf{F}(z_1, z_2)\|_2^2. \quad (11)$$

By the Cauchy–Schwartz inequality, we have

$$\begin{aligned} & \|W_1(z_1, z_2)\mathbf{G}^t(z_1, z_2)W_2(z_1, z_2)\mathbf{F}^t(z_1, z_2)\|_1^2 \\ & \leq \|W_1(z_1, z_2)\mathbf{G}^t(z_1, z_2)\|_2^2 \|W_2(z_1, z_2)\mathbf{F}^t(z_1, z_2)\|_2^2 \end{aligned} \quad (12a)$$

where the equality sign holds if and only if

$$\begin{aligned} & |W_1(z_1, z_2)|^2 \mathbf{G}(z_1, z_2) \mathbf{G}^*(z_1, z_2) \\ & = \rho^2 |W_2(z_1, z_2)|^2 \mathbf{F}^*(z_1, z_2) \mathbf{F}(z_1, z_2) \end{aligned} \quad (12b)$$

for some nonzero real number ρ . To facilitate the mathematical treatment, an upper bound of $m_{1/2}$ is employed as follows:

$$\begin{aligned} M_{1/2} & = \|W_1(z_1, z_2)\mathbf{G}^t(z_1, z_2)\|_2^2 \|W_2(z_1, z_2)\mathbf{F}(z_1, z_2)\|_2^2 \\ & + \|W_B(z_1, z_2)\mathbf{G}^t(z_1, z_2)\|_2^2 \\ & + \|W_C(z_1, z_2)\mathbf{F}(z_1, z_2)\|_2^2 \end{aligned} \quad (13)$$

where $m_{1/2} \leq M_{1/2}$. This upper bound can be viewed as a 2-D extension of the upper bound $m_{1/2} \leq M_{1/2}$ for the 1-D case introduced by Thiele [2]. From (9) it is easy to show that

$$M_{1/2} = \text{tr} \mathbf{K}_{o1} + \text{tr} \mathbf{K}_{c2} + \text{tr} \mathbf{K}_{oB} + \text{tr} \mathbf{K}_{cC} \quad (14)$$

where \mathbf{K}_{o1} , \mathbf{K}_{c2} , \mathbf{K}_{oB} , and \mathbf{K}_{cC} are often referred to as weighted observability (for those with subindex o) and controllability (for those with subindex c) Gramians, and can be obtained by the following general expression:

$$\mathbf{K} = \frac{1}{(2\pi j)^2} \oint_{\Gamma^2} \mathbf{Y}(z_1, z_2) \mathbf{Y}^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \quad (15)$$

with $\mathbf{Y}(z_1, z_2) = W_1^*(z_1, z_2)\mathbf{G}^*(z_1, z_2)$, $W_2(z_1, z_2)\mathbf{F}(z_1, z_2)$, $W_B^*(z_1, z_2)\mathbf{G}^*(z_1, z_2)$, and $W_C(z_1, z_2)\mathbf{F}(z_1, z_2)$, respectively.

The coordinate transformation defined by (3) transforms the weighted Gramians (\mathbf{K}_{o1} , \mathbf{K}_{c2} , \mathbf{K}_{oB} , \mathbf{K}_{cC}) into ($\bar{\mathbf{K}}_{o1}$, $\bar{\mathbf{K}}_{c2}$, $\bar{\mathbf{K}}_{oB}$, $\bar{\mathbf{K}}_{cC}$). Then (14) is changed to

$$\bar{M}_{1/2} = \text{tr} \bar{\mathbf{K}}_{o1} + \text{tr} \bar{\mathbf{K}}_{c2} + \text{tr} \bar{\mathbf{K}}_{oB} + \text{tr} \bar{\mathbf{K}}_{cC} \quad (16)$$

where

$$\begin{aligned} \bar{\mathbf{K}}_{o1} & = \mathbf{T}^t \mathbf{K}_{o1} \mathbf{T}, & \bar{\mathbf{K}}_{c2} & = \mathbf{T}^{-1} \mathbf{K}_{c2} \mathbf{T}^{-t} \\ \bar{\mathbf{K}}_{oB} & = \mathbf{T}^t \mathbf{K}_{oB} \mathbf{T}, & \bar{\mathbf{K}}_{cC} & = \mathbf{T}^{-1} \mathbf{K}_{cC} \mathbf{T}^{-t}. \end{aligned}$$

Moreover, (16) is written as

$$\bar{M}_{1/2}(\mathbf{P}) = J(\mathbf{P}) + L(\mathbf{P}) \quad (17)$$

where

$$\mathbf{P} = \mathbf{T} \mathbf{T}^t, \quad \mathbf{T} = \begin{bmatrix} \mathbf{T}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_4 \end{bmatrix}$$

$$J(\mathbf{P}) = \text{tr}[\mathbf{K}_{o1} \mathbf{P}] + \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}],$$

$$L(\mathbf{P}) = \text{tr}[\mathbf{K}_{oB} \mathbf{P}] + \text{tr}[\mathbf{K}_{cC} \mathbf{P}^{-1}].$$

The problem being considered here is to obtain the symmetric and positive-definite matrix \mathbf{P} that minimizes (17), subject to the minimization of $J(\mathbf{P})$.

Remark 1: In order to effectively control the upper bound of the L_1 -norm term in (11), while minimizing (17), we seek to find a 2-D coordinate transformation \mathbf{T} (3) that minimizes (17) subject to the minimization of $J(\mathbf{P})$.

III. FILTER SYNTHESIS WITH VERY LOW WEIGHTED L_1/L_2 MIXED SENSITIVITY

In this section, we consider the problem of obtaining the matrix $\mathbf{P} = \mathbf{T} \mathbf{T}^t$ that minimizes (17), subject to the minimization of $J(\mathbf{P})$, where \mathbf{T} is block diagonal. An analytical method will be developed for obtaining such a matrix \mathbf{P} . The problem of iteratively minimizing (17) with respect to $\mathbf{P} = \mathbf{T} \mathbf{T}^t$ for any nonsingular \mathbf{T} that is not block diagonal has been solved in [5]. However, apart from whether the \mathbf{T} matrix is block diagonal or not, the two problems mentioned above are similar, yet different.

According to the partition

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_4 \end{bmatrix} \quad \mathbf{P}_i = \mathbf{T}_i \mathbf{T}_i^t, \quad i = 1, 4 \quad (18)$$

the weighted Gramians \mathbf{K}_{o1} , \mathbf{K}_{c2} , \mathbf{K}_{oB} , and \mathbf{K}_{cC} can now be represented as

$$\begin{aligned} \mathbf{K}_{o1} & = \begin{bmatrix} \mathbf{K}_{o1}^{(1)} & \mathbf{K}_{o1}^{(2)} \\ \mathbf{K}_{o1}^{(3)} & \mathbf{K}_{o1}^{(4)} \end{bmatrix} & \mathbf{K}_{c2} & = \begin{bmatrix} \mathbf{K}_{c2}^{(1)} & \mathbf{K}_{c2}^{(2)} \\ \mathbf{K}_{c2}^{(3)} & \mathbf{K}_{c2}^{(4)} \end{bmatrix} \\ \mathbf{K}_{oB} & = \begin{bmatrix} \mathbf{K}_{oB}^{(1)} & \mathbf{K}_{oB}^{(2)} \\ \mathbf{K}_{oB}^{(3)} & \mathbf{K}_{oB}^{(4)} \end{bmatrix} & \mathbf{K}_{cC} & = \begin{bmatrix} \mathbf{K}_{cC}^{(1)} & \mathbf{K}_{cC}^{(2)} \\ \mathbf{K}_{cC}^{(3)} & \mathbf{K}_{cC}^{(4)} \end{bmatrix}. \end{aligned} \quad (19a)$$

If we introduce

$$\begin{aligned} \hat{\mathbf{K}}_{o1} & = \begin{bmatrix} \mathbf{K}_{o1}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{o1}^{(4)} \end{bmatrix} & \hat{\mathbf{K}}_{c2} & = \begin{bmatrix} \mathbf{K}_{c2}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{c2}^{(4)} \end{bmatrix} \\ \hat{\mathbf{K}}_{oB} & = \begin{bmatrix} \mathbf{K}_{oB}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{oB}^{(4)} \end{bmatrix} & \hat{\mathbf{K}}_{cC} & = \begin{bmatrix} \mathbf{K}_{cC}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{cC}^{(4)} \end{bmatrix} \end{aligned} \quad (19b)$$

$J(\mathbf{P})$ and $L(\mathbf{P})$ in (17) can be expressed as

$$\begin{aligned} J(\mathbf{P}) & = \text{tr}[\hat{\mathbf{K}}_{o1} \mathbf{P}] + \text{tr}[\hat{\mathbf{K}}_{c2} \mathbf{P}^{-1}], \\ L(\mathbf{P}) & = \text{tr}[\hat{\mathbf{K}}_{oB} \mathbf{P}] + \text{tr}[\hat{\mathbf{K}}_{cC} \mathbf{P}^{-1}]. \end{aligned} \quad (20)$$

Hence it suffices to deal with the matrices $\hat{\mathbf{K}}$ instead of \mathbf{K} . To make the exposition simple, we omit the hat and write \mathbf{K} for $\hat{\mathbf{K}}$ in the following.

First, we minimize $J(\mathbf{P})$ and then minimize (17) subject to the minimization of $J(\mathbf{P})$. Using the formula for evaluating the matrix gradient [19, p. 275]

$$\begin{aligned} \frac{\partial[\text{tr}(\mathbf{M}\mathbf{X})]}{\partial \mathbf{X}} & = \mathbf{M}^t \\ \frac{\partial[\text{tr}(\mathbf{M}\mathbf{X}^{-1})]}{\partial \mathbf{X}} & = -(\mathbf{X}^{-1} \mathbf{M} \mathbf{X}^{-1})^t \end{aligned} \quad (21)$$

we obtain the equation for extrema of $J(\mathbf{P})$

$$\begin{aligned} \frac{\partial J(\mathbf{P})}{\partial \mathbf{P}} & = \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}] \mathbf{K}_{o1} \\ & - \text{tr}[\mathbf{K}_{o1} \mathbf{P}] \mathbf{P}^{-1} \mathbf{K}_{c2} \mathbf{P}^{-1} = \mathbf{0}. \end{aligned} \quad (22)$$

All the solutions of this equation take the form

$$\mathbf{P} = \rho \mathbf{P}_b \quad (23)$$

where \mathbf{P}_b is the unique solution of the equation

$$\mathbf{P} \mathbf{K}_{o1} \mathbf{P} = \mathbf{K}_{c2}$$

and ρ is an arbitrary positive number. Moreover, $J(\mathbf{P})$ has the single extremum

$$J^\circ = J(\rho \mathbf{P}_b) = (\text{tr}[\mathbf{K}_{o1} \mathbf{P}_b])^2 = (\text{tr}[\mathbf{K}_{c2} \mathbf{P}_b^{-1}])^2 \quad (24)$$

which is independent of ρ .¹ Noting that $\mathbf{PWP} = \mathbf{M}$ has the unique solution [5]

$$\mathbf{P} = \mathbf{W}^{-(1/2)} [\mathbf{W}^{1/2} \mathbf{M} \mathbf{W}^{1/2}]^{1/2} \mathbf{W}^{-(1/2)} \quad (25)$$

where $\mathbf{W} > 0$ and $\mathbf{M} \geq 0$ are symmetric, the \mathbf{P}_b matrix is given by

$$\mathbf{P}_b = (\mathbf{K}_{o1})^{-(1/2)} [(\mathbf{K}_{o1})^{1/2} \mathbf{K}_{c2} (\mathbf{K}_{o1})^{1/2}]^{1/2} (\mathbf{K}_{o1})^{-(1/2)} \quad (26)$$

and J° is described by

$$J^\circ = (\text{tr}[\mathbf{K}_{c2} \mathbf{K}_{o1}]^{1/2})^2. \quad (27)$$

Theorem 1: If $\mathbf{K}_{o1} = \mathbf{K}_{o1}^{(1)} \oplus \mathbf{K}_{o1}^{(4)}$ and $\mathbf{K}_{c2} = \mathbf{K}_{c2}^{(1)} \oplus \mathbf{K}_{c2}^{(4)}$ are $(m+n) \times (m+n)$ real symmetric positive-definite matrices, then the extremum J° in (27) is really the minimum of $J(\mathbf{P})$. Furthermore, J° can be expressed in terms of the square roots of the eigenvalues $\{\sigma_1^2, \sigma_2^2, \dots, \sigma_{m+n}^2\}$ of $\mathbf{K}_{c2} \mathbf{K}_{o1}$ as follows:

$$J^\circ = J_{\min} = \left(\sum_{i=1}^{m+n} \sigma_i \right)^2. \quad (28)$$

Proof: The proof relies on the following inequality [20, p. 556]. If \mathbf{D} is a real symmetric positive-definite matrix and if \mathbf{Q} is any nonsingular real matrix, then

$$\text{tr}[\mathbf{Q} \mathbf{D} \mathbf{Q}^t] \text{tr}[\mathbf{Q}^{-t} \mathbf{D} \mathbf{Q}^{-1}] \geq (\text{tr} \mathbf{D})^2 \quad (29)$$

where the equality sign holds if and only if

$$\rho \mathbf{Q} \mathbf{Q}^t = \mathbf{I} \quad (30)$$

for some positive real number ρ .

Choosing the above \mathbf{D} and \mathbf{Q} matrices as

$$\begin{aligned} \mathbf{D} &= [(\mathbf{K}_{o1})^{1/2} \mathbf{K}_{c2} (\mathbf{K}_{o1})^{1/2}]^{1/2} \\ \mathbf{Q} &= \mathbf{T}^t (\mathbf{K}_{o1})^{1/2} [(\mathbf{K}_{o1})^{1/2} \mathbf{K}_{c2} (\mathbf{K}_{o1})^{1/2}]^{-(1/4)} \end{aligned} \quad (31)$$

inequality (29) can be written as

$$\text{tr}[\mathbf{K}_{o1} \mathbf{P}] \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}] \geq (\text{tr}[\mathbf{K}_{c2} \mathbf{K}_{o1}]^{1/2})^2 \quad (32)$$

where $\mathbf{P} = \mathbf{T} \mathbf{T}^t$. On the other hand, taking $\mathbf{P} = \mathbf{T} \mathbf{T}^t$ and $\mathbf{P}_b = \mathbf{T}_b \mathbf{T}_b^t$ into account and using (26), it is clear that (23) is equivalent to

$$\mathbf{T} = \sqrt{\rho} \mathbf{T}_b \quad (33)$$

¹ Suppose \mathbf{P} is a solution of (22). Since $\text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}]$ and $\text{tr}[\mathbf{K}_{o1} \mathbf{P}]$ are positive, we can take $\rho > 0$ such that $\rho^2 = \text{tr}[\mathbf{K}_{o1} \mathbf{P}] / \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}]$. Then $\mathbf{P}_b = (1/\rho) \mathbf{P}$ satisfies $\mathbf{P}_b \mathbf{K}_{o1} \mathbf{P}_b = \mathbf{K}_{c2}$. Indeed, $\mathbf{P}_b^{-1} = \rho \mathbf{P}^{-1}$ and

$$\begin{aligned} 0 &= \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}] \mathbf{K}_{o1} - \text{tr}[\mathbf{K}_{o1} \mathbf{P}] \mathbf{P}^{-1} \mathbf{K}_{c2} \mathbf{P}^{-1} \\ &= \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}] (\mathbf{K}_{o1} - \rho^2 \mathbf{P}^{-1} \mathbf{K}_{c2} \mathbf{P}^{-1}) \\ &= \text{tr}[\mathbf{K}_{c2} \mathbf{P}^{-1}] (\mathbf{K}_{o1} - \mathbf{P}_b^{-1} \mathbf{K}_{c2} \mathbf{P}_b^{-1}). \end{aligned}$$

Thus, the solutions of (22) are exhausted by the solution of the form (23).

where

$$\begin{aligned} \mathbf{T}_b &= (\mathbf{K}_{o1})^{-(1/2)} [(\mathbf{K}_{o1})^{1/2} \mathbf{K}_{c2} (\mathbf{K}_{o1})^{1/2}]^{1/4} \mathbf{U} \\ \mathbf{U} &= \mathbf{U}_1 \oplus \mathbf{U}_4 \end{aligned}$$

and $\mathbf{U}_1 (\mathbf{U}_4)$ is an arbitrary $m \times m$ ($n \times n$) orthogonal matrix. Substituting (33) into (31) gives

$$\mathbf{Q} \mathbf{Q}^t = \rho \mathbf{I}_{m+n}. \quad (34)$$

This implies that equality in (32) holds, that is, the extremum J° in (27) is actually the minimum of $J(\mathbf{P})$.

Let $\sigma_1^2, \sigma_2^2, \dots, \sigma_{m+n}^2$ be the eigenvalues of $\mathbf{K}_{c2} \mathbf{K}_{o1}$. Then there exists a nonsingular matrix \mathbf{R} such that $\mathbf{K}_{c2} \mathbf{K}_{o1} = \mathbf{R}^{-1} \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_{m+n}^2) \mathbf{R}$. Hence

$$\begin{aligned} J^\circ &= J_{\min} \\ &= (\text{tr}[\mathbf{K}_{c2} \mathbf{K}_{o1}]^{1/2})^2 \\ &= (\text{tr}[\text{diag}(\sigma_1, \sigma_2, \dots, \sigma_{m+n})])^2 \\ &= \left(\sum_{i=1}^{m+n} \sigma_i \right)^2. \end{aligned} \quad (35)$$

This completes the proof of Theorem 1.

From (16), (19), and (33) we obtain

$$\rho^2 \overline{\mathbf{K}}_{c2}^{(i)} = \overline{\mathbf{K}}_{o1}^{(i)}, \quad i = 1, 4. \quad (36)$$

This shows that the weighted Gramians $\overline{\mathbf{K}}_{o1}$ and $\overline{\mathbf{K}}_{c2}$ are block balanced [21] when $\rho = 1$.

Remark 2: If (12b) can be derived from (36), then (12a) becomes an equality. However, unlike the 1-D case [2], the derivation is impossible in the 2-D case.

It turns out that the minimization of $J(\mathbf{P})$ forms a family of a matrix \mathbf{P} parameterized by $\rho > 0$. We now proceed to determine ρ that minimizes $\overline{M}_{1/2}(\mathbf{P})$ in (17).

Theorem 2: The optimal solution $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$ that minimizes the weighted sensitivity measure $\overline{M}_{1/2}(\mathbf{P})$ in (17) subject to the minimization of $J(\mathbf{P})$ is given by

$$\mathbf{P} = \sqrt{\frac{\text{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]}{\text{tr}[\mathbf{K}_{oB} \mathbf{P}_b]}} \mathbf{P}_b \quad (37a)$$

or equivalently

$$\mathbf{T} = \left(\frac{\text{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]}{\text{tr}[\mathbf{K}_{oB} \mathbf{P}_b]} \right)^{1/4} \mathbf{T}_b \mathbf{U} \quad (37b)$$

where $\mathbf{U} = \mathbf{U}_1 \oplus \mathbf{U}_2$ and $\mathbf{U}_1 (\mathbf{U}_4)$ is an arbitrary $m \times m$ ($n \times n$) orthogonal matrix. The minimum of $\overline{M}_{1/2}(\mathbf{P})$ is

$$\overline{M}_{1/2}(\mathbf{P}) = \left(\sum_{i=1}^{m+n} \sigma_i \right)^2 + 2 \sqrt{\text{tr}[\mathbf{K}_{oB} \mathbf{P}_b] \text{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]}. \quad (38)$$

Proof: Substituting (23) into (17) gives

$$\overline{M}_{1/2}(\rho \mathbf{P}_b) = J_{\min} + \rho \operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_b] + \rho^{-1} \operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]. \quad (39)$$

Here, the arithmetic-geometric inequality says that

$$\begin{aligned} & \rho \operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_b] + \rho^{-1} \operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}] \\ & \geq 2\sqrt{\operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_b] \operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]} \end{aligned} \quad (40a)$$

where the equality is valid if and only if

$$\rho \operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_b] = \rho^{-1} \operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]$$

or equivalently

$$\rho = \sqrt{\frac{\operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]}{\operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_b]}}. \quad (40b)$$

Substituting (40b) into (23) yields (37a). Substituting (28) into (39) with (40a) produces (38).

This completes the proof of Theorem 2.

The next theorem describes the relation between the second term in (38) and the minimum of $L(\mathbf{P})$.

Theorem 3: If $\mathbf{K}_{oB} = \mathbf{K}_{oB}^{(1)} \oplus \mathbf{K}_{oB}^{(4)}$ and $\mathbf{K}_{cC} = \mathbf{K}_{cC}^{(1)} \oplus \mathbf{K}_{cC}^{(4)}$ are $(m+n) \times (m+n)$ real symmetric positive-definite matrices, then

$$\sqrt{\operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_b] \operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_b^{-1}]} \geq \sum_{i=1}^{m+n} \lambda_i \quad (41)$$

where $\lambda_1, \lambda_2, \dots, \lambda_{m+n}$ are the square roots of the eigenvalues of $\mathbf{K}_{cC} \mathbf{K}_{oB}$. The equality in (41) holds if and only if the system satisfies

$$\mathbf{K}_{oB} = \alpha \mathbf{K}_{o1}, \quad \mathbf{K}_{cC} = \beta \mathbf{K}_{c2} \quad (42)$$

where α and β are some positive real numbers, $\mathbf{K}_{o1} = \mathbf{K}_{o1}^{(1)} \oplus \mathbf{K}_{o1}^{(4)}$, and $\mathbf{K}_{c2} = \mathbf{K}_{c2}^{(1)} \oplus \mathbf{K}_{c2}^{(4)}$.

Proof: To minimize $L(\mathbf{P})$ in (20), we carry out computations similar to those done in (22)–(27). The result is that $L(\mathbf{P})$ has the extremum

$$L^\circ = L(\mathbf{P}_c) = 2 \operatorname{tr}[\mathbf{K}_{oB} \mathbf{P}_c] = 2 \operatorname{tr}[\mathbf{K}_{cC} \mathbf{P}_c^{-1}] \quad (43)$$

at the matrix \mathbf{P}_c , which is the unique solution of the equation

$$\mathbf{P} \mathbf{K}_{oB} \mathbf{P} = \mathbf{K}_{cC}.$$

Since \mathbf{P}_c is solved as

$$\begin{aligned} \mathbf{P}_c &= (\mathbf{K}_{oB})^{-(1/2)} [(\mathbf{K}_{oB})^{1/2} \mathbf{K}_{cC} (\mathbf{K}_{oB})^{1/2}]^{1/2} \\ &\quad \cdot (\mathbf{K}_{oB})^{-(1/2)} \end{aligned} \quad (44)$$

we obtain

$$L^\circ = 2 \operatorname{tr}[\mathbf{K}_{cC} \mathbf{K}_{oB}]^{1/2}. \quad (45)$$

By an argument similar to those in Theorem 1, it can be shown that the extremum L° in (45) is really the minimum of $L(\mathbf{P})$ and is expressed in the form

$$L^\circ = L_{\min} = 2 \left(\sum_{i=1}^{m+n} \lambda_i \right). \quad (46)$$

Hence, the inequality (41) is proved.

(Necessity) Assume that the equality in (41) holds. Then, for a positive number ρ

$$\rho \mathbf{P}_b = \mathbf{P}_c \quad (47a)$$

must hold where \mathbf{P}_b and \mathbf{P}_c satisfy

$$\alpha \mathbf{P}_b \mathbf{K}_{o1} \mathbf{P}_b = \alpha \mathbf{K}_{c2}, \quad \mathbf{P}_c \mathbf{K}_{oB} \mathbf{P}_c = \mathbf{K}_{cC} \quad (47b)$$

and α is any positive real number. Using (47a) enables one to change (47b) to

$$\alpha \mathbf{P}_c \mathbf{K}_{o1} \mathbf{P}_c = \alpha \rho^2 \mathbf{K}_{c2}, \quad \mathbf{P}_b \mathbf{K}_{oB} \mathbf{P}_b = \frac{1}{\rho^2} \mathbf{K}_{cC} \quad (48)$$

Since \mathbf{P}_b and \mathbf{P}_c are the unique solutions, comparing (47b) with (48) concludes that

$$\mathbf{K}_{oB} = \alpha \mathbf{K}_{o1}, \quad \mathbf{K}_{cC} = \beta \mathbf{K}_{c2} \quad (49)$$

where $\beta = \alpha \rho^2$.

(Sufficiency) Assume that (42) holds. Substituting (42) into (26), we obtain

$$\mathbf{P}_b = \sqrt{\frac{\alpha}{\beta}} \mathbf{P}_c \quad (50)$$

where \mathbf{P}_c is given by (44). It is obvious that the equality in (41) holds.

This completes the proof of Theorem 3.

It should be noted that

$$\begin{aligned} \min_{\mathbf{P}} \overline{M}_{1/2}(\mathbf{P}) &\geq \min_{\mathbf{P}} J(\mathbf{P}) + \min_{\mathbf{P}} L(\mathbf{P}) \\ &= J_{\min} + L_{\min}. \end{aligned} \quad (51)$$

Corollary 1: The relation (42) holds provided

$$\begin{aligned} |W_B(z_1, z_2)| &= \sqrt{\alpha} |W_1(z_1, z_2)| \\ |W_C(z_1, z_2)| &= \sqrt{\beta} |W_2(z_1, z_2)|. \end{aligned} \quad (52)$$

Corollary 2: If (42) holds, then (37) is changed to

$$\mathbf{P} = \sqrt{\frac{\beta}{\alpha}} \mathbf{P}_b \quad (53a)$$

or equivalently

$$\mathbf{T} = \left(\frac{\beta}{\alpha} \right)^{1/4} \mathbf{T}_b \mathbf{U} \quad (53b)$$

and the equality sign in (41) holds. Moreover, (38) becomes

$$\overline{M}_{1/2}(\mathbf{P}) = \left(\sum_{i=1}^{m+n} \sigma_i \right) \left(\sum_{i=1}^{m+n} \sigma_i + 2\sqrt{\alpha\beta} \right). \quad (54)$$

The optimal filter structures that minimize $\overline{M}_{1/2}(\mathbf{P})$ (17) subject to the minimization of $J(\mathbf{P})$ can readily be synthesized by substituting (37b) into (4).

Remark 3: Notice that (53) can be considered to be an extension of the 1-D closed-form solution reported in [5] to the 2-D case. In Corollary 2 it is mentioned that (53) can be derived from (37) as a special case for the system such that (42) is satisfied. In other words, unlike the solution given by (53), (37) can be applied to the general systems where (42) is not always satisfied. It should be pointed out that neither the 1-D version of the closed-form solution (37) stated in Theorem 2 nor the 1-D counterpart of arguments stated in Theorem 3 has been reported in [5].

IV. FILTER SYNTHESIS WITH MINIMUM WEIGHTED L_2 -SENSITIVITY

In this section, we synthesize the 2-D filter structures that minimize a weighted L_2 -sensitivity measure defined by

$$m_2 = \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{A}} \right\|_2^2 + \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{b}} \right\|_2^2 + \left\| \frac{\delta H(z_1, z_2)}{\delta \mathbf{c}^t} \right\|_2^2 \quad (55)$$

instead of (10). Referring to (11) and (14), we can write (55) as

$$\begin{aligned} m_2 &= \text{tr}[\mathbf{K}_A] + \text{tr}[\mathbf{K}_{oB}] + \text{tr}[\mathbf{K}_{cC}] \\ &= \text{tr} \left[\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{M}_A(i, j) \mathbf{M}_A^t(i, j) \right] + \text{tr}[\mathbf{K}_{oB}] \\ &\quad + \text{tr}[\mathbf{K}_{cC}] \end{aligned} \quad (56)$$

where \mathbf{K}_A is obtained by the general expression of (15) with

$$\mathbf{Y}(z_1, z_2) = W_A(z_1, z_2) \mathbf{G}^t(z_1, z_2) \mathbf{F}^t(z_1, z_2)$$

and $\mathbf{M}_A(i, j)$ is derived from

$$\mathbf{M}_A(i, j) = \sum_{(0,0) \leq (k,r) < (i,j)} w_A(k, r) \mathbf{M}(i-k, j-r)$$

$$\mathbf{M}(i, j) = \sum_{(0,0) \leq (k,r) < (i,j)} \mathbf{g}^t(k, r) \mathbf{f}^t(i-k, j-r).$$

$$\mathbf{f}(i, j) = \mathbf{A}^{(i-1, j)} \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{b} + \mathbf{A}^{(i, j-1)} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \mathbf{b}$$

$$\mathbf{g}(i, j) = \mathbf{c} \mathbf{A}^{(i-1, j)} \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \mathbf{c} \mathbf{A}^{(i, j-1)} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}$$

$$\mathbf{A}^{(1,0)} = \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{A},$$

$$\mathbf{A}^{(0,1)} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \mathbf{A}$$

$$\mathbf{A}^{(0,0)} = \mathbf{I}_{m+n}, \quad \mathbf{A}^{(-i, j)} = \mathbf{A}^{(i, -j)} = \mathbf{0}, \quad (i, j \geq 1)$$

$$\begin{aligned} \mathbf{A}^{(i, j)} &= \mathbf{A}^{(1,0)} \mathbf{A}^{(i-1, j)} + \mathbf{A}^{(0,1)} \mathbf{A}^{(i, j-1)} \\ &= \mathbf{A}^{(i-1, j)} \mathbf{A}^{(1,0)} + \mathbf{A}^{(i, j-1)} \mathbf{A}^{(0,1)}, \quad (i, j) > (0, 0) \end{aligned}$$

with $w_A(k, r)$ being the unit-sample response of $W_A(z_1, z_2)$. Applying the coordinate transformation defined by (3) to the original filter, (56) becomes

$$\begin{aligned} \bar{m}_2(\mathbf{P}) &= \text{tr} \left[\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{M}_A(i, j) \mathbf{P}^{-1} \mathbf{M}_A^t(i, j) \mathbf{P} \right] \\ &\quad + \text{tr}[\mathbf{K}_{oB} \mathbf{P}] + \text{tr}[\mathbf{K}_{cC} \mathbf{P}^{-1}] \end{aligned} \quad (57)$$

where \mathbf{P} is as defined in (17). According to the partition of (18), $\mathbf{M}_A(i, j)$ can be represented as

$$\mathbf{M}_A(i, j) = \begin{bmatrix} \mathbf{M}_A^{(1)}(i, j) & \mathbf{M}_A^{(2)}(i, j) \\ \mathbf{M}_A^{(3)}(i, j) & \mathbf{M}_A^{(4)}(i, j) \end{bmatrix}. \quad (58)$$

Taking (21) into account, it follows from (57) that

$$\begin{aligned} \frac{\partial \bar{m}_2(\mathbf{P})}{\partial \mathbf{P}_1} &= \mathbf{F}_1(\mathbf{P}) - \mathbf{P}_1^{-1} \mathbf{F}_2(\mathbf{P}) \mathbf{P}_1^{-1} \\ \frac{\partial \bar{m}_2(\mathbf{P})}{\partial \mathbf{P}_4} &= \mathbf{F}_3(\mathbf{P}) - \mathbf{P}_4^{-1} \mathbf{F}_4(\mathbf{P}) \mathbf{P}_4^{-1} \end{aligned} \quad (59)$$

where

$$\begin{aligned} \mathbf{F}_1(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\mathbf{M}_A^{(1)}(i, j) \mathbf{P}_1^{-1} \mathbf{M}_A^{(1)}(i, j)^t \\ &\quad + \mathbf{M}_A^{(2)}(i, j) \mathbf{P}_4^{-1} \mathbf{M}_A^{(2)}(i, j)^t] + \mathbf{K}_{oB}^{(1)} \\ \mathbf{F}_2(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\mathbf{M}_A^{(1)}(i, j)^t \mathbf{P}_1 \mathbf{M}_A^{(1)}(i, j) \\ &\quad + \mathbf{M}_A^{(3)}(i, j)^t \mathbf{P}_4 \mathbf{M}_A^{(3)}(i, j)] + \mathbf{K}_{cC}^{(1)} \\ \mathbf{F}_3(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\mathbf{M}_A^{(3)}(i, j) \mathbf{P}_1^{-1} \mathbf{M}_A^{(3)}(i, j)^t \\ &\quad + \mathbf{M}_A^{(4)}(i, j) \mathbf{P}_4^{-1} \mathbf{M}_A^{(4)}(i, j)^t] + \mathbf{K}_{oB}^{(4)} \\ \mathbf{F}_4(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\mathbf{M}_A^{(2)}(i, j)^t \mathbf{P}_1 \mathbf{M}_A^{(2)}(i, j) \\ &\quad + \mathbf{M}_A^{(4)}(i, j)^t \mathbf{P}_4 \mathbf{M}_A^{(4)}(i, j)] + \mathbf{K}_{cC}^{(4)}. \end{aligned}$$

Letting the two equations in (59) be null yields

$$\begin{aligned} \mathbf{P}_1 \mathbf{F}_1(\mathbf{P}) \mathbf{P}_1 &= \mathbf{F}_2(\mathbf{P}) \\ \mathbf{P}_4 \mathbf{F}_3(\mathbf{P}) \mathbf{P}_4 &= \mathbf{F}_4(\mathbf{P}) \end{aligned} \quad (60)$$

respectively. From (25) and (60) it follows that the values $\mathbf{P}_1^{(i+1)}$ and $\mathbf{P}_4^{(i+1)}$ satisfying

$$\begin{aligned} \mathbf{P}_1^{(i+1)} \mathbf{F}_1(\mathbf{P}^{(i)}) \mathbf{P}_1^{(i+1)} &= \mathbf{F}_2(\mathbf{P}^{(i)}) \\ \mathbf{P}_4^{(i+1)} \mathbf{F}_3(\mathbf{P}^{(i)}) \mathbf{P}_4^{(i+1)} &= \mathbf{F}_4(\mathbf{P}^{(i)}) \end{aligned} \quad (61)$$

respectively, are given by

$$\begin{aligned} \mathbf{P}_1^{(i+1)} &= \mathbf{F}_1^{-1(1/2)}(\mathbf{P}^{(i)}) [\mathbf{F}_1^{1/2}(\mathbf{P}^{(i)}) \mathbf{F}_2(\mathbf{P}^{(i)}) \mathbf{F}_1^{1/2}(\mathbf{P}^{(i)})]^{1/2} \\ &\quad \cdot \mathbf{F}_1^{-1(1/2)}(\mathbf{P}^{(i)}) \\ \mathbf{P}_4^{(i+1)} &= \mathbf{F}_3^{-1(1/2)}(\mathbf{P}^{(i)}) [\mathbf{F}_3^{1/2}(\mathbf{P}^{(i)}) \mathbf{F}_4(\mathbf{P}^{(i)}) \mathbf{F}_3^{1/2}(\mathbf{P}^{(i)})]^{1/2} \\ &\quad \cdot \mathbf{F}_3^{-1(1/2)}(\mathbf{P}^{(i)}) \end{aligned} \quad (62)$$

where $\mathbf{P}^{(i)}$ is the solution of the previous iteration. The initial estimate $\mathbf{P}^{(0)}$ in the above iteration is given by (37a). This iteration process continues until

$$|\bar{m}_2(\mathbf{P}^{(i+1)}) - \bar{m}_2(\mathbf{P}^{(i)})| < \varepsilon \quad (63)$$

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

While the convergence of the iterative algorithm described in (62) remains to be proved, the algorithm was applied to quite a number of simulation examples and fast convergence was observed in all the cases. A sample of these examples will be

illustrated in the next section. As a remark on this convergence issue, we note that an interesting iterative algorithm, based on the concept of gradient flow for frequency weighted sensitivity minimization of 1-D discrete-time systems, was proposed in [7] and extended to 2-D Fornasini–Marchesini model in [15]. Although the nonlinear setting in (62) differs from that of [7] and [15], the technique employed there to show the convergence of the algorithms appears worthwhile to analyze in order to show the convergence of (62) or similar algorithms.

Given the L_2 -optimal matrix $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$ which is positive-definite and symmetric, the corresponding L_2 -optimal transformation matrix can be constructed as

$$\mathbf{T} = [\mathbf{P}_1^{1/2} \oplus \mathbf{P}_4^{1/2}][\mathbf{U}_1 \oplus \mathbf{U}_4] \quad (64)$$

where \mathbf{U}_1 and \mathbf{U}_4 are arbitrary $m \times m$ and $n \times n$ orthogonal matrices, respectively. It is possible to synthesize the L_2 -optimal filter structures such that (57) is minimum by substituting (64) into (4).

Remark 4: As was shown in [22], the orthogonal matrices \mathbf{U}_1 and \mathbf{U}_4 in (64) can be used to obtain a state-space realization with more zero or one entries, which further reduces the L_2 sensitivity. An alternative approach to accomplish this is to use singular value decomposition (SVD) [23], [24] as follows.

Let us denote

$$\tilde{\mathbf{A}} = \tilde{\mathbf{T}}^{-1} \mathbf{A} \tilde{\mathbf{T}} = \begin{bmatrix} \tilde{\mathbf{A}}_1 & \tilde{\mathbf{A}}_2 \\ \tilde{\mathbf{A}}_3 & \tilde{\mathbf{A}}_4 \end{bmatrix}, \quad \tilde{\mathbf{T}} = \mathbf{P}_1^{1/2} \oplus \mathbf{P}_4^{1/2} \quad (65)$$

and apply SVD to $\tilde{\mathbf{A}}_2$

$$\tilde{\mathbf{A}}_2 = \mathbf{R} \mathbf{S} \mathbf{Q}^T \quad (66)$$

where \mathbf{R} and \mathbf{Q} are $m \times m$ and $n \times n$ orthogonal matrices, respectively, and

$$\mathbf{S}_2 = \left[\begin{array}{ccc|c} \sigma_1 & & & \mathbf{0} \\ & \ddots & & \\ & & \sigma_{r_2} & \\ \hline \mathbf{0} & & & \mathbf{0} \end{array} \right]$$

with r_2 being the rank of $\tilde{\mathbf{A}}_2$. Evidently, if we let $\mathbf{U}_1 = \mathbf{R}$, $\mathbf{U}_2 = \mathbf{Q}$, then $\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}$ has the form

$$\bar{\mathbf{A}} = \begin{bmatrix} \bar{\mathbf{A}}_1 & \mathbf{S}_2 \\ \bar{\mathbf{A}}_3 & \bar{\mathbf{A}}_4 \end{bmatrix} \quad (67)$$

where \mathbf{S}_2 has $(mn - r_2)$ zero entries. Alternatively, SVD may be applied to the matrix $\tilde{\mathbf{A}}_3$ to yield $(mn - r_3)$ zero entries where r_3 is the rank of $\tilde{\mathbf{A}}_3$.

V. ILLUSTRATIVE EXAMPLES

The frequency weighting functions $W_A(z_1, z_2)$, $W_B(z_1, z_2)$, and $W_C(z_1, z_2)$ can be either of 2-D finite impulse response (FIR) or infinite impulse response (IIR) digital filters. For simplicity, let these be given by the

following 2-D lowpass filters:

$$W_A(z_1, z_2) = \sum_{i=0}^{20} \sum_{j=0}^{20} w_a(i, j) z_1^{-i} z_2^{-j}$$

$$W_B(z_1, z_2) = \sum_{i=0}^{20} \sum_{j=0}^{20} w_b(i, j) z_1^{-i} z_2^{-j}$$

$$W_C(z_1, z_2) = \frac{N(z_1, z_2)}{D_1(z_1) D_2(z_2)}$$

where

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

$$w_b(i, j) = 0.256322 \exp[-0.103203\{(i-5)^2 + (j-i)^2\}]$$

$$N(z_1, z_2) = \sum_{i=0}^4 \sum_{j=0}^4 b_{ij} z_1^{-i} z_2^{-j}$$

$$D(z_k) = 1.0 - 1.11425z_k^{-1} + 0.75745z_k^{-2} - 0.34255z_k^{-3} + 0.10171z_k^{-4}, \quad k = 1, 2$$

$$\begin{bmatrix} b_{00} & b_{01} & \cdots & b_{04} \\ b_{10} & b_{11} & \cdots & b_{14} \\ \vdots & \vdots & \ddots & \vdots \\ b_{40} & b_{41} & \cdots & b_{44} \end{bmatrix} = \begin{bmatrix} 0.12814 & 0.64232 & 0.74979 & 0.64232 & 0.12814 \\ 0.64232 & 0.33077 & 0.68889 & 0.33077 & 0.64232 \\ 0.74979 & 0.68889 & 1.34339 & 0.68889 & 0.74979 \\ 0.64232 & 0.33077 & 0.68889 & 0.33077 & 0.64232 \\ 0.12814 & 0.64232 & 0.74979 & 0.64232 & 0.12814 \end{bmatrix} \times 10^{-2}.$$

A factorization of (8) is now assumed to be

$$W_1(z_1, z_2) = 1$$

$$W_2(z_1, z_2) = W_A(z_1, z_2)$$

Example 1: (2, 2)th-Order Filter

Consider the LSS model (1) specified by

$$\mathbf{A} = \begin{bmatrix} 1.88899 & -0.91219 & -1.0 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 \\ 0.02771 & -0.02580 & 1.88899 & 1.0 \\ -0.02580 & 0.02431 & -0.91219 & 0.0 \end{bmatrix}$$

$$\mathbf{b} = [0.219089 \quad 0.0 \quad -0.028889 \quad 0.091219]^t$$

$$\mathbf{c} = [0.28889 \quad -0.091219 \quad -0.219089 \quad 0.0]$$

where $m = n = 2$.

Applying Parseval's relation to (14) and (15), it follows that

$$\mathbf{K}_{o1} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{g}_{o1}^t(i, j) \mathbf{g}_{o1}(i, j)$$

$$\mathbf{K}_{c2} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}_{c2}(i, j) \mathbf{f}_{c2}^t(i, j) \quad (68)$$

where

$$g_{o1}(i, j) = \sum_{(0,0) \leq (k,r) < (i,j)} w_1(k, r) g(i - k, j - r)$$

$$f_{c2}(i, j) = \sum_{(0,0) \leq (k,r) < (i,j)} w_2(k, r) f(i - k, j - r)$$

$$W_k(z_1, z_2) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} w_k(i, j) z_1^{-i} z_2^{-j}, \quad k = 1, 2.$$

Here, $f(i, j)$ and $g(i, j)$ are as defined in (56). From (56) and (68), the 2-D weighted Gramians $K_A^{(i)}$, $K_{o1}^{(i)}$, $K_{c2}^{(i)}$, $K_{oB}^{(i)}$, and $K_{cC}^{(i)}$, $i = 1, 4$ can be calculated by truncation $0 \leq i \leq 100$ and $0 \leq j \leq 100^2$ as

$$K_A^{(1)} = \begin{bmatrix} 460405.02074980 & -406380.17548504 \\ -406380.17548504 & 364565.10231906 \end{bmatrix}$$

$$K_A^{(4)} = \begin{bmatrix} 390597141.53364521 & 383598953.44964486 \\ 383598953.44964486 & 390423379.40672630 \end{bmatrix}$$

$$K_{o1}^{(1)} = \begin{bmatrix} 11.33639752 & -10.35083304 \\ -10.35083304 & 9.66068618 \end{bmatrix}$$

$$K_{o1}^{(4)} = \begin{bmatrix} 638.95921778 & 622.90731381 \\ 622.90731381 & 638.94448493 \end{bmatrix}$$

$$K_{c2}^{(1)} = \begin{bmatrix} 2869.13457721 & 2822.98523270 \\ 2822.98523270 & 2869.02675750 \end{bmatrix}$$

$$K_{c2}^{(4)} = \begin{bmatrix} 4.92296703 & -4.33804824 \\ -4.33804824 & 3.94010274 \end{bmatrix}$$

$$K_{oB}^{(1)} = \begin{bmatrix} 23.29761522 & -20.35185194 \\ -20.35185194 & 18.37114758 \end{bmatrix}$$

$$K_{oB}^{(4)} = \begin{bmatrix} 21436.38888158 & 21164.41125664 \\ 21164.41125664 & 21436.07017164 \end{bmatrix}$$

$$K_{cC}^{(1)} = \begin{bmatrix} 50.99062638 & 50.04288786 \\ 50.04288786 & 50.98841828 \end{bmatrix}$$

$$K_{cC}^{(4)} = \begin{bmatrix} 0.12967208 & -0.11585595 \\ -0.11585595 & 0.10657265 \end{bmatrix}.$$

Then, the original weighted sensitivities (14) and (56) become

$$M_{1/2} = J + L = 7507830.8616$$

$$m_2 = 781888507.4065$$

respectively, where $J = 7464814.5185$ and $L = 43016.3431$.

²This region was chosen according to the memory capacity of computers in the laboratory as well as the approximation accuracy in the truncation.

The suboptimal $P_c = P_1 \oplus P_4$ matrix that minimizes $L(P)$ in (17) is calculated from (44) as

$$P_1 = \begin{bmatrix} 7.84351089 & 8.18028539 \\ 8.18028539 & 8.86708460 \end{bmatrix}$$

$$P_4 = \begin{bmatrix} 0.01483796 & -0.01394208 \\ -0.01394208 & 0.01349485 \end{bmatrix}$$

or equivalently

$$T_1 = \begin{bmatrix} 2.06433731 & 1.89262314 \\ 1.76150446 & 2.40087206 \end{bmatrix}$$

$$T_4 = \begin{bmatrix} 0.08854411 & -0.08365345 \\ -0.06955674 & 0.09304145 \end{bmatrix}.$$

In this case, from (17) and (57) we have

$$\bar{M}_{1/2} = J + L_{\min} = 8057.9652$$

$$\bar{m}_2 = 980226506$$

respectively, where $J = 7998.2396$ and $L_{\min} = 59.7256$.

The L_1/L_2 -optimal P matrix which minimizes (17) subject to the minimization of $J(P)$ can be computed from (37a) as

$$P_1 = \begin{bmatrix} 5.14975986 & 5.30489866 \\ 5.30489866 & 5.62128787 \end{bmatrix}$$

$$P_4 = \begin{bmatrix} 0.01884643 & -0.01733058 \\ -0.01733058 & 0.01662501 \end{bmatrix}$$

or equivalently

$$T_1 = \begin{bmatrix} 1.41877744 & 1.77110995 \\ 1.15269142 & 2.07185674 \end{bmatrix}$$

$$T_4 = \begin{bmatrix} 0.12303452 & -0.06090101 \\ -0.10149961 & 0.07951626 \end{bmatrix}.$$

As a result, the L_1/L_2 -optimal filter structure is synthesized from (4) as shown at the bottom of this page and it follows from (17) and (57) that

$$\bar{M}_{1/2} = J_{\min} + L = 6459.9130$$

$$\bar{m}_2 = 163012.2215$$

respectively, where $J_{\min} = 6391.6953$ and $L = 68.2177$.

Applying the iterative procedure (62) produces

$$P_1 = \begin{bmatrix} 12.78717263 & 13.55220273 \\ 13.55220273 & 14.75846399 \end{bmatrix}$$

$$P_4 = \begin{bmatrix} 0.00869668 & -0.00794560 \\ -0.00794560 & 0.00751555 \end{bmatrix}$$

$$\bar{A} = \begin{bmatrix} 0.95926221 & -0.13460051 & -0.28387657 & 0.14051640 \\ 0.15109352 & 0.92972779 & 0.15793663 & -0.07817724 \\ 0.06626498 & -0.01761805 & 0.99244647 & 0.15503995 \\ -0.02334944 & 0.03626901 & -0.14459846 & 0.89654353 \end{bmatrix}$$

$$\bar{b} = [0.50550229 \quad -0.28123959 \quad 0.90459477 \quad 2.30185664]^t$$

$$\bar{c} = [0.30472326 \quad 0.32266325 \quad -0.02695551 \quad 0.01334274]$$

after 20 iterations or equivalently

$$\mathbf{T}_1 = \begin{bmatrix} 2.21652433 & 2.80609916 \\ 1.85565665 & 3.36377799 \end{bmatrix}$$

$$\mathbf{T}_4 = \begin{bmatrix} 0.07803397 & -0.05106247 \\ -0.06253108 & 0.06004513 \end{bmatrix}.$$

where truncation $0 \leq i \leq 100$ and $0 \leq j \leq 100$ was used to compute $F_k(\mathbf{P})$, $k = 1, 2, 3, 4$. Substituting the above coordinate transformation into (4) provides the L_2 -optimal filter structure as shown at the bottom of this page and (17) and (57) were used to calculate

$$\overline{M}_{1/2} = J + L = 11546.8608$$

$$\overline{m}_2 = 84193.9719$$

respectively, where $J = 11481.6550$ and $L = 65.2058$.

Example 2: (3, 3)th Order Filter

Let the LSS model (1) be given by

$$\mathbf{A}_1 = \begin{bmatrix} 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \\ 0.38315 & -1.38605 & 1.90670 \end{bmatrix}$$

$$\mathbf{A}_2 = \begin{bmatrix} -0.06280 & 0.06190 & 0.00654 \\ -0.02810 & 0.03956 & -0.02248 \\ 1.24452 & -0.57092 & 2.05865 \end{bmatrix}$$

$$\mathbf{A}_3 = \begin{bmatrix} -0.00003 & 0.00038 & -0.00053 \\ -0.00001 & 0.00018 & -0.00026 \\ -0.00008 & 0.00023 & -0.00017 \end{bmatrix}$$

$$\mathbf{A}_4 = \begin{bmatrix} 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \\ 0.38238 & -1.38178 & 1.90253 \end{bmatrix}$$

$$\mathbf{b}_1 = \mathbf{b}_2 = [0.0 \ 0.0 \ 1.0]^t$$

$$\mathbf{c}_1 = [0.01141 \ -0.00540 \ 0.01956]$$

$$\mathbf{c}_2 = [0.01164 \ -0.00545 \ 0.01960]$$

where $m = n = 3$.

Using (56) and (68), the submatrices of the 2-D weighted Gramians $\mathbf{K}_A^{(i)}$, $\mathbf{K}_{o1}^{(i)}$, $\mathbf{K}_{c2}^{(i)}$, $\mathbf{K}_{oB}^{(i)}$, and $\mathbf{K}_{cC}^{(i)}$, $i = 1, 4$ can be calculated by truncation $0 \leq i \leq 100$ and $0 \leq j \leq 100$ as shown at the bottom of the next page. Then (14) and (56) becomes

$$M_{1/2} = J + L = 299673.7313$$

$$m_2 = 10441330.7603$$

respectively, where $J = 296773.3396$ and $L = 2900.3916$.

Applying (44), the suboptimal filter structures are realized from

$$\mathbf{P}_1 = \begin{bmatrix} 1994.65983573 & 1090.47959350 & 349.84385826 \\ 1090.47959350 & 810.29935629 & 464.59809968 \\ 349.84385826 & 464.59809968 & 447.86957497 \end{bmatrix}$$

$$\mathbf{P}_4 = \begin{bmatrix} 6.11907288 & 3.09433408 & 0.66823909 \\ 3.09433408 & 2.08175279 & 0.97490842 \\ 0.66823909 & 0.97490842 & 0.93506167 \end{bmatrix}$$

or equivalently,

$$\mathbf{T}_1 = \begin{bmatrix} 13.82537475 & 30.61643239 & 29.43047598 \\ 16.78176757 & 22.16623497 & 6.10980039 \\ 18.56399712 & 8.45750822 & -5.63188606 \end{bmatrix}$$

$$\mathbf{T}_4 = \begin{bmatrix} 0.59412325 & 1.56120892 & 1.82447723 \\ 0.72216789 & 1.15695186 & 0.47083831 \\ 0.79947389 & 0.46108297 & -0.28862721 \end{bmatrix}.$$

From (17) and (57), this gives

$$\overline{M}_{1/2} = J + L_{\min} = 364.6126$$

$$\overline{m}_2 = 2186.2190$$

respectively, where $J = 351.6035$ and $L_{\min} = 13.0091$.

Making use of (37a), the L_1/L_2 -optimal filter structures that minimize (17) subject to the minimization of $J(\mathbf{P})$ are constructed from

$$\mathbf{P}_1 = \begin{bmatrix} 695.06511011 & 443.59300157 & 226.86622553 \\ 443.59300157 & 382.35826949 & 282.87830911 \\ 226.86622553 & 282.87830911 & 287.77432092 \end{bmatrix}$$

$$\mathbf{P}_4 = \begin{bmatrix} 2.71258093 & 1.73686609 & 0.89019654 \\ 1.73686609 & 1.50929838 & 1.12229674 \\ 0.89019654 & 1.12229674 & 1.14720037 \end{bmatrix}$$

or equivalently

$$\mathbf{T}_1 = \begin{bmatrix} 13.73007262 & 19.38689860 & 11.43233915 \\ 15.36067722 & 12.09880780 & -0.16344619 \\ 16.16489972 & 2.79936292 & -4.31670073 \end{bmatrix}$$

$$\mathbf{T}_4 = \begin{bmatrix} 0.86378782 & 1.21374840 & 0.70232925 \\ 0.96989606 & 0.75372228 & -0.02242183 \\ 1.02128137 & 0.16658800 & -0.27646548 \end{bmatrix}.$$

$$\overline{\mathbf{A}} = \begin{bmatrix} 0.96517195 & -0.16243289 & -0.11672716 & 0.07638182 \\ 0.12649366 & 0.92381805 & 0.06439353 & -0.04213668 \\ 0.13175457 & -0.04244315 & 0.97923151 & 0.12870234 \\ -0.06389460 & 0.11195028 & -0.16569878 & 0.90975849 \end{bmatrix}$$

$$\overline{\mathbf{b}} = [0.32772442 \ -0.18079196 \ 1.95851839 \ 3.55877769]^t$$

$$\overline{\mathbf{c}} = [0.47106057 \ 0.50381352 \ -0.01709638 \ 0.01118722]$$

Then, making use of (4), we get

$$\begin{aligned} \bar{A}_1 &= \begin{bmatrix} 0.83377989 & -0.21625751 & -0.01783948 \\ 0.27580984 & 0.51201201 & -0.32656076 \\ -0.12545803 & 0.44975172 & 0.56090810 \end{bmatrix} \\ \bar{A}_2 &= \begin{bmatrix} 0.30544369 & 0.16449039 & 0.03618117 \\ -0.38465345 & -0.20745137 & -0.04669960 \\ 0.28655056 & 0.15175330 & 0.03160229 \end{bmatrix} \\ \bar{A}_3 &= \begin{bmatrix} -0.00034161 & 0.00120968 & -0.00058481 \\ -0.00168438 & 0.00016443 & 0.00204681 \\ -0.00114298 & 0.00183361 & -0.00013724 \end{bmatrix} \\ \bar{A}_4 &= \begin{bmatrix} 0.83476698 & -0.21581510 & -0.01797649 \\ 0.27709299 & 0.51222077 & -0.32714155 \\ -0.12456645 & 0.45339646 & 0.55554225 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \bar{b}_1 &= [0.11589425 \quad -0.14568253 \quad 0.10786033]^t \\ \bar{b}_2 &= [1.81737136 \quad -2.28750015 \quad 1.71803526]^t \\ \bar{c}_1 &= [0.38989791 \quad 0.21062649 \quad 0.04689093] \\ \bar{c}_2 &= [0.02478567 \quad 0.01328537 \quad 0.00287859] \end{aligned}$$

and from (17) and (57) it follows that

$$\bar{M}_{1/2} = J_{\min} + L = 290.7477$$

$$\bar{m}_2 = 2806.3110$$

respectively, where $J_{\min} = 276.0360$ and $L = 14.7116$.

Applying the iterative procedure (62) provides

$$P_1 = \begin{bmatrix} 2402.69261014 & 1326.43635801 & 446.26772332 \\ 1326.43635801 & 927.17620872 & 512.42934385 \\ 446.26772332 & 512.42934385 & 480.37920997 \end{bmatrix}$$

$$P_4 = \begin{bmatrix} 5.13089004 & 2.66738136 & 0.67461188 \\ 2.66738136 & 1.78259209 & 0.86291979 \\ 0.67461188 & 0.86291979 & 0.81547874 \end{bmatrix}$$

after 20 iterations or equivalently

$$T_1 = \begin{bmatrix} 17.53074796 & 34.34813573 & 30.25840475 \\ 19.02504678 & 22.70962425 & 7.03539409 \\ 19.82782611 & 7.62476990 & -5.39438650 \end{bmatrix}$$

$$T_4 = \begin{bmatrix} 0.65589244 & 1.51045141 & 1.55538795 \\ 0.73604761 & 1.04267141 & 0.39199789 \\ 0.78027301 & 0.37396246 & -0.25846633 \end{bmatrix}.$$

where truncation $0 \leq i \leq 100$ and $0 \leq j \leq 100$ was used to compute $F_k(P)$, $k = 1, 2, 3, 4$. Substituting the above

$$K_A^{(1)} = \begin{bmatrix} 1251.75449772 & -2979.86329982 & 2999.81755684 \\ -2979.86329982 & 7122.30587652 & -7133.50326166 \\ 2999.81755684 & -7133.50326166 & 7299.95917720 \end{bmatrix}$$

$$K_A^{(4)} = \begin{bmatrix} 845585.48203204 & -1993349.60116238 & 2015318.08545784 \\ -1993349.60116238 & 4714069.08436384 & -4744374.88497449 \\ 2015318.08545784 & -4744374.88497449 & 4863101.78275287 \end{bmatrix}$$

$$K_{o1}^{(1)} = \begin{bmatrix} 0.00149858 & -0.00364930 & 0.00356937 \\ -0.00364930 & 0.00956165 & -0.00905795 \\ 0.00356937 & -0.00905795 & 0.00933207 \end{bmatrix}$$

$$K_{o1}^{(4)} = \begin{bmatrix} 0.15953189 & -0.38560521 & 0.37922176 \\ -0.38560521 & 1.00770155 & -0.95804046 \\ 0.37922176 & -0.95804046 & 0.99147667 \end{bmatrix}$$

$$K_{c2}^{(1)} = \begin{bmatrix} 45322.52827960 & 43822.10189449 & 39634.12075737 \\ 43822.10189449 & 45307.52406750 & 43818.67539592 \\ 39634.12075737 & 43818.67539592 & 45327.66655788 \end{bmatrix}$$

$$K_{c2}^{(4)} = \begin{bmatrix} 77.13949821 & 74.86433871 & 67.80484994 \\ 74.86433871 & 77.73814927 & 75.31371509 \\ 67.80484994 & 75.31371509 & 78.03763354 \end{bmatrix}$$

$$K_{oB}^{(1)} = \begin{bmatrix} 0.00657555 & -0.01610628 & 0.01597188 \\ -0.01610628 & 0.04000610 & -0.03922995 \\ 0.01597188 & -0.03922995 & 0.04019986 \end{bmatrix}$$

$$K_{oB}^{(4)} = \begin{bmatrix} 4.24697220 & -10.09219961 & 10.16484848 \\ -10.09219961 & 24.09804943 & -24.12544059 \\ 10.16484848 & -24.12544059 & 24.73648669 \end{bmatrix}$$

$$K_{cC}^{(1)} = \begin{bmatrix} 947.12826344 & 900.34470992 & 772.59639177 \\ 900.34470992 & 946.91813917 & 900.22072602 \\ 772.59639177 & 900.22072602 & 947.21153262 \end{bmatrix}$$

$$K_{cC}^{(4)} = \begin{bmatrix} 1.97861983 & 1.88693064 & 1.62312130 \\ 1.88693064 & 1.99044115 & 1.89578584 \\ 1.62312130 & 1.89578584 & 1.99633146 \end{bmatrix}.$$

TABLE I
WEIGHTED SENSITIVITY ANALYSIS

Realization	Example 1	Example 2
Original	$M_{1/2} = 7.5078 \times 10^6$ $m_2 = 7.8189 \times 10^8$	$M_{1/2} = 2.9967 \times 10^5$ $m_2 = 1.0441 \times 10^7$
Suboptimal	$\bar{M}_{1/2} = 8.0580 \times 10^3$ $\bar{m}_2 = 9.8023 \times 10^4$	$\bar{M}_{1/2} = 3.6461 \times 10^2$ $\bar{m}_2 = 2.1862 \times 10^3$
L_1/L_2 -Optimal	$\bar{M}_{1/2} = 6.4599 \times 10^3$ $\bar{m}_2 = 1.6301 \times 10^5$	$\bar{M}_{1/2} = 2.9075 \times 10^2$ $\bar{m}_2 = 2.8063 \times 10^3$
L_2 -Optimal	$\bar{M}_{1/2} = 1.1547 \times 10^4$ $\bar{m}_2 = 8.4194 \times 10^4$	$\bar{M}_{1/2} = 3.8044 \times 10^2$ $\bar{m}_2 = 2.1519 \times 10^3$
$(J_{\min} + L_{\min})$	6.4514×10^3	2.8905×10^2

coordinate transformation into (4) results in

$$\bar{\mathbf{A}}_1 = \begin{bmatrix} 0.79612923 & -0.21332680 & -0.06943544 \\ 0.23792156 & 0.37583635 & -0.40698693 \\ -0.10257831 & 0.44748310 & 0.73473442 \end{bmatrix}$$

$$\bar{\mathbf{A}}_2 = \begin{bmatrix} 0.22375543 & 0.22818147 & 0.13019843 \\ -0.22779964 & -0.23189051 & -0.13252835 \\ 0.12926548 & 0.13011081 & 0.07252603 \end{bmatrix}$$

$$\bar{\mathbf{A}}_3 = \begin{bmatrix} -0.00038721 & 0.00272144 & -0.00071499 \\ -0.00109674 & -0.00103953 & 0.00243030 \\ -0.00121810 & 0.00214948 & 0.00091477 \end{bmatrix}$$

$$\bar{\mathbf{A}}_4 = \begin{bmatrix} 0.79627366 & -0.23797139 & -0.10339016 \\ 0.21193424 & 0.37311787 & -0.45053023 \\ -0.06836792 & 0.40837312 & 0.73313847 \end{bmatrix}$$

$$\bar{\mathbf{b}}_1 = [0.11118275 \quad -0.11288667 \quad 0.06372873]^t$$

$$\bar{\mathbf{b}}_2 = [2.85345887 \quad -2.46012520 \quad 1.18577329]^t$$

$$\bar{\mathbf{c}}_1 = [0.48512286 \quad 0.41842076 \quad 0.20174307]$$

$$\bar{\mathbf{c}}_2 = [0.01891648 \quad 0.01922876 \quad 0.01090239]$$

which is L_2 optimal and from (17) and (57) we have

$$\bar{M}_{1/2} = J + L = 380.4402$$

$$\bar{m}_2 = 2151.9307$$

respectively, where $J = 367.3009$ and $L = 13.1394$.

The simulation results of the above examples are summarized in terms of the weighted sensitivities $\bar{M}_{1/2}$ and \bar{m}_2 in Table I. It is observed that the weighted sensitivity $\bar{M}_{1/2}$ of the L_1/L_2 -optimal filter structures is very close to the value of $J_{\min} + L_{\min}$. Also, \bar{m}_2 of the L_1/L_2 -optimal filter structures is not far away from \bar{m}_2 of the L_2 -optimal filter structures.

VI. CONCLUSION

Two frequency-weighted sensitivity measures have been defined as a generalization of those reported in [10] and [13]. To construct the 2-D coordinate transformation matrix such that the weighted L_1/L_2 mixed sensitivity is optimal in a certain sense, an analytical method has been developed to obtain the closed-form solution. The 1-D version of the

analytical method can be viewed as an alternative to the weighted sensitivity minimization algorithm reported in [4] and is much simpler than the algorithm which relies on the Lagrange multiplier method. In addition, the 1-D version of this closed-form solution has not been reported in [5]. An iterative procedure has been proposed to find the optimal coordinate transformation that minimizes the weighted L_2 -sensitivity measure. The merit of this procedure is that the estimate at each iteration can be derived analytically. Our first contribution in this paper has been the introduction of general unconstrained frequency weights for 2-D state-space digital filters. The second is to present a novel closed-form solution for obtaining the 2-D filter structures that minimize $\bar{M}_{1/2}(\mathbf{P})$, subject to the minimization of $J(\mathbf{P})$. The third is to develop a procedure for iteratively finding the optimal coordinate transformation that yields the filter structures with minimum weighted L_2 -sensitivity. We have illustrated the utility of the proposed technique with two numerical examples.

It should be noted that the approach presented here can be extended to the M -dimensional case where $M > 2$ in a straightforward manner, provided the multidimensional LSS model reported in [25] is employed. In addition, similar arguments can be applied to the Fornasini–Marchesini second LSS model [17], [15].

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] ———, "On the sensitivity of linear state-space systems," *IEEE Trans. Circuits Syst.*, vol. CAS-33, pp. 502–510, May 1986.
- [3] 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.
- [4] G. Li and M. Gevers, "Optimal finite precision implementation of a state-estimate feedback controller," *IEEE Trans. Circuits Syst.*, vol. 37, pp. 1487–1498, Dec. 1990.
- [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. Acoustics, Speech, Signal Processing*, vol. 4, pp. 429–432.
- [8] W. J. Lutz and S. L. Hakimi, "Design of multi-input multi-output systems with minimum sensitivity," *IEEE Trans. Circuits Syst.*, vol. 35, pp. 1114–1122, Sept. 1988.
- [9] A. Zilouchian and R. L. Carroll, "A coefficient sensitivity bound in 2-D state-space digital filtering," *IEEE Trans. Circuits Syst.*, vol. CAS-33, pp. 665–667, June 1986.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] ———, "Minimization of frequency-weighting sensitivity in 2-D systems based on the Fornasini–Marchesini second model," in *Proc. 1992*

- IEEE Int. Conf. Acoust., Speech, Signal Processing*, pp. 401–404; also “On the frequency-weighting sensitivity of 2-D state-space digital filters based on the Fornasini-Marchesini second model,” *IEICE Trans. Fundamentals*, vol. E75-A, pp. 813–820, July 1992.
- [15] 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.
- [16] R. P. Roesser, “A discrete state-space model for linear image processing,” *IEEE Trans. Automat. Contr.*, vol. AC-20, pp. 1–10, Feb. 1975.
- [17] E. Fornasini and G. Marchesini, “Doubly-indexed dynamical systems: State-space models and structural properties,” *Math. Syst. Theory*, vol. 12, pp. 59–72, 1978.
- [18] S. Kung, B. C. Lévy, M. Morf, and T. Kailath, “New results in 2-D systems theory, Part II: 2-D state-space models—Realization and the notions of controllability, observability, and minimality,” *Proc. IEEE*, vol. 65, pp. 945–961, June 1977.
- [19] L. L. Scharf, *Statistical Signal Processing*. Reading, MA: Addison-Wesley, 1991.
- [20] 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, 1976.
- [21] S. Shokoohi, “Block-balanced realizations,” in *Proc. 1984 IEEE Int. Symp. Circuits Systems*, pp. 825–829.
- [22] C. Xiao, P. Agathoklis, and D. J. Hill, “Coefficient sensitivity and structure optimization of multidimensional state-space digital filters,” in *Proc. 1997 IEEE Int. Symp. Circuits Systems*, pp. 2465–2468.
- [23] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 2nd ed. Baltimore, MD: Johns Hopkins Univ. Press, 1989.
- [24] B. W. Bowmar and J. C. Hung, “Minimum roundoff noise digital filters with some power-of-two coefficients,” *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp. 833–840, Oct. 1984.
- [25] D. S. K. Chan, “The structure of recursive multidimensional discrete systems,” *IEEE Trans. Automat. Contr.*, vol. AC-25, pp. 663–673, Aug. 1980.



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

From 1972 to 1988, he was with the Faculty of Engineering, Kobe University, Kobe, Japan. From September 1979 to March 1981, he was on leave from Kobe University as a Visiting Member of Staff

in the Department of Electrical Engineering, Queen's University, Kingston, ONT, Canada. During 1988–1991, he was Professor of Electronic Circuits in the Faculty of Engineering, Tottori University, Tottori, Japan. Since January 1992, he has been Professor of Electronic Control in the Department of Electrical Engineering, Hiroshima University, Higashi-Hiroshima, Japan. His teaching and research interests include digital signal processing, system theory, and control systems engineering. He has coedited and coauthored the book *Two-Dimensional Signal and Image Processing* (Tokyo, Japan: SICE, 1996).

Dr. Hinamoto was an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II during 1993–1995 and served as the Chair of the 12th Digital Signal Processing (DSP) Symposium held in Hiroshima in November 1997, which is the 12th of a series of annual symposia sponsored by the DSP Technical Committee of the IEICE. He also served as the Guest Editor of the special section on DSP in the August 1998 issue of the *IEICE Transactions on Fundamentals*.



Yoshitaka Zempo received the B.E. and M.E. degrees in electrical engineering from Hiroshima University in 1993 and 1995, respectively.

In April 1995 he joined Yamaguchi Television and Radio Broadcasting Station, Inc., Tokuyama, Japan. His research interests are in digital signal processing and system theory.



Yoshio Nishino received the B.E., M.E., and Dr. Eng. degrees in applied mathematics and physics from Kyoto University in 1971, 1973, and 1979, respectively.

He was a Lecturer in the Faculty of Engineering, Kyoto University, Kyoto, Japan, from 1979 to 1981. Since 1982 he has been an Associate Professor of the Faculty of Engineering, Hiroshima University, Hiroshima, Japan. His study is concerned with applications of differential geometry and group theory to problems arising in physics and engineering.



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

He was a post-doctoral fellow at the University of Victoria, Victoria, BC, Canada in 1985 and a visiting assistant professor at the University of Minnesota in 1986. Since 1987 he has been with the University of Victoria, where he is currently Professor. His

teaching and research interests are in the areas of digital signal processing and robotics. He is the co-author with A. Antoniou of *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 was the Editor of the same journal from 1990 to 1992. He was an Associate Editor of IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II from 1993 to 1995 and is presently an Associate Editor of *Multidimensional Systems and Signal Processing*.