

# Minimization of $L_2$ -Sensitivity for State-Space Digital Filters Subject to $L_2$ -Dynamic-Range-Scaling Constraints

Takao Hinamoto, Hiroaki Ohnishi and Wu-Sheng Lu

**Abstract** The problem of minimizing an  $L_2$ -sensitivity measure subject to  $L_2$ -norm dynamic-range scaling constraints for state-space digital filters is formulated. It is shown that the problem can be converted into an unconstrained optimization problem by using linear-algebraic techniques. Next, the unconstrained optimization problem is solved by applying an efficient quasi-Newton algorithm with closed-form formula for gradient evaluation. The coordinate transformation matrix obtained is then used to construct the optimal state-space filter structure that minimizes the  $L_2$ -sensitivity measure subject to the scaling constraints. A numerical example is presented to illustrate the utility of the proposed technique.

**Keywords**  $L_2$ -sensitivity, dynamic-range scaling constraints, scaling constrained sensitivity minimization, optimal realization, state-space digital filters

---

T. Hinamoto and H. Ohnishi are with the Graduate School of Engineering, Hiroshima University, Higashi-Hiroshima, Japan 739-8527.

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

## I. INTRODUCTION

The problem of realizing a fixed-point state-space digital filter with finite word length (FWL) is a significant research topic, since the efficiency and performance of the filter are directly affected by the choice of its state-space filter structure. When a transfer function with infinite accuracy coefficients is designed so as to meet the filter specification requirements and realized by a state-space model, in order to implement the filter in a finite binary representation, the coefficients in the state-space model must be truncated or rounded to fit the FWL constraints. This coefficient quantization usually alters the characteristics of the filter. For instance, it may change a stable filter to an unstable one. This motivates the study of the coefficient sensitivity minimization problem. In the literature, two main classes of techniques have been proposed for constructing state-space digital filters that minimize the coefficient sensitivity in [1]-[10]:  $L_1/L_2$ -sensitivity minimization [1]-[5] and  $L_2$ -sensitivity minimization [6]-[10]. It has been argued [6]-[10] that the sensitivity measure based on the  $L_2$  norm is more natural and reasonable relative to that based on the  $L_1/L_2$ -sensitivity minimization. The  $L_1/L_2$ -sensitivity minimization and  $L_2$ -sensitivity minimization have also been considered in linear continuous-time systems in [11] and [10], respectively. However, to our best knowledge, there is no study for the minimization of the  $L_2$ -sensitivity subject to the  $L_2$ -norm dynamic-range scaling constraints for state-space digital filters, although it has been known that the use of scaling constraints can be beneficial for suppressing overflow oscillation [12],[13].

This paper investigates the problem of minimizing the  $L_2$ -sensitivity measure subject to  $L_2$ -norm dynamic-range scaling constraints for state-space digital filters. To this end, we introduce an expression for evaluating the  $L_2$ -sensitivity and formulate the  $L_2$ -sensitivity minimization problem subject to the scaling constraints. Next, the constrained optimization problem is converted into an unconstrained optimization problem by using linear-algebraic techniques. The unconstrained optimization problem is then solved using an efficient quasi-Newton algorithm [14]. A numerical example is presented to demonstrate that the proposed algorithms offer reduced  $L_2$ -sensitivity compared with that obtained using the conventional methods.

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

## II. $L_2$ -SENSITIVITY ANALYSIS

Consider a stable, controllable and observable state-space digital filter

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

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

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

*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

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

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

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

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

From (2), *Definition 1* and *Definition 2*, the overall  $L_2$ -sensitivity measure for the state-space digital filter in (1) is defined as

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

where

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

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

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

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

where

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

The matrices  $\mathbf{K}_c$  and  $\mathbf{W}_o$  are called the controllability and observability Gramians, respectively. The Gramians  $\mathbf{M}_A$ ,  $\mathbf{K}_c$  and  $\mathbf{W}_o$  can be obtained by solving the Lyapunov equations [15]:

$$\begin{aligned} \begin{bmatrix} * & * \\ * & \mathbf{M}_A \end{bmatrix} &= \begin{bmatrix} \mathbf{A} & \mathbf{bc} \\ \mathbf{0} & \mathbf{A} \end{bmatrix}^T \begin{bmatrix} * & * \\ * & \mathbf{M}_A \end{bmatrix} \\ &\quad \cdot \begin{bmatrix} \mathbf{A} & \mathbf{bc} \\ \mathbf{0} & \mathbf{A} \end{bmatrix} + \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{K}_c &= \mathbf{AK}_c\mathbf{A}^T + \mathbf{bb}^T \\ \mathbf{W}_o &= \mathbf{A}^T\mathbf{W}_o\mathbf{A} + \mathbf{c}^T\mathbf{c} \end{aligned} \quad (7)$$

Utilizing the *Cauchy integral theorem*, matrix  $\mathbf{M}_A$  in (6) can be written as

$$\mathbf{M}_A = \sum_{k=0}^{\infty} \mathbf{H}^T(k)\mathbf{H}(k) \quad (8)$$

where

$$\mathbf{H}(k) = \sum_{p=0}^k \mathbf{A}^p \mathbf{bc} \mathbf{A}^{k-p}.$$

If a coordinate transformation defined by

$$\bar{\mathbf{x}}(k) = \mathbf{T}^{-1}\mathbf{x}(k) \quad (9)$$

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

$$\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T}. \quad (10)$$

From (2) and (10), it is clear that the transfer function  $H(z)$  is invariant under the coordinate transformation in (9). The coordinate transformation defined by (9) transforms the Gramians  $\{\mathbf{M}_A, \mathbf{K}_c, \mathbf{W}_o\}$  to  $\{\bar{\mathbf{M}}_A, \bar{\mathbf{K}}_c, \bar{\mathbf{W}}_o\}$ , and changes (6) to

$$S(\mathbf{T}) = \text{tr}[\bar{\mathbf{M}}_A] + \text{tr}[\bar{\mathbf{W}}_o] + \text{tr}[\bar{\mathbf{K}}_c] \quad (11)$$

where

$$\begin{aligned} \bar{\mathbf{M}}_A &= \sum_{k=0}^{\infty} \mathbf{T}^T \mathbf{H}^T(k) \mathbf{T}^{-T} \mathbf{T}^{-1} \mathbf{H}(k) \mathbf{T} \\ &= \mathbf{T}^T \hat{\mathbf{M}}_A \mathbf{T} \\ \begin{bmatrix} * & * \\ * & \hat{\mathbf{M}}_A \end{bmatrix} &= \begin{bmatrix} \mathbf{A} & \mathbf{b} \mathbf{c} \\ \mathbf{0} & \mathbf{A} \end{bmatrix}^T \begin{bmatrix} * & * \\ * & \hat{\mathbf{M}}_A \end{bmatrix} \\ &\quad \cdot \begin{bmatrix} \mathbf{A} & \mathbf{b} \mathbf{c} \\ \mathbf{0} & \mathbf{A} \end{bmatrix} + \begin{bmatrix} \mathbf{T}^{-T} \mathbf{T}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \bar{\mathbf{W}}_o &= \mathbf{T}^T \mathbf{W}_o \mathbf{T}, \quad \bar{\mathbf{K}}_c = \mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T}. \end{aligned}$$

Moreover, if the  $L_2$ -norm dynamic-range scaling constraints are imposed on the new state-variable vector  $\bar{\mathbf{x}}(k)$ , then it is required that for  $i = 1, 2, \dots, n$

$$(\bar{\mathbf{K}}_c)_{ii} = (\mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T})_{ii} = 1. \quad (12)$$

The problem of  $L_2$ -sensitivity minimization subject to  $L_2$ -norm dynamic-range scaling constraints is now formulated as follows: For given  $\mathbf{A}$ ,  $\mathbf{b}$  and  $\mathbf{c}$  (and therefore,  $\mathbf{M}_A$ ,  $\mathbf{K}_c$  and  $\mathbf{W}_o$ ), obtain an  $n \times n$  nonsingular matrix  $\mathbf{T}$  which minimizes (11) subject to the constraints in (12).

**Remark 1:** By definition, the sensitivity measure described in this paper has no upper bound but does have a lower bound (say zero) because the sensitivity measure is always nonnegative. This observation, in conjunction with the fact that both the objective function (defined by (11)) and constraint functions (defined by (12)) are continuously differentiable, leads us to conclude that the problem of minimizing  $S(\mathbf{T})$  subject to (12) admits local solutions. On the other hand, since  $S(\mathbf{T})$  and

the constraint functions are nonconvex, the best that an optimization algorithm can claim is that it identifies a local minimizer of the problem at hand.

**Remark 2:** It is well known that constrained minimization problems, especially the minimization of nonconvex function subject to nonconvex constraints, are in general considerably more involved than their unconstrained counterpart, but the problem is that it is not always possible to convert a nonconvex constrained problem to an unconstrained one. What is done in the next section is to show that in the present case we do have an equivalent unconstrained problem to work with.

### III. $L_2$ -SENSITIVITY MINIMIZATION UNDER SCALING CONSTRAINTS

When the state-space model (1) is assumed to be stable and controllable, the controllability Gramian  $\mathbf{K}_c$  is symmetric and positive-definite [15]. This implies that  $\mathbf{K}_c^{1/2}$  satisfying  $\mathbf{K}_c = \mathbf{K}_c^{1/2}\mathbf{K}_c^{1/2}$  is also symmetric and positive-definite. Defining

$$\hat{\mathbf{T}} = \mathbf{T}^T \mathbf{K}_c^{-\frac{1}{2}}, \quad (13)$$

the constraints in (12) can be expressed as

$$(\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1})_{ii} = 1, \quad i = 1, 2, \dots, n. \quad (14)$$

The constraints in (14) simply state that each column in  $\hat{\mathbf{T}}^{-1}$  must be a unity vector. If matrix  $\hat{\mathbf{T}}^{-1}$  is assumed to have the form

$$\hat{\mathbf{T}}^{-1} = \left[ \frac{\mathbf{t}_1}{\|\mathbf{t}_1\|}, \frac{\mathbf{t}_2}{\|\mathbf{t}_2\|}, \dots, \frac{\mathbf{t}_n}{\|\mathbf{t}_n\|} \right], \quad (15)$$

then (14) is always satisfied. From (13), it follows that (11) is changed to

$$\begin{aligned} J_o(\hat{\mathbf{T}}) = & \text{tr} \left[ \sum_{k=0}^{\infty} \hat{\mathbf{T}} \hat{\mathbf{H}}^T(k) \hat{\mathbf{T}}^{-1} \hat{\mathbf{T}}^{-T} \hat{\mathbf{H}}(k) \hat{\mathbf{T}}^T \right] \\ & + \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{W}}_o \hat{\mathbf{T}}^T] + \text{tr}[\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1}] \end{aligned} \quad (16)$$

where

$$\hat{\mathbf{H}}(k) = \mathbf{K}_c^{-\frac{1}{2}} \mathbf{H}(k) \mathbf{K}_c^{\frac{1}{2}}, \quad \hat{\mathbf{W}}_o = \mathbf{K}_c^{\frac{1}{2}} \mathbf{W}_o \mathbf{K}_c^{\frac{1}{2}}.$$

From the foregoing arguments, the problem of obtaining an  $n \times n$  nonsingular matrix  $\mathbf{T}$  which minimizes (11) subject to the constraints in (12) can be converted into an unconstrained optimization problem of obtaining an  $n \times n$  nonsingular matrix  $\hat{\mathbf{T}}$  which minimizes (16).

**Remark 3:** The use of a structured  $\mathbf{T}$  matrix as specified in (15) is merely to eliminate the constraints in (12). Needless to say, this increases the degree of nonlinearity for the objective function in (16). However, due to the equivalence of the two optimization problem, it can be concluded that the unconstrained minimization of the objective function in (16) admits at least local solutions.

Now we apply a quasi-Newton algorithm to minimize (16) with respect to matrix  $\hat{\mathbf{T}}$  given by (15). Let  $\mathbf{x}$  be the column vector that collects the variables in matrix  $\hat{\mathbf{T}}$ . Then  $J_o(\hat{\mathbf{T}})$  is a function of  $\mathbf{x}$ , which we denote by  $J(\mathbf{x})$ . The algorithm starts with a trivial initial point  $\mathbf{x}_0$  obtained from an initial assignment  $\hat{\mathbf{T}} = \mathbf{I}_n$ . Then, in the  $k$ th iteration a quasi-Newton algorithm updates the most recent point  $\mathbf{x}_k$  to point  $\mathbf{x}_{k+1}$  as

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad (17)$$

where [14]

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

Here,  $\nabla J(\mathbf{x})$  is the gradient of  $J(\mathbf{x})$  with respect to  $\mathbf{x}$ , and  $\mathbf{S}_k$  is a positive-definite approximation of the inverse Hessian matrix of  $J(\mathbf{x})$ . This iteration process continues until

$$|J(\mathbf{x}_{k+1}) - J(\mathbf{x}_k)| < \varepsilon \quad (18)$$

where  $\varepsilon > 0$  is a prescribed tolerance. If the iteration is terminated at step  $k$ , then  $\mathbf{x}_k$  is viewed as a solution point.

The implementation of (17) requires the gradient of  $J(\mathbf{x})$ . Closed-form expressions for  $\nabla J(\mathbf{x})$  are given below.

$$\begin{aligned}
\frac{\partial J_o(\hat{\mathbf{T}})}{\partial t_{ij}} &= \lim_{\Delta \rightarrow \infty} \frac{J_o(\hat{\mathbf{T}}_{ij}) - J_o(\hat{\mathbf{T}})}{\Delta} \\
&= 2\beta_1 - \beta_2 + 2\beta_3
\end{aligned} \tag{19}$$

where  $\hat{\mathbf{T}}_{ij}$  is the matrix obtained from  $\hat{\mathbf{T}}$  with a perturbed  $(i, j)$ th component, which is given by

$$\begin{aligned}
\hat{\mathbf{T}}_{ij} &= \hat{\mathbf{T}} + \frac{\Delta \hat{\mathbf{T}} \mathbf{g}_{ij} \mathbf{e}_j^T \hat{\mathbf{T}}}{1 - \Delta \mathbf{e}_j^T \hat{\mathbf{T}} \mathbf{g}_{ij}} \\
\beta_1 &= \mathbf{e}_j^T \sum_{k=0}^{\infty} \hat{\mathbf{T}} \hat{\mathbf{H}}^T(k) \hat{\mathbf{T}}^{-1} \hat{\mathbf{T}}^{-T} \hat{\mathbf{H}}(k) \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij} \\
\beta_2 &= \mathbf{e}_j^T \sum_{k=0}^{\infty} \hat{\mathbf{T}}^{-T} \hat{\mathbf{H}}(k) \hat{\mathbf{T}}^T \hat{\mathbf{T}} \hat{\mathbf{H}}^T(k) \mathbf{g}_{ij} \\
\beta_3 &= \mathbf{e}_j^T \hat{\mathbf{T}} \hat{\mathbf{W}}_o \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij} \\
\mathbf{g}_{ij} &= \partial \left\{ \frac{\mathbf{t}_j}{\|\mathbf{t}_j\|} \right\} / \partial t_{ij} \\
&= \frac{1}{\|\mathbf{t}_j\|^3} (t_{ij} \mathbf{t}_j - \|\mathbf{t}_j\|^2 \mathbf{e}_i)
\end{aligned}$$

where  $\mathbf{e}_i$  denotes an  $n \times 1$  unit vector whose  $i$ th element equals unity.

**Remark 4:** Although a global solution cannot be claimed, it may be worthwhile to report that we have applied the proposed algorithm to a significant number of state-space filters of various sizes, and without exception the algorithm converges to a solution that outperforms their counterparts obtained by the methods in [10] and [13]. In this regard, it is interesting to note that the unconstrained solution offered by the method in [13] is known to be globally optimal.

#### IV. NUMERICAL EXAMPLE

Consider a state-space digital filter, (1), specified by

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

$$d = 0.015940$$

whose poles are at  $z = 0.6578817$  and  $z = 0.6584892 \pm j0.5060989$ .

Using (7), the Gramians  $\mathbf{M}_A$ ,  $\mathbf{K}_c$  and  $\mathbf{W}_o$  are calculated as

$$\mathbf{M}_A = \begin{bmatrix} 8.921380 & -22.046457 & 17.916285 \\ -22.046457 & 55.671710 & -46.052011 \\ 17.916285 & -46.052011 & 42.522082 \end{bmatrix}$$

$$\mathbf{K}_c = \begin{bmatrix} 1.000000 & 0.872501 & 0.562821 \\ 0.872501 & 1.000000 & 0.872501 \\ 0.562821 & 0.872501 & 1.000000 \end{bmatrix}$$

$$\mathbf{W}_o = \begin{bmatrix} 0.820741 & -2.035328 & 1.628161 \\ -2.035328 & 5.307273 & -4.264903 \\ 1.628161 & -4.264903 & 3.941491 \end{bmatrix}$$

and the  $L_2$ -sensitivity measure  $S$  in (6) is found to be

$$S = 120.184677.$$

Choosing  $\hat{\mathbf{T}} = \mathbf{I}_n$  (therefore  $\mathbf{T} = \mathbf{K}_c^{1/2}$  in (13)) as the initial estimate and  $\varepsilon = 10^{-7}$ , it took the proposed quasi-Newton algorithm 15 iterations to converge to

$$\hat{\mathbf{T}}^{opt} = \begin{bmatrix} 0.376709 & -0.319162 & 0.256412 \\ 0.910212 & 0.154833 & -0.218993 \\ 0.172058 & 0.934967 & 0.941433 \end{bmatrix}$$

or equivalently,

$$\mathbf{T}^{opt} = \begin{bmatrix} 0.913655 & -0.857313 & 0.877296 \\ 0.905773 & -0.121938 & 0.493844 \\ 0.576905 & 0.415235 & 0.377361 \end{bmatrix}$$

where (16) is used by truncating the infinite sum with  $k = 100$ . In this case, the  $L_2$ -sensitivity measure in (11) is minimized subject to the scaling constraints in (12) to

$$S(\mathbf{T}) = 8.683279.$$

The optimal state-space filter structure that minimizes the  $L_2$ -sensitivity measure, (11), subject to the scaling constraints in (12) is then constructed by substituting  $\mathbf{T} = \mathbf{T}^{opt}$  into (10) as

$$\bar{\mathbf{A}} = \begin{bmatrix} 0.586086 & 0.567626 & 0.080361 \\ -0.450289 & 0.725542 & 0.188298 \\ -0.017947 & -0.021129 & 0.663233 \end{bmatrix}$$

$$\bar{\mathbf{b}} = \begin{bmatrix} -0.362927 \\ 0.393927 \\ 0.762923 \end{bmatrix}$$

$$\bar{\mathbf{c}} = [ 0.362538 \quad 0.042368 \quad 0.254527 ]$$

$$d = 0.015940.$$

The  $L_2$ -sensitivity profile of first 18 iterations is given in Table I, where  $J_o(\hat{\mathbf{T}})$  in (16) is used to evaluate the  $L_2$ -sensitivity under the same truncation of the infinite sum. From Table I, it is seen that with a tolerance  $\varepsilon = 10^{-7}$  the algorithm converges with 15 iterations.

TABLE I  
 $L_2$ -SENSITIVITY PROFILE OF FIRST 18 ITERATIONS

| $k$ | $L_2$ -Sensitivity | $k$ | $L_2$ -Sensitivity |
|-----|--------------------|-----|--------------------|
|     | 120.18467700       | 9   | 8.80368098         |
| 0   | 10.71346288        | 10  | 8.70867352         |
| 1   | 10.70375499        | 11  | 8.70588841         |
| 2   | 10.70375483        | 12  | 8.69200561         |
| 3   | 10.33912612        | 13  | 8.69014918         |
| 4   | 10.27174422        | 14  | 8.68509811         |
| 5   | 9.59476520         | 15  | 8.68327357         |
| 6   | 8.99317571         | 16  | 8.68327359         |
| 7   | 8.93852568         | 17  | 8.68327362         |
| 8   | 8.88777828         | 18  | 8.68327366         |

For comparison purposes, the method reported in [10] is applied to minimize the  $L_2$ -sensitivity measure in (11) (without considering the scaling constraints in (12)) and the resulting optimal coordinate transformation matrix is scaled by an appropriate nonsingular diagonal matrix, so that (12) is satisfied. The result is

$$S(\mathbf{T}) = 9.817579$$

where

$$\mathbf{T} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.594723 & 0.562052 & 0.0 \\ 0.221714 & 0.736136 & 0.306792 \end{bmatrix}$$

Moreover, by applying the method reported in [13], the optimal coordinate transformation matrix  $\mathbf{T}$  is constructed as

$$\mathbf{T} = \begin{bmatrix} -0.605406 & -0.119653 & 1.219423 \\ 0.107851 & 0.097317 & 0.941720 \\ 0.540830 & -0.071898 & 0.569047 \end{bmatrix}$$

which minimizes the roundoff noise at the filter output subject to the scaling constraints in (12). The  $L_2$ -sensitivity of the resulting filter with minimum roundoff noise is computed as

$$S(\mathbf{T}) = 8.797931.$$

It is noted that these values of the  $L_2$ -sensitivity are larger than  $S(\mathbf{T}) = 8.683279$  obtained by the proposed method.

Finally, the results obtained in the simulation are summarized in Table II. From this table, it is observed that the proposed technique offers the smallest  $L_2$ -sensitivity subject to the  $L_2$ -norm dynamic-range scaling constraints relative to the existing methods presented in [10] and [13] for state-space digital filters.

TABLE II  
 $L_2$ -SENSITIVITY COMPARISON

| Realization    | $L_2$ -Sensitivity |
|----------------|--------------------|
| Original       | 120.184661         |
| Proposed       | 8.683279           |
| Method in [10] | 9.817579           |
| Method in [13] | 8.797931           |

## V. CONCLUSION

The problem of minimizing the  $L_2$ -sensitivity of a state-space digital filter subject to  $L_2$ -norm dynamic-range scaling constraints has been investigated. It has been shown that the  $L_2$ -sensitivity minimization problem subject to the scaling constraints can be converted into an unconstrained optimization problem by using linear algebraic techniques. An efficient quasi-Newton algorithm has been applied to solve the unconstrained optimization problem. The coordinate transformation matrix obtained has allowed us to construct the optimal state-space filter structure. Our computer simulation results have demonstrated the effectiveness of the proposed technique compared with the existing methods [10],[13].

## 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. 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] U. Helmke and J. B. Moore, *Optimization and Dynamical Systems*, Springer-Verlag, London, 1994.
- [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] 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.
- [12] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed-point digital filters," *IEEE Trans. Circuits Syst.*, vol. 23, pp. 551-562, Sept. 1976.
- [13] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 25, pp. 273-281, Aug. 1977.
- [14] R. Fletcher, *Practical Methods of Optimization*, 2nd ed. Wiley, New York, 1987.
- [15] T. Kailath, *Linear System*, Englewood Cliffs, N.J.: Prentice Hall, 1980.