

MINIMIZED ROUND-OFF NOISE AND POLE SENSITIVITY SUBJECT TO L_2 -SCALING CONSTRAINTS FOR IIR FILTERS

Yoichi Hinamoto

Akimitsu Doi

Department of Electrical and Computer Engineering
National Institute of Technology, Kagawa College
Takamatsu, Kagawa 761-8058, Japan
Email: hinamoto@t.kagawa-nct.ac.jp

Department of Computer Science
Hiroshima Institute of Technology
Hiroshima 731-5193, Japan
Emails: doi@cc.it-hiroshima.ac.jp

ABSTRACT

This paper investigates the minimization problem of weighted roundoff noise and pole sensitivity subject to l_2 -scaling constraints for state-space digital filters. A new measure for evaluating roundoff noise and pole sensitivity is proposed, and an efficient technique for minimizing this measure is developed. It is shown that the problem can be converted into an unconstrained optimization problem by using linear-algebraic techniques. The unconstrained optimization problem at hand is then solved iteratively by employing an efficient quasi-Newton algorithm with closed-form formulas for key gradient evaluation. Finally a numerical example is presented to demonstrate the validity and effectiveness of the proposed technique.

1. INTRODUCTION

In fixed-point IIR digital filters, roundoff noise occurs due to finite-precision nature of computer arithmetic and in turn degrades the filters' performance. It is well known that the roundoff noise is critically dependent on the internal structure of an IIR digital filter, and several techniques for synthesizing a state-space filter structure that minimizes the roundoff noise at the filter output subject to l_2 -scaling constraints have been explored by appropriately choosing a linear transformation to state-space coordinates [1]-[4]. When implementing a state-space model, the coefficients must be truncated or rounded to fit the finite-word-length (FWL) constraints. This coefficient quantization usually changes the characteristics of the filter. For instance, it may alter a stable filter to unstable one. This motivates the study of minimizing coefficients sensitivity. There exist several ways to define sensitivity of a filter with respect to its realization coefficients. Two of them rely on a mixed l_1/l_2 norm or a pure l_2 norm that measures changes of a certain transfer function, while the other is defined in terms of the poles and zeros of a filter. Several techniques for minimizing the l_1/l_2 -sensitivity measure [5]-[9] and the l_2 -sensitivity measure [10]-[13] have been proposed. Alternatively, the pole and zero sensitivity of a filter with respect to state-space parameters has been analyzed, and its reduction or minimization have been considered [14], [15]. It

has been stated in [14] that the sensitivities of poles and zeros of a transfer function with respect to the variations in coefficients of a state-space model (due to FWL) are closely related to the sensitivity characteristics of its frequency transfer function, and that minimal pole sensitivity is achieved when matrix \mathbf{A} is normal.

In this paper, the problem of minimizing weighted roundoff noise and pole sensitivity subject to l_2 -scaling constraints for state-space digital filters is investigated.

2. PROBLEM FORMULATION

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

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{b}u(k), \quad y(k) = \mathbf{c}\mathbf{x}(k) \quad (1)$$

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} and \mathbf{c} are real constant matrices of appropriate dimensions.

A. Roundoff Noise Analysis

By taking the quantizations performed before matrix-vector multiplication into account, an FWL implementation of the filter in (1) can be obtained as [14]

$$\begin{aligned} \tilde{\mathbf{x}}(k+1) &= \mathbf{A}\mathbf{Q}[\tilde{\mathbf{x}}(k)] + \mathbf{b}u(k) \\ \tilde{y}(k) &= \mathbf{c}\mathbf{Q}[\tilde{\mathbf{x}}(k)]. \end{aligned} \quad (2)$$

The matrices \mathbf{A} , \mathbf{b} , and \mathbf{c} in (2) are assumed to have exact fractional B_c -bit representations. The FWL state-variable vector $\tilde{\mathbf{x}}(k)$ and the output $\tilde{y}(k)$ all have B -bit fractional representations, while the input $u(k)$ is a $(B - B_c)$ -bit fraction. The quantizer $\mathbf{Q}[\cdot]$ in (2) rounds the B -bit fraction $\tilde{\mathbf{x}}(k)$ to $(B - B_c)$ -bit after the multiplications and additions, where the sign bit is not counted. Subtracting (2) from (1) yields

$$\begin{aligned} \Delta\mathbf{x}(k+1) &= \mathbf{A}\Delta\mathbf{x}(k) + \mathbf{A}\mathbf{e}(k) \\ \Delta y(k) &= \mathbf{c}\Delta\mathbf{x}(k) + \mathbf{c}\mathbf{e}(k) \end{aligned} \quad (3)$$

where $\mathbf{e}(k) = \tilde{\mathbf{x}}(k) - \mathbf{Q}[\tilde{\mathbf{x}}(k)]$, $\Delta\mathbf{x}(k) = \mathbf{x}(k) - \tilde{\mathbf{x}}(k)$ and $\Delta y(k) = y(k) - \tilde{y}(k)$. It is assumed that the roundoff error

$e(k)$ can be modeled as a zero-mean noise process with covariance $\sigma^2 \mathbf{I}_n$. By taking the z -transform on both sides of (3) and setting $\Delta \mathbf{x}(0) = \mathbf{0}$, we have

$$\begin{aligned} \Delta Y(z) &= \mathbf{H}_e(z) \mathbf{E}(z) \\ \mathbf{H}_e(z) &= \mathbf{c}(z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{A} + \mathbf{c} = \sum_{k=0}^{\infty} \mathbf{c} \mathbf{A}^k z^{-k} \end{aligned} \quad (4)$$

where $\Delta Y(z)$ and $\mathbf{E}(z)$ stand for the z -transforms of $\Delta y(k)$ and $e(k)$, respectively. Next, the normalized noise gain $J_R = \sigma_{out}^2 / \sigma^2$ is defined as

$$J_R = \text{tr} \left[\frac{1}{2\pi j} \oint_{|z|=1} \mathbf{H}_e^*(z) \mathbf{H}_e(z) \frac{dz}{z} \right]. \quad (5)$$

Substituting (4) into (5) yields

$$J_R = \text{tr} [\mathbf{W}] \quad (6)$$

where \mathbf{W} is the observability Gramian of the digital filter in (1) that can be obtained by solving the Lyapunov equation

$$\mathbf{W} = \mathbf{A}^T \mathbf{W} \mathbf{A} + \mathbf{c}^T \mathbf{c}. \quad (7)$$

Applying a coordinate transformation defined by

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

to the digital filter in (1), a different yet equivalent state-space realization of (1), $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}})_n$, can be characterized by

$$\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T}. \quad (9)$$

With an equivalent state-space realization as specified in (9), the normalized noise gain in (6) is changed to

$$J_R(\mathbf{T}) = \text{tr} [\bar{\mathbf{W}}] = \text{tr} [\mathbf{T}^T \mathbf{W} \mathbf{T}] \quad (10)$$

where $\bar{\mathbf{W}}$ denotes the observability Gramian for an equivalent state-space realization as specified in (9).

B. Pole Sensitivity Analysis

The transfer function of the digital filter in (1) can be expressed as

$$H(z) = \mathbf{c}(z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b} \quad (11)$$

Suppose that the poles $\{\lambda_l\}$ of $H(z)$ are denoted by $\{\lambda_l\} = \lambda(\mathbf{A})$. Moreover, let $\mathbf{x}_p(l)$ be a right eigenvector corresponding to λ_l and let $\mathbf{y}_p(l)$ be the reciprocal left eigenvector that corresponds to $\mathbf{x}_p(l)$. Then, assuming that \mathbf{A} has a full set of linearly independent eigenvectors, it follows that [15]

$$\left(\frac{\partial \lambda_l}{\partial \mathbf{A}} \right)^T = \mathbf{x}_p(l) \mathbf{y}_p^H(l) \quad \text{for } l = 1, 2, \dots, n. \quad (12)$$

Using the Frobenius norm for an $m \times n$ complex matrix \mathbf{M} defined by

$$\|\mathbf{M}\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |(\mathbf{M})_{ij}|^2 \right)^{\frac{1}{2}} \quad (13)$$

the pole sensitivity measure for the digital filter in (1) is defined as [15]

$$J_p = \sum_{l=1}^n \left\| \frac{\partial \lambda_l}{\partial \mathbf{A}} \right\|_F^2 = \sum_{l=1}^n \left\| \left(\frac{\partial \lambda_l}{\partial \mathbf{A}} \right)^T \right\|_F^2. \quad (14)$$

By virtue of

$$\|\mathbf{M}\|_F^2 = \text{tr}[\mathbf{M} \mathbf{M}^H] = \text{tr}[\mathbf{M}^H \mathbf{M}] = \|\mathbf{M}^H\|_F^2 \quad (15)$$

and using (12), the pole sensitivity measure in (14) is written as

$$J_p = \sum_{l=1}^n \left(\mathbf{x}_p^H(l) \mathbf{x}_p(l) \right) \left(\mathbf{y}_p^H(l) \mathbf{y}_p(l) \right). \quad (16)$$

With an equivalent state-space realization as specified in (9), the right eigenvector $\mathbf{x}_p(l)$ corresponding to λ_l and the reciprocal left eigenvector $\mathbf{y}_p(l)$ corresponding to $\mathbf{x}_p(l)$ are changed to

$$\bar{\mathbf{x}}_p(l) = \mathbf{T}^{-1} \mathbf{x}_p(l), \quad \bar{\mathbf{y}}_p(l) = \mathbf{T}^T \mathbf{y}_p(l) \quad (17)$$

respectively, for $l = 1, 2, \dots, n$. Thus, referring to (16), the pole sensitivity measure $J_p(\mathbf{T})$ for an equivalent state-space realization can be expressed as

$$J_p(\mathbf{T}) = \sum_{l=1}^n \left(\mathbf{x}_p^H(l) \mathbf{T}^{-T} \mathbf{T}^{-1} \mathbf{x}_p(l) \right) \left(\mathbf{y}_p^H(l) \mathbf{T} \mathbf{T}^T \mathbf{y}_p(l) \right). \quad (18)$$

It is known that an inequality $J_p(\mathbf{T}) \geq n$ are always satisfied for any $n \times n$ nonsingular matrix \mathbf{T} [15].

Remark 1 [14], [15]: It is noted that zero sensitivity for the digital filter in (1) can be examined in the case where a direct path from the input to the output in (1) exists.

C. Problem Statement

It should be noted that the l_2 -scaling constraints on the state-variable vector $\mathbf{x}(k)$ involve the controllability Gramian \mathbf{K} of the digital filter in (1), which can be computed by solving the Lyapunov equation

$$\mathbf{K} = \mathbf{A} \mathbf{K} \mathbf{A}^T + \mathbf{b} \mathbf{b}^T. \quad (19)$$

In order to suppress overflow oscillations, l_2 -scaling constraints are imposed on the state-variable vector $\bar{\mathbf{x}}(k)$ so that

$$(\bar{\mathbf{K}})_{ii} = (\mathbf{T}^{-1} \mathbf{K} \mathbf{T}^{-T})_{ii} = 1, \quad i = 1, 2, \dots, n \quad (20)$$

where $\bar{\mathbf{K}}$ indicates the controllability Gramian for an equivalent state-space realization as specified in (9).

In this paper, we consider a weighted roundoff noise and pole sensitivity measure $J_\gamma(\mathbf{T})$ defined by

$$J_\gamma(\mathbf{T}) = (1 - \gamma) J_R(\mathbf{T}) + \gamma J_p(\mathbf{T}) \quad (21)$$

where $0 \leq \gamma \leq 1$ is a weighting factor. It is evident that a $J_\gamma(\mathbf{T})$ with a greater γ represents a measure that places more

emphasis on pole sensitivity while a $J_\gamma(\mathbf{T})$ with a smaller γ serves as a measure that weights more heavily on roundoff noise. In addition, by setting γ to unity or zero $J_\gamma(\mathbf{T})$ becomes $J_p(\mathbf{T})$ or $J_R(\mathbf{T})$, respectively.

The problem being considered here is to design the optimal coordinate transformation matrix \mathbf{T} that minimize (21) subject to l_2 -scaling constraints in (20) where $0 \leq \gamma \leq 1$ is a weighting factor specified by the designer.

Remark 2 [10]: It is noted that an l_2 -sensitivity measure for the equivalent state-space realization $(\overline{\mathbf{A}}, \overline{\mathbf{b}}, \overline{\mathbf{c}})_n$ in (9) can be evaluated by

$$\begin{aligned} S_2(\overline{\mathbf{A}}, \overline{\mathbf{b}}, \overline{\mathbf{c}}) &= \left\| \frac{\partial H(z)}{\partial \overline{\mathbf{A}}} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial \overline{\mathbf{b}}} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial \overline{\mathbf{c}}^T} \right\|_2^2 \\ &= \text{tr} \left\{ [\mathbf{T}^{-1} \ \mathbf{0}] \mathbf{R} \begin{bmatrix} \mathbf{T}^{-T} \\ \mathbf{0} \end{bmatrix} \right\} + \text{tr}[\overline{\mathbf{W}}] + \text{tr}[\overline{\mathbf{K}}] \end{aligned} \quad (22)$$

where matrix \mathbf{R} is obtained by solving the Lyapunov equation

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

3. WEIGHTED ROUND OFF NOISE AND POLE SENSITIVITY MINIMIZATION

In this section, we define

$$\hat{\mathbf{T}} = \mathbf{T}^T \mathbf{K}^{-\frac{1}{2}} \quad (23)$$

which leads (20) to

$$(\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1})_{ii} = 1 \text{ for } i = 1, 2, \dots, n. \quad (24)$$

It is obvious that the conditions in (24) are always satisfied by choosing matrix $\hat{\mathbf{T}}^{-1}$ as

$$\hat{\mathbf{T}}^{-1} = \begin{bmatrix} \frac{\mathbf{t}_1}{\|\mathbf{t}_1\|} & \frac{\mathbf{t}_2}{\|\mathbf{t}_2\|} & \dots & \frac{\mathbf{t}_n}{\|\mathbf{t}_n\|} \end{bmatrix}. \quad (25)$$

By defining an $n^2 \times 1$ vector $\mathbf{x} = (\mathbf{t}_1^T, \mathbf{t}_2^T, \dots, \mathbf{t}_n^T)^T$, the normalized noise gain $J_R(\mathbf{T})$ in (10) can be expressed as

$$J_R(\mathbf{x}) = \text{tr} \left[\hat{\mathbf{T}} \hat{\mathbf{W}} \hat{\mathbf{T}}^T \right] \quad (26)$$

where

$$\hat{\mathbf{W}} = \mathbf{K}^{\frac{1}{2}} \mathbf{W} \mathbf{K}^{\frac{1}{2}}.$$

Moreover, the pole sensitivity measure $J_p(\mathbf{T})$ in (18) can be written as

$$J_p(\mathbf{x}) = \sum_{l=1}^n \left(\hat{\mathbf{x}}_p^H(l) \hat{\mathbf{T}}^{-1} \hat{\mathbf{T}}^{-T} \hat{\mathbf{x}}_p(l) \right) \left(\hat{\mathbf{y}}_p^H(l) \hat{\mathbf{T}}^T \hat{\mathbf{T}} \hat{\mathbf{y}}_p(l) \right) \quad (27)$$

where

$$\hat{\mathbf{x}}_p(l) = \mathbf{K}^{-\frac{1}{2}} \mathbf{x}_p(l), \quad \hat{\mathbf{y}}_p(l) = \mathbf{K}^{\frac{1}{2}} \mathbf{y}_p(l).$$

In this way, the original constrained optimization problem can be converted into an unconstrained optimization problem of obtaining an $n^2 \times 1$ vector \mathbf{x} which minimizes

$$J_\gamma(\mathbf{x}) = (1 - \gamma) J_R(\mathbf{x}) + \gamma J_p(\mathbf{x}). \quad (28)$$

Applying a quasi-Newton algorithm to minimize $J_\gamma(\mathbf{x})$ in (28), in the k th iteration the most recent point \mathbf{x}_k is updated to point \mathbf{x}_{k+1} as [16]

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad (29)$$

where

$$\mathbf{d}_k = -\mathbf{S}_k \nabla J_\gamma(\mathbf{x}_k), \quad \alpha_k = \arg \min_{\alpha} J_\gamma(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

$$\mathbf{S}_{k+1} = \mathbf{S}_k + \left(1 + \frac{\boldsymbol{\varphi}_k^T \mathbf{S}_k \boldsymbol{\varphi}_k}{\boldsymbol{\varphi}_k^T \boldsymbol{\delta}_k} \right) \frac{\boldsymbol{\delta}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\varphi}_k^T \boldsymbol{\delta}_k} - \frac{\boldsymbol{\delta}_k \boldsymbol{\varphi}_k^T \mathbf{S}_k + \mathbf{S}_k \boldsymbol{\varphi}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\varphi}_k^T \boldsymbol{\delta}_k}$$

$$\mathbf{S}_0 = \mathbf{I}, \quad \boldsymbol{\delta}_k = \mathbf{x}_{k+1} - \mathbf{x}_k, \quad \boldsymbol{\varphi}_k = \nabla J_\gamma(\mathbf{x}_{k+1}) - \nabla J_\gamma(\mathbf{x}_k).$$

In the above, $\nabla J_\gamma(\mathbf{x})$ is the gradient of $J_\gamma(\mathbf{x})$ with respect to \mathbf{x} , and \mathbf{S}_k is a positive-definite approximation of the inverse Hessian matrix of $J_\gamma(\mathbf{x})$. The above algorithm starts with a trivial initial point \mathbf{x}_0 obtained from an initial assignment $\hat{\mathbf{T}}^{-1} = \mathbf{I}_n$, and this iteration process continues until

$$|J_\gamma(\mathbf{x}_{k+1}) - J_\gamma(\mathbf{x}_k)| < \varepsilon \quad (30)$$

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

The gradient of $J_\gamma(\mathbf{x})$ can be evaluated using closed-form expressions as

$$\begin{aligned} \frac{\partial J_\gamma(\mathbf{x})}{\partial t_{ij}} &= \lim_{\Delta \rightarrow 0} \frac{J_\gamma(\hat{\mathbf{T}}_{ij}) - J_\gamma(\hat{\mathbf{T}})}{\Delta} \\ &= (1 - \gamma) \frac{\partial J_R(\mathbf{x})}{\partial t_{ij}} + \gamma \frac{\partial J_p(\mathbf{x})}{\partial t_{ij}} \end{aligned} \quad (31)$$

where $\hat{\mathbf{T}}_{ij}$ is the matrix obtained from $\hat{\mathbf{T}}$ with a perturbed (i, j) th component, it follows that [17, p. 655]

$$\hat{\mathbf{T}}_{ij} = \hat{\mathbf{T}} + \frac{\Delta \hat{\mathbf{T}} \mathbf{g}_{ij} \mathbf{e}_j^T \hat{\mathbf{T}}}{1 - \Delta \mathbf{e}_j^T \hat{\mathbf{T}} \mathbf{g}_{ij}}, \quad \hat{\mathbf{T}}_{ij}^{-1} = \hat{\mathbf{T}}^{-1} - \Delta \mathbf{g}_{ij} \mathbf{e}_j^T$$

$$\mathbf{g}_{ij} = -\partial \left\{ \frac{\mathbf{t}_j}{\|\mathbf{t}_j\|} \right\} / \partial t_{ij} = \frac{1}{\|\mathbf{t}_j\|^3} (t_{ij} \mathbf{t}_j - \|\mathbf{t}_j\|^2 \mathbf{e}_i)$$

and

$$\frac{\partial J_R(\mathbf{x})}{\partial t_{ij}} = 2 \mathbf{e}_j^T \hat{\mathbf{T}} \hat{\mathbf{W}} \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij}$$

$$\begin{aligned} \frac{\partial J_p(\mathbf{x})}{\partial t_{ij}} &= \sum_{l=1}^n \left\{ \left(\hat{\mathbf{x}}_p^H(l) \hat{\mathbf{M}}(\hat{\mathbf{T}}) \hat{\mathbf{x}}_p(l) \right) \left(\hat{\mathbf{y}}_p^H(l) \hat{\mathbf{T}}^T \hat{\mathbf{T}} \hat{\mathbf{y}}_p(l) \right) \right. \\ &\quad \left. + \left(\hat{\mathbf{x}}_p^H(l) \hat{\mathbf{T}}^{-1} \hat{\mathbf{T}}^{-T} \hat{\mathbf{x}}_p(l) \right) \left(\hat{\mathbf{y}}_p^H(l) \hat{\mathbf{N}}(\hat{\mathbf{T}}) \hat{\mathbf{y}}_p(l) \right) \right\} \end{aligned}$$

with

$$\hat{\mathbf{M}}(\hat{\mathbf{T}}) = -[\mathbf{g}_{ij} \mathbf{e}_j^T \hat{\mathbf{T}}^{-T} + \hat{\mathbf{T}}^{-1} \mathbf{e}_j \mathbf{g}_{ij}^T]$$

$$\hat{\mathbf{N}}(\hat{\mathbf{T}}) = \hat{\mathbf{T}}^T [\mathbf{e}_j \mathbf{g}_{ij}^T \hat{\mathbf{T}}^T + \hat{\mathbf{T}} \mathbf{g}_{ij} \mathbf{e}_j^T] \hat{\mathbf{T}}.$$

4. AN ILLUSTRATIVE EXAMPLE

Consider the 4th-order Butterworth lowpass filter $(\mathbf{A}, \mathbf{b}, \mathbf{c})_4$ with a narrow normalized bassband of 0.05, described by

$$\mathbf{A} = \begin{bmatrix} 3.589734 & 1 & 0 & 0 \\ -4.851276 & 0 & 1 & 0 \\ 2.924053 & 0 & 0 & 1 \\ -0.663010 & 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{b} = 10^{-3} \begin{bmatrix} 0.237096 \\ 0.035885 \\ 0.216300 \\ 0.010527 \end{bmatrix}$$

$$\mathbf{c} = [1 \ 0 \ 0 \ 0]$$

whose numerator b and denominator a was found using MATLAB as $[b, a] = \text{butter}(4, 0.05)$.

When a coordinate transformation defined by

$$\mathbf{T}_o = \text{diag}\{0.226458, 0.588059, 0.513017, 0.150144\}$$

was applied to the above filter $(\mathbf{A}, \mathbf{b}, \mathbf{c})_4$, the resulting controllability and observability Gramians $\mathbf{K}_o = \mathbf{T}_o^{-1} \mathbf{K} \mathbf{T}_o^{-T}$ and $\mathbf{W}_o = \mathbf{T}_o^T \mathbf{W} \mathbf{T}_o$ were computed using (19) and (7) as

$$\mathbf{K}_o = \begin{bmatrix} 1.000000 & -0.999248 & 0.997433 & -0.994918 \\ -0.999248 & 1.000000 & -0.999452 & 0.998047 \\ 0.997433 & -0.999452 & 1.000000 & -0.999566 \\ -0.994918 & 0.998047 & -0.999566 & 1.000000 \end{bmatrix}$$

$$\mathbf{W}_o = 10^4 \begin{bmatrix} 1.063597 & 2.747677 & 2.360090 & 0.672994 \\ 2.747677 & 7.172055 & 6.224575 & 1.793652 \\ 2.360090 & 6.224575 & 5.458399 & 1.589267 \\ 0.672994 & 1.793652 & 1.589267 & 0.467539 \end{bmatrix}$$

Then, the normalized noise gain in (10) subject to l_2 -scaling constraints: $(\mathbf{K}_o)_{ii} = 1$ for $i = 1, 2, \dots, n$ was found to be

$$J_R(\mathbf{T}_o) = \text{tr}[\mathbf{W}_o] = 1.416159 \times 10^5.$$

Moreover, the l_2 -sensitivity measure in (22) was computed as

$$S_2(\mathbf{T}_o^{-1} \mathbf{A} \mathbf{T}_o, \mathbf{T}_o^{-1} \mathbf{b}, \mathbf{c} \mathbf{T}_o) = 9.779175 \times 10^6.$$

Next, the eigenvalues of matrix \mathbf{A} were found to be

$$\lambda = 0.931900 \pm j0.136363, \quad 0.862967 \pm j0.052305.$$

The original pole sensitivity measures in (16) and (18) were computed as

$$J_p = 1.863101 \times 10^7 \quad \text{and} \quad J_p(\mathbf{T}_o) = 1.774671 \times 10^7$$

respectively.

The quasi-Newton algorithm was applied to minimize (28) with $\gamma = 0.7$ by choosing $\hat{\mathbf{T}}^{-1} = \mathbf{I}_4$ as an initial assignment, and setting tolerance to $\varepsilon = 10^{-8}$ in (30). It took the algorithm 67 iterations to converge to the solution

$$\hat{\mathbf{T}} = \begin{bmatrix} -0.005208 & 1.163747 & -0.111659 & -0.787204 \\ -0.037351 & -0.404560 & 1.799691 & 2.040456 \\ -0.164510 & -0.518875 & -1.552980 & 4.982677 \\ 1.193078 & 0.558202 & -0.064418 & -6.380411 \end{bmatrix}$$

which is equivalent to

$$\mathbf{T} = \begin{bmatrix} -0.178135 & 0.230044 & -0.317367 & 0.217107 \\ 0.449984 & -0.602281 & 0.885719 & -0.599864 \\ -0.380296 & 0.526722 & -0.824583 & 0.559049 \\ 0.107642 & -0.153385 & 0.256254 & -0.175300 \end{bmatrix}$$

and the objective function in (21) was found to be

$$J_\gamma(\mathbf{T}) = 3.246633$$

with $J_R(\mathbf{T}) = 1.327903$ and $J_p(\mathbf{T}) = 4.068946$.

The profile of $J_\gamma(\mathbf{x})$ during the first 67 iterations of the algorithm is depicted in Fig. 1. From this figure it is observed that the weighted roundoff noise and pole sensitivity for a state-space digital filter were drastically reduced subject to l_2 -scaling constraints. In this case, the coefficient matrices

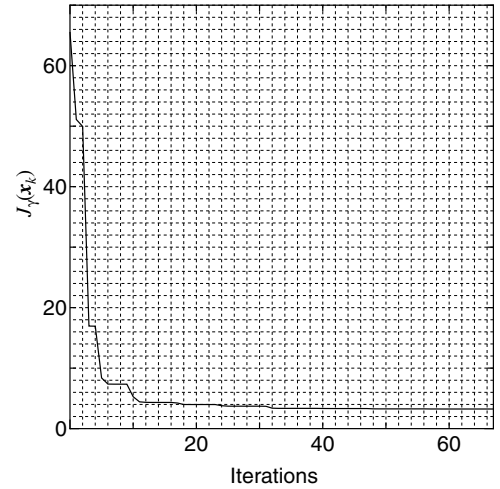


Fig. 1. Profile of $J_\gamma(\mathbf{x})$ during the first 67 iterations.

$\bar{\mathbf{A}}$, $\bar{\mathbf{b}}$ and $\bar{\mathbf{c}}$ for the equivalent state-space realization specified by (9) were obtained as

$$\bar{\mathbf{A}} = \begin{bmatrix} 0.925785 & 0.005903 & -0.123828 & -0.066614 \\ -0.042551 & 0.870369 & 0.037577 & -0.044076 \\ 0.097415 & -0.053644 & 0.909232 & 0.044846 \\ 0.074370 & 0.033706 & 0.019872 & 0.884347 \end{bmatrix}$$

$$\bar{\mathbf{b}} = [0.472882 \ 0.535104 \ 0.318385 \ 0.287515]^T$$

$$\bar{\mathbf{c}} = [-0.178135 \ 0.230044 \ -0.317367 \ 0.217107]$$

and the controllability and observability Gramians $\bar{\mathbf{K}}$ and $\bar{\mathbf{W}}$ were computed from (20) and (10) as

$$\bar{\mathbf{K}} = \begin{bmatrix} 1.000000 & 0.430444 & 0.368731 & 0.281793 \\ 0.430444 & 1.000000 & 0.825645 & 0.884131 \\ 0.368731 & 0.825645 & 1.000000 & 0.969484 \\ 0.281793 & 0.884131 & 0.969484 & 1.000000 \end{bmatrix}$$

$$\bar{\mathbf{W}} = \begin{bmatrix} 0.402694 & -0.345058 & 0.155350 & -0.241403 \\ -0.345058 & 0.339744 & -0.177349 & 0.255916 \\ 0.155350 & -0.177349 & 0.357751 & -0.087177 \\ -0.241403 & 0.255916 & -0.087177 & 0.227714 \end{bmatrix}$$

respectively. For the above matrix $\bar{\mathbf{A}}$, it follows from (13) that

$$\|\bar{\mathbf{A}}\bar{\mathbf{A}}^T - \bar{\mathbf{A}}^T\bar{\mathbf{A}}\|_F = 7.731601 \times 10^{-3}$$

which implies that matrix $\bar{\mathbf{A}}$ is practically normal. Moreover, the controllability Gramian $\bar{\mathbf{K}}$ satisfies l_2 -scaling constraints. For reference, the l_2 -sensitivity measure in (22) was computed as

$$S_2(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}) = 45.179954.$$

This reveals that minimization of weighted roundoff noise and pole sensitivity subject to l_2 -scaling constraints also reduce the l_2 -sensitivity for a state-space digital filter considerably.

The whole numerical results of the roundoff noise measure $J_R(\mathbf{T})$ in (10), the pole sensitivity measure $J_p(\mathbf{T})$ in (18), $J_R(\mathbf{T}) + J_p(\mathbf{T})$, and the objective function $J_\gamma(\mathbf{T})$ in (21) are summarized in Table 1.

Table 1. Performance comparison

γ	$J_\gamma(\mathbf{T})$	$J_R(\mathbf{T})$	$J_p(\mathbf{T})$	$J_R(\mathbf{T}) + J_p(\mathbf{T})$
1.0	4.000000	6.678752	4.000000	10.678752
0.9	3.765801	1.565948	4.010229	5.576178
0.8	3.513441	1.435279	4.032982	5.468261
0.7	3.246633	1.327903	4.068946	5.396849
0.6	2.965042	1.227214	4.123594	5.350808
0.5	2.666454	1.126471	4.206436	5.332908
0.4	2.347839	1.027726	4.328008	5.355734
0.3	2.004220	0.928102	4.515163	5.443265
0.2	1.625958	0.823537	4.835642	5.659179
0.1	1.189538	0.706245	5.539181	6.245426
0.0	0.555541	0.555541	13.604578	14.160119

5. CONCLUSION

The problem of minimizing a weighted roundoff noise and pole sensitivity measure subject to l_2 -scaling constraints for state-space filters has been investigated. An efficient iterative technique for minimizing a weighted roundoff noise and pole sensitivity measure subject to l_2 -scaling constraints has been developed by employing a quasi-Newton algorithm. Computer simulation results have demonstrated the validity and effectiveness of the proposed technique.

REFERENCES

- [1] S. Y. Hwang, "Roundoff noise in state-space digital filtering: A general analysis," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 24, pp. 256-262, June 1976.
- [2] 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.
- [3] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 25, pp. 273-281, Aug. 1977.
- [4] L. B. Jackson, A. G. Lindgren and Y. Kim, "Optimal synthesis of second-order state-space structures for digital filters," *IEEE Trans. Circuits Syst.*, vol. 26, pp. 149-153, Mar. 1979.
- [5] 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.
- [6] V. Tavsanoglu and L. Thiele, "Optimal design of state-space digital filters by simultaneous minimization of sensitivity and roundoff noise," *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp.884-888, Oct. 1984.
- [7] L. Thiele, "On the sensitivity of linear state-space systems," *IEEE Trans. Circuits Syst.*, vol.CAS-33, pp.502-510, May 1986.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] G. Li and M. Gevers, "Optimal synthetic FWL design of state-space digital filters," in *Proc. ICASSP 1992*, San Francisco, CA, August 24-28, 1992, vol.4, pp.429-432.
- [12] M. Gevers and G. Li, *Parameterizations in Control, Estimation and Filtering Problems: Accuracy Aspects*, Springer-Verlag, 1993.
- [13] T. Hinamoto, H. Ohnishi and W.-S. Lu, "Minimization of L_2 -sensitivity for state-space digital filters subject to L_2 -dynamic-range scaling constraints," *IEEE Trans. II*, vol.52, pp.641-645, Oct. 2005.
- [14] D. Williamson, "Roundoff noise minimization and pole-zero sensitivity in fixed-point digital filters using residue feedback," *IEEE Trans. Acoust. Speech, Signal Process.*, vol. ASSP-34, pp.1210-1220, Oct. 1986.
- [15] G. Li, "On pole and zero sensitivity of linear systems," *IEEE Trans. Circuits Syst. I*, vol. 44, pp.583-590, Jul. 1997.
- [16] R. Fletcher, *Practical Methods of Optimization*, 2nd ed. Wiley, New York, 1987.
- [17] T. Kailath, *Linear Systems*, Prentice Hall, 1980.