

JOINTLY OPTIMIZED ERROR-FEEDBACK AND REALIZATION FOR ROUND-OFF NOISE MINIMIZATION IN STATE-ESTIMATE FEEDBACK DIGITAL CONTROLLERS

*Takao Hinamoto, Keiji Kawai
and Masayoshi Nakamoto*

Graduate School of Engineering, Hiroshima University
Higashi-Hiroshima 739-8527, Japan
phone: +81-82-424-7672, fax: +81-82-422-7195
email: {hinamoto, msy}@hiroshima-u.ac.jp

Wu-Sheng Lu

Dept. of Elect. & Comput. Engineering,
University of Victoria
Victoria, B.C., V8W 3P6, Canada
phone: +1-250-721-8692, fax: +1-250-721-6052
email: wslu@ece.uvic.ca

ABSTRACT

The joint optimization problem of error-feedback and realization for the closed-loop system with a state-estimate feedback digital controller is investigated where the main objective is to minimize the effects of round-off noise at the closed-loop system output subject to l_2 -norm dynamic-range scaling constraints. 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. Analytical details are given as to how the proposed technique can be applied to the cases where the error-feedback matrix is a general, diagonal, or scalar matrix. A numerical example is presented to illustrate the utility of the proposed technique.

1. INTRODUCTION

Due to the finite precision nature of computer arithmetic, the output roundoff noise of a fixed-point IIR digital filter usually arises. This noise is critically dependent on the internal structure of an IIR digital filter [1],[2]. Error feedback (EF) is known as an effective technique for reducing the output roundoff noise in an IIR digital filter [3]-[5]. Williamson [6] has reduced the output roundoff noise more effectively by choosing the filter structure and applying EF to the filter. Lu and Hinamoto [7] have developed a jointly optimized technique of EF and realization to minimize the effects of roundoff noise at the filter output subject to l_2 -scaling constraints. Li and Gevers [8] have analyzed the output roundoff noise of the closed-loop system with a state-estimate feedback controller, and presented an algorithm for realizing the state-estimate feedback controller with minimum output roundoff noise under l_2 -norm dynamic-range scaling constraints. Hinamoto and Yamamoto [9] have proposed a method for applying EF to a given closed-loop system with a state-estimate feedback controller.

This paper investigates the problem of jointly optimizing EF and realization for the closed-loop system with a state-estimate feedback controller so as to

minimize the output roundoff noise subject to l_2 -norm dynamic-range scaling constraints. To this end, an iterative technique which relies on an efficient quasi-Newton algorithm [10] is developed. Our computer simulation results demonstrate the validity and effectiveness of the proposed technique.

2. ROUND-OFF NOISE ANALYSIS

Let a stable, controllable and observable linear discrete-time system be described by

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

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}_o , \mathbf{b}_o and \mathbf{c}_o are real constant matrices of appropriate dimensions. The transfer function of the linear system in (1) is given by

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

If a regulator is designed by using the full-order state observer, we obtain a state-estimate feedback controller as

$$\begin{aligned} \tilde{\mathbf{x}}(k+1) &= \mathbf{F}_o \tilde{\mathbf{x}}(k) + \mathbf{b}_o u(k) + \mathbf{g}_o y(k) \\ &= \mathbf{R}_o \tilde{\mathbf{x}}(k) + \mathbf{b}_o r(k) + \mathbf{g}_o y(k) \\ u(k) &= -\mathbf{k}_o \tilde{\mathbf{x}}(k) + r(k) \end{aligned} \quad (3)$$

where $\tilde{\mathbf{x}}(k)$ is an $n \times 1$ state-variable vector in the full-order state observer, \mathbf{g}_o is an $n \times 1$ gain vector chosen so that all the eigenvalues of $\mathbf{F}_o = \mathbf{A}_o - \mathbf{g}_o \mathbf{c}_o$ are inside the unit circle in the complex plane, \mathbf{k}_o is a $1 \times n$ state-feedback gain vector chosen so that each of the eigenvalues of $\mathbf{A}_o - \mathbf{b}_o \mathbf{k}_o$ is at a desirable location within the unit circle, $r(k)$ is a scalar reference signal, and $\mathbf{R}_o = \mathbf{F}_o - \mathbf{b}_o \mathbf{k}_o$.

Performing quantization before matrix-vector multiplication, we can express the finite-word-length (FWL) implementation of (3) with error feedback as

$$\begin{aligned} \hat{\mathbf{x}}(k+1) &= \mathbf{R} \mathbf{Q}[\hat{\mathbf{x}}(k)] + \mathbf{b} r(k) + \mathbf{g} y(k) + \mathbf{D} e(k) \\ u(k) &= -\mathbf{k} \mathbf{Q}[\hat{\mathbf{x}}(k)] + r(k) \end{aligned} \quad (4)$$

where $e(k) = \hat{\mathbf{x}}(k) - \mathbf{Q}[\hat{\mathbf{x}}(k)]$ and \mathbf{D} is an $n \times n$ error feedback matrix. All coefficient matrices \mathbf{R} , \mathbf{b} , \mathbf{g} and \mathbf{k}

are assumed to have an exact fractional B_c bit representation. The FWL state-variable vector $\hat{\mathbf{x}}(k)$ and signal $u(k)$ all have a B bit fractional representation, while the reference input $r(k)$ is a $(B - B_c)$ bit fraction. The vector quantizer $\mathbf{Q}[\cdot]$ in (4) rounds the B bit fraction $\tilde{\mathbf{x}}(k)$ to $(B - B_c)$ bits after completing the multiplications and additions, where the sign bit is not counted. It is assumed that the roundoff error $\mathbf{e}(k)$ can be modeled as a zero-mean noise process with covariance $\sigma^2 \mathbf{I}_n$.

The closed-loop control system consisting of the linear system in (1) and the state-estimate feedback controller in (4) is illustrated in Fig. 1. This closed-loop

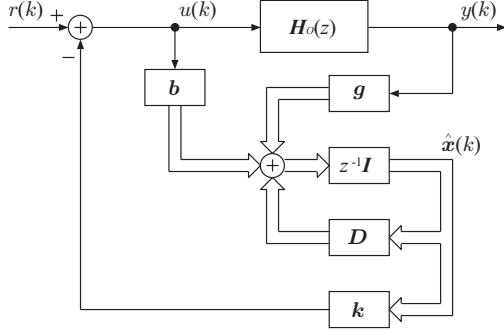


Fig. 1. The closed-loop control system with a state-estimate feedback controller.

system is described by

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \hat{\mathbf{x}}(k+1) \end{bmatrix} = \bar{\mathbf{A}} \begin{bmatrix} \mathbf{x}(k) \\ \hat{\mathbf{x}}(k) \end{bmatrix} + \bar{\mathbf{b}}r(k) + \bar{\mathbf{B}}\mathbf{e}(k) \quad (5)$$

$$y(k) = \bar{\mathbf{c}} \begin{bmatrix} \mathbf{x}(k) \\ \hat{\mathbf{x}}(k) \end{bmatrix}$$

where

$$\bar{\mathbf{A}} = \begin{bmatrix} \mathbf{A}_o & -\mathbf{b}_o \mathbf{k} \\ \mathbf{g} \mathbf{c}_o & \mathbf{R} \end{bmatrix}, \quad \bar{\mathbf{b}} = \begin{bmatrix} \mathbf{b}_o \\ \mathbf{b} \end{bmatrix}$$

$$\bar{\mathbf{B}} = \begin{bmatrix} \mathbf{b}_o \mathbf{k} \\ \mathbf{D} - \mathbf{R} \end{bmatrix}, \quad \bar{\mathbf{c}} = [\mathbf{c}_o \quad \mathbf{0}].$$

Let the transfer function from the roundoff noise $\mathbf{e}(k)$ to the output $y(k)$ in (5) be defined by $\mathbf{G}_D(z)$. Then

$$\mathbf{G}_D(z) = \bar{\mathbf{c}}(z\mathbf{I}_{2n} - \bar{\mathbf{A}})^{-1}\bar{\mathbf{B}}. \quad (6)$$

The noise gain $J(\mathbf{D}) = \sigma_{out}^2/\sigma^2$ is then computed as

$$J(\mathbf{D}) = \text{tr}[\mathbf{W}_D] \quad (7)$$

with

$$\mathbf{W}_D = \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{G}_D^*(z)\mathbf{G}_D(z) \frac{dz}{z} \quad (8)$$

where σ_{out}^2 stands for the noise variance of the output. For tractability, we evaluate $J(\mathbf{D})$ in (7) by replacing \mathbf{R} , \mathbf{b} , \mathbf{g} and \mathbf{k} by \mathbf{R}_o , \mathbf{b}_o , \mathbf{g}_o and \mathbf{k}_o , respectively. Define

$$\mathbf{S} = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{I}_n & -\mathbf{I}_n \end{bmatrix}, \quad (9)$$

the transfer function $\mathbf{G}_D(z)$ in (6) can be expressed as

$$\begin{aligned} \mathbf{G}_D(z) &= \bar{\mathbf{c}}\mathbf{S}(z\mathbf{I}_{2n} - \mathbf{S}^{-1}\bar{\mathbf{A}}\mathbf{S})^{-1}\mathbf{S}^{-1}\bar{\mathbf{B}} \\ &= \bar{\mathbf{c}}(z\mathbf{I}_{2n} - \Phi)^{-1} \begin{bmatrix} \mathbf{b}_o \mathbf{k}_o \\ \mathbf{F}_o - \mathbf{D} \end{bmatrix} \\ &= \mathbf{c}_o(z\mathbf{I}_n - \mathbf{A}_o + \mathbf{b}_o \mathbf{k}_o)^{-1} \mathbf{b}_o \mathbf{k}_o (z\mathbf{I}_n - \mathbf{F}_o)^{-1} \\ &\quad \cdot (z\mathbf{I}_n - \mathbf{D}) \\ &= \bar{\mathbf{c}}(z\mathbf{I}_{2n} - \Phi)^{-1} \mathbf{U}(z\mathbf{I}_n - \mathbf{D}) \end{aligned} \quad (10)$$

where

$$\Phi = \begin{bmatrix} \mathbf{A}_o - \mathbf{b}_o \mathbf{k}_o & \mathbf{b}_o \mathbf{k}_o \\ \mathbf{0} & \mathbf{F}_o \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_n \end{bmatrix}.$$

It is noted that the stability of the closed-loop control system is determined by the eigenvalues of matrix $\bar{\mathbf{A}}$ in (5), or equivalently, those of matrix Φ in (10). This means that neither of the quantization error $\mathbf{e}(k)$ and the error-feedback matrix \mathbf{D} affects the stability.

Substituting (10) into matrix \mathbf{W}_D in (8) gives

$$\begin{aligned} \mathbf{W}_D &= (\mathbf{b}_o \mathbf{k}_o)^T \mathbf{W}_1 \mathbf{b}_o \mathbf{k}_o + (\mathbf{b}_o \mathbf{k}_o)^T \mathbf{W}_2 (\mathbf{F}_o - \mathbf{D}) \\ &\quad + (\mathbf{F}_o - \mathbf{D})^T \mathbf{W}_3 \mathbf{b}_o \mathbf{k}_o \\ &\quad + (\mathbf{F}_o - \mathbf{D})^T \mathbf{W}_4 (\mathbf{F}_o - \mathbf{D}) \end{aligned} \quad (11)$$

where

$$\mathbf{W} = \Phi^T \mathbf{W} \Phi + \bar{\mathbf{c}}^T \bar{\mathbf{c}}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{W}_1 & \mathbf{W}_2 \\ \mathbf{W}_3 & \mathbf{W}_4 \end{bmatrix}.$$

Since \mathbf{W} is positive semidefinite, it can be shown that there exists an $n \times n$ matrix \mathbf{P} such that $\mathbf{W}_3 = \mathbf{W}_4 \mathbf{P}$. In addition, (11) can be written by virtue of $\mathbf{W}_2 = \mathbf{W}_3^T$ as

$$\begin{aligned} \mathbf{W}_D &= (\mathbf{F}_o + \mathbf{P} \mathbf{b}_o \mathbf{k}_o - \mathbf{D})^T \mathbf{W}_4 (\mathbf{F}_o + \mathbf{P} \mathbf{b}_o \mathbf{k}_o - \mathbf{D}) \\ &\quad + (\mathbf{b}_o \mathbf{k}_o)^T (\mathbf{W}_1 - \mathbf{P}^T \mathbf{W}_4 \mathbf{P}) \mathbf{b}_o \mathbf{k}_o. \end{aligned} \quad (12)$$

Alternatively, applying z -transform to the first equation in (5) under the assumption that $\mathbf{e}(k) = \mathbf{0}$, we obtain

$$\begin{bmatrix} \mathbf{X}(z) \\ \hat{\mathbf{X}}(z) \end{bmatrix} = (z\mathbf{I} - \bar{\mathbf{A}})^{-1} \bar{\mathbf{b}} R(z) \quad (13)$$

where $\mathbf{X}(z)$, $\hat{\mathbf{X}}(z)$ and $R(z)$ represent the z -transforms of $\mathbf{x}(k)$, $\hat{\mathbf{x}}(k)$ and $r(k)$, respectively. Replacing \mathbf{R} , \mathbf{b} , \mathbf{k} and \mathbf{g} by \mathbf{R}_o , \mathbf{b}_o , \mathbf{k}_o and \mathbf{g}_o , respectively, and then using

$$\mathbf{S}^{-1} \begin{bmatrix} \mathbf{X}(z) \\ \hat{\mathbf{X}}(z) \end{bmatrix} = (z\mathbf{I}_{2n} - \mathbf{S}^{-1}\bar{\mathbf{A}}\mathbf{S})^{-1} \mathbf{S}^{-1} \bar{\mathbf{b}}$$

yield

$$\hat{\mathbf{X}}(z) = \mathbf{X}(z) = \mathbf{F}(z)R(z) \quad (14)$$

where

$$\mathbf{F}(z) = [z\mathbf{I}_n - (\mathbf{A}_o - \mathbf{b}_o \mathbf{k}_o)]^{-1} \mathbf{b}_o.$$

The controllability Gramian \mathbf{K} defined by

$$\mathbf{K} = \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{F}(z)\mathbf{F}^*(z) \frac{dz}{z} \quad (15)$$

can be obtained by solving the Lyapunov equation

$$\mathbf{K} = (\mathbf{A}_o - \mathbf{b}_o\mathbf{k}_o)\mathbf{K}(\mathbf{A}_o - \mathbf{b}_o\mathbf{k}_o)^T + \mathbf{b}_o\mathbf{b}_o^T. \quad (16)$$

3. ROUND-OFF NOISE MINIMIZATION

Consider the system in (4) with $\mathbf{D} = \mathbf{0}$ and denote it by $(\mathbf{R}, \mathbf{b}, \mathbf{g}, \mathbf{k})_n$. By applying a coordinate transformation $\tilde{\mathbf{x}}'(k) = \mathbf{T}^{-1}\hat{\mathbf{x}}(k)$ to the above system $(\mathbf{R}, \mathbf{b}, \mathbf{g}, \mathbf{k})_n$, we obtain a new realization characterized by $(\tilde{\mathbf{R}}, \tilde{\mathbf{b}}, \tilde{\mathbf{g}}, \tilde{\mathbf{k}})_n$ where

$$\begin{aligned} \tilde{\mathbf{R}} &= \mathbf{T}^{-1}\mathbf{R}\mathbf{T}, & \tilde{\mathbf{b}} &= \mathbf{T}^{-1}\mathbf{b} \\ \tilde{\mathbf{g}} &= \mathbf{T}^{-1}\mathbf{g}, & \tilde{\mathbf{k}} &= \mathbf{k}\mathbf{T}. \end{aligned} \quad (17)$$

For the system in (17), the counterparts of \mathbf{W}_i for $i = 1, 2, 3, 4$ are given by

$$\tilde{\mathbf{W}}_i = \mathbf{T}^T\mathbf{W}_i\mathbf{T} \quad (18)$$

and the corresponding noise gain is given by

$$J(\mathbf{D}, \mathbf{T}) = \text{tr}[\tilde{\mathbf{W}}_D] \quad (19)$$

where $\tilde{\mathbf{W}}_D$ can be obtained using (11) as

$$\begin{aligned} \tilde{\mathbf{W}}_D &= [\mathbf{T}^{-1}(\mathbf{F}_0 + \mathbf{P}\mathbf{b}_0\mathbf{k}_0)\mathbf{T} - \mathbf{D}]^T \\ &\cdot \mathbf{T}^T\mathbf{W}_4\mathbf{T} [\mathbf{T}^{-1}(\mathbf{F}_0 + \mathbf{P}\mathbf{b}_0\mathbf{k}_0)\mathbf{T} - \mathbf{D}] \\ &+ \mathbf{T}^T(\mathbf{b}_0\mathbf{k}_0)^T(\mathbf{W}_1 - \mathbf{P}^T\mathbf{W}_4\mathbf{P})\mathbf{b}_0\mathbf{k}_0\mathbf{T}. \end{aligned}$$

In addition, (15) can be written as

$$\tilde{\mathbf{K}} = \mathbf{T}^{-1}\mathbf{K}\mathbf{T}^{-T}. \quad (20)$$

As a result, the output roundoff noise minimization problem amounts to obtaining matrices \mathbf{D} and \mathbf{T} which jointly minimize $J(\mathbf{D}, \mathbf{T})$ in (19) subject to the l_2 -scaling constraints specified by

$$(\tilde{\mathbf{K}})_{ii} = (\mathbf{T}^{-1}\mathbf{K}\mathbf{T}^{-T})_{ii} = 1, \quad i = 1, 2, \dots, n. \quad (21)$$

To deal with (21), we define

$$\hat{\mathbf{T}} = \mathbf{T}^T\mathbf{K}^{-\frac{1}{2}}. \quad (22)$$

Then the l_2 -scaling constraints in (21) can be written as

$$(\hat{\mathbf{T}}^{-T}\hat{\mathbf{T}}^{-1})_{ii} = 1, \quad i = 1, 2, \dots, n. \quad (23)$$

These constraints are always satisfied if $\hat{\mathbf{T}}^{-1}$ assumes the form

$$\hat{\mathbf{T}}^{-1} = \left[\frac{\mathbf{t}_1}{\|\mathbf{t}_1\|}, \frac{\mathbf{t}_2}{\|\mathbf{t}_2\|}, \dots, \frac{\mathbf{t}_n}{\|\mathbf{t}_n\|} \right]. \quad (24)$$

Substituting (22) into (19), we obtain

$$\begin{aligned} J(\mathbf{D}, \hat{\mathbf{T}}) &= \text{tr} \left[\hat{\mathbf{T}}(\hat{\mathbf{A}} - \hat{\mathbf{T}}^T\mathbf{D}\hat{\mathbf{T}}^{-T})^T\hat{\mathbf{W}}_4 \right. \\ &\cdot (\hat{\mathbf{A}} - \hat{\mathbf{T}}^T\mathbf{D}\hat{\mathbf{T}}^{-T})\hat{\mathbf{T}}^T + \hat{\mathbf{T}}\hat{\mathbf{C}}\hat{\mathbf{T}}^T \left. \right] \end{aligned} \quad (25)$$

where

$$\begin{aligned} \hat{\mathbf{A}} &= \mathbf{K}^{-\frac{1}{2}}(\mathbf{F}_0 + \mathbf{P}\mathbf{b}_0\mathbf{k}_0)\mathbf{K}^{\frac{1}{2}}, & \hat{\mathbf{W}}_4 &= \mathbf{K}^{\frac{1}{2}}\mathbf{W}_4\mathbf{K}^{\frac{1}{2}} \\ \hat{\mathbf{C}} &= \mathbf{K}^{\frac{1}{2}}(\mathbf{b}_0\mathbf{k}_0)^T(\mathbf{W}_1 - \mathbf{P}^T\mathbf{W}_4\mathbf{P})\mathbf{b}_0\mathbf{k}_0\mathbf{K}^{\frac{1}{2}}. \end{aligned}$$

From the foregoing arguments, the problem of obtaining matrices \mathbf{D} and \mathbf{T} that minimize (19) subject to the scaling constraints in (21) is now converted into an unconstrained optimization problem of obtaining \mathbf{D} and $\hat{\mathbf{T}}$ that jointly minimize $J(\mathbf{D}, \hat{\mathbf{T}})$ in (25).

Let \mathbf{x} be the column vector that collects the variables in matrices \mathbf{D} and $[\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n]$. Then $J(\mathbf{D}, \hat{\mathbf{T}})$ is a function of \mathbf{x} , denoted by $J(\mathbf{x})$. The proposed algorithm starts with an initial point \mathbf{x}_0 obtained from an initial assignment $\mathbf{D} = \hat{\mathbf{T}} = \mathbf{I}_n$. In the k th iteration, a quasi-Newton algorithm updates the most recent point \mathbf{x}_k to point \mathbf{x}_{k+1} as [10]

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k\mathbf{d}_k, \quad (26)$$

where

$$\mathbf{d}_k = -\mathbf{S}_k\nabla J(\mathbf{x}_k)$$

$$\alpha_k = \arg \left[\min_{\alpha} J(\mathbf{x}_k + \alpha\mathbf{d}_k) \right]$$

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

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 $J(\mathbf{x}_k)$. This iteration process continues until

$$|J(\mathbf{x}_{k+1}) - J(\mathbf{x}_k)| < \varepsilon \quad (27)$$

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

In what follows, we derive closed-form expressions of $\nabla J(\mathbf{x})$ for the cases where \mathbf{D} assumes the form of a general, diagonal, or scalar matrix.

1) *Case 1: \mathbf{D} Is a General Matrix:* From (25), the optimal choice of \mathbf{D} is given by

$$\mathbf{D} = \hat{\mathbf{T}}^{-T}\hat{\mathbf{A}}\hat{\mathbf{T}}^T, \quad (28)$$

which leads to

$$J(\hat{\mathbf{T}}^{-T}\hat{\mathbf{A}}\hat{\mathbf{T}}^T, \hat{\mathbf{T}}) = \text{tr} \left[\hat{\mathbf{T}}\hat{\mathbf{C}}\hat{\mathbf{T}}^T \right]. \quad (29)$$

In this case, the number of elements in vector \mathbf{x} consisting of $\hat{\mathbf{T}}$ is equal to n^2 and the gradient of $J(\mathbf{x})$ is found to be

$$\begin{aligned} \frac{\partial J(\mathbf{x})}{\partial t_{ij}} &= \lim_{\Delta \rightarrow 0} \frac{J(\hat{\mathbf{T}}_{ij}) - J(\hat{\mathbf{T}})}{\Delta} \\ &= 2\mathbf{e}_j^T\hat{\mathbf{T}}\hat{\mathbf{C}}\hat{\mathbf{T}}^T\hat{\mathbf{T}}\mathbf{g}_{ij}, \quad i, j = 1, 2, \dots, n \end{aligned} \quad (30)$$

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

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

and \mathbf{g}_{ij} is computed using

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

2) *Case 2: \mathbf{D} Is a Diagonal Matrix:* Here, matrix \mathbf{D} assumes the form

$$\mathbf{D} = \text{diag}\{d_1, d_2, \dots, d_n\}. \quad (31)$$

In this case, (25) becomes

$$J(\mathbf{D}, \hat{\mathbf{T}}) = \text{tr} \left[\hat{\mathbf{T}} \mathbf{M}_d \hat{\mathbf{T}}^T \right] \quad (32)$$

where

$$\begin{aligned} \mathbf{M}_d = & \hat{\mathbf{C}} + \hat{\mathbf{A}}^T \hat{\mathbf{W}}_4 \hat{\mathbf{A}} + \hat{\mathbf{W}}_4 \hat{\mathbf{T}}^T \mathbf{D}^2 \hat{\mathbf{T}}^{-T} \\ & - \hat{\mathbf{A}}^T \hat{\mathbf{W}}_4 \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T} - \hat{\mathbf{W}}_4 \hat{\mathbf{A}} \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T}. \end{aligned}$$

It follows that

$$\begin{aligned} \frac{\partial J(\mathbf{x})}{\partial t_{ij}} &= 2\mathbf{e}_j^T \hat{\mathbf{T}} \mathbf{M}_d \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij}, \quad i, j = 1, 2, \dots, n \\ \frac{\partial J(\mathbf{x})}{\partial d_i} &= 2\mathbf{e}_i^T (\mathbf{D} \hat{\mathbf{T}} - \hat{\mathbf{T}} \hat{\mathbf{A}}^T) \hat{\mathbf{W}}_4 \hat{\mathbf{T}}^T \mathbf{e}_i, \quad i = 1, 2, \dots, n. \end{aligned} \quad (33)$$

3) *Case 3: \mathbf{D} Is a Scalar Matrix:* It is assumed here that $\mathbf{D} = \alpha \mathbf{I}_n$ with a scalar α . The gradient of $J(\mathbf{x})$ can then be calculated as

$$\begin{aligned} \frac{\partial J(\mathbf{x})}{\partial t_{ij}} &= 2\mathbf{e}_j^T \hat{\mathbf{T}} \mathbf{M}_s \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij}, \quad i, j = 1, 2, \dots, n \\ \frac{\partial J(\mathbf{x})}{\partial \alpha} &= \text{tr} \left[\hat{\mathbf{T}} (2\alpha \hat{\mathbf{W}}_4 - \hat{\mathbf{A}}^T \hat{\mathbf{W}}_4 - \hat{\mathbf{W}}_4 \hat{\mathbf{A}}) \hat{\mathbf{T}}^T \right] \end{aligned} \quad (34)$$

where

$$\mathbf{M}_s = (\hat{\mathbf{A}} - \alpha \mathbf{I}_n)^T \hat{\mathbf{W}}_4 (\hat{\mathbf{A}} - \alpha \mathbf{I}_n) + \hat{\mathbf{C}}.$$

4. A NUMERICAL EXAMPLE

In this section we illustrate the proposed method by considering a linear discrete-time system specified by

$$\begin{aligned} \mathbf{A}_o &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0.339377 & -1.152652 & 1.520167 \end{bmatrix}, \quad \mathbf{b}_o = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \mathbf{c}_o &= [0.093253 \quad 0.128620 \quad 0.314713]. \end{aligned}$$

Suppose that the poles of the observer and regulator in the system are required to be located at $z = 0.1532$, 0.2861 , 0.1137 , and $z = 0.5067$, 0.6023 , 0.4331 , respectively. This can be achieved by choosing

$$\begin{aligned} \mathbf{k}_o &= [0.471552 \quad -0.367158 \quad 3.062267] \\ \mathbf{g}_o &= [-0.006436 \quad 3.683651 \quad 5.083920]^T. \end{aligned}$$

Performing the l_2 -scaling to the state-estimate feedback controller, we obtain $J(\mathbf{0}) = 686.4121$ in (7) where $\mathbf{D} = \mathbf{0}$. Next, the controller is transformed into the optimal realization that minimizes $J(\mathbf{0})$ in (7) under the l_2 -scaling constraints. This leads to $J_{\min}(\mathbf{0}) = 28.6187$. Finally, EF and state-variable coordinate transformation are applied to the above optimal realization so as to jointly minimize the output roundoff noise. The profiles of $J(\mathbf{x})$ during the first 20 iteration for the cases of \mathbf{D} being a general, diagonal, and scalar matrix are depicted in Fig. 2.

1) *Case 1: \mathbf{D} Is a General Matrix:* The quasi-Newton algorithm was applied to minimize (25). It took the algorithm 20 iterations to converge to the solution

$$\begin{aligned} \mathbf{D} &= \begin{bmatrix} 0.211191 & -3.078211 & -3.344596 \\ -1.321589 & 1.897308 & 3.243515 \\ 1.917916 & -1.890027 & -3.807473 \end{bmatrix} \\ \mathbf{T} &= \begin{bmatrix} -11.039974 & -43.683697 & -30.131793 \\ -3.231505 & 8.919473 & 9.118205 \\ 2.620911 & 6.462685 & 7.032260 \end{bmatrix} \end{aligned}$$

and the minimized noise gain was found to be $J(\mathbf{D}, \hat{\mathbf{T}}) = 4.8823$. Next, the above optimal EF matrix \mathbf{D} was rounded to a power-of-two representation with 3 bits after the binary point, which resulted in

$$\mathbf{D}_{3bit} = \begin{bmatrix} 0.250 & -3.125 & -3.375 \\ -1.375 & 1.875 & 3.250 \\ 1.875 & -1.875 & -3.750 \end{bmatrix}$$

and a noise gain $J(\mathbf{D}_{3bit}, \hat{\mathbf{T}}) = 23.4873$. Furthermore, when the optimal EF matrix \mathbf{D} was rounded to the integer representation

$$\mathbf{D}_{int} = \begin{bmatrix} 0 & -3 & -3 \\ -1 & 2 & 3 \\ 2 & -2 & -4 \end{bmatrix},$$

the noise gain was found to be $J(\mathbf{D}_{int}, \hat{\mathbf{T}}) = 293.0187$.

2) *Case 2: \mathbf{D} Is a Diagonal Matrix:* Again, the quasi-Newton algorithm was applied to minimize $J(\mathbf{D}, \hat{\mathbf{T}})$ in (25) for a diagonal EF matrix \mathbf{D} . It took the algorithm 20 iterations to converge to the solution

$$\begin{aligned} \mathbf{D} &= \text{diag}\{0.050638, -0.608845, -0.951572\} \\ \mathbf{T} &= \begin{bmatrix} 3.588878 & 0.735966 & 0.010417 \\ -2.457241 & 0.728171 & 0.556762 \\ 1.514232 & -2.058856 & 0.142204 \end{bmatrix} \end{aligned}$$

and the minimized noise gain was found to be $J(\mathbf{D}, \hat{\mathbf{T}}) = 12.7097$. Next, the above optimal diagonal EF matrix \mathbf{D} was rounded to a power-of-two representation with 3 bits after the binary point to yield $\mathbf{D}_{3bit} = \text{diag}\{0.000, -0.625, -1.000\}$, which leads to a noise gain $J(\mathbf{D}_{3bit}, \hat{\mathbf{T}}) = 12.7722$. Furthermore, when the optimized diagonal EF matrix \mathbf{D} was rounded to the integer representation $\mathbf{D}_{int} = \text{diag}\{0, -1, -1\}$, the noise gain was found to be $J(\mathbf{D}_{int}, \hat{\mathbf{T}}) = 13.7535$.

3) *Case 3: \mathbf{D} Is a Scalar Matrix:* In this case, the quasi-Newton algorithm was applied to minimize (25)

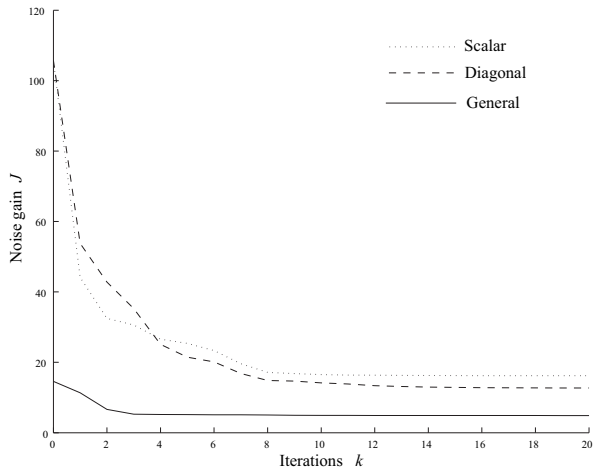


Fig. 2. Profiles of iterative noise gain minimization.

for $\mathbf{D} = \alpha \mathbf{I}_3$ with a scalar α . The algorithm converges after 20 iterations to converge to the solution

$$\mathbf{D} = -0.779678 \mathbf{I}_3$$

$$\mathbf{T} = \begin{bmatrix} 3.252790 & -0.081745 & -0.198376 \\ -1.717225 & 1.220068 & -0.792487 \\ 0.546599 & -0.854316 & 2.295944 \end{bmatrix}$$

and the minimized noise gain was found to be $J(\mathbf{D}, \hat{\mathbf{T}}) = 16.2006$. Next, the EF matrix $\mathbf{D} = \alpha \mathbf{I}_3$ was rounded to a power-of-two representation with 3 bits after the binary point as well as an integer representation. It was found that these representations were given by $\mathbf{D}_{3bit} = \text{diag}\{0.750, 0.750, 0.750\}$ and $\mathbf{D}_{int} = \text{diag}\{1, 1, 1\}$, respectively. The corresponding noise gains were obtained as $J(\mathbf{D}_{3bit}, \hat{\mathbf{T}}) = 16.2370$ and $J(\mathbf{D}_{int}, \hat{\mathbf{T}}) = 18.2063$, respectively.

The above simulation results in terms of noise gain $J(\mathbf{D}, \hat{\mathbf{T}})$ in (25) are summarized in Table 1. For comparison purpose, their counterparts obtained using the method in [9] are also included in the table, where the minimization of the roundoff noise was carried out using EF and state-variable coordinate transformation, but in a separate manner. From the table, it is observed that the proposed joint optimization offers improved reduction in roundoff noise gain for the cases of a scalar EF matrix and a diagonal EF matrix when compared with those obtained by using *separate* optimization. However, in the case of a general EF matrix, the optimal solution with infinite precision appears to be quite sensitive to the parameter perturbations.

More reduction of the noise gain might be possible by re-designing the coordinate transformation matrix \mathbf{T} for the optimally quantized \mathbf{D} .

5. CONCLUSION

The joint optimization problem of EF and realization to minimize the effects of roundoff noise of the closed-loop system with a state-estimate feedback controller subject to l_2 -scaling constraints has been investigated. The problem at hand has been converted into an unconstrained optimization problem by using linear alge-

Table 1: Noise gain $J(\mathbf{D}, \hat{\mathbf{T}})$ for different EF schemes.

Error-Feedback Scheme	Accuracy of \mathbf{D}		
	Infinite Precision	3 Bit Quantization	Integer Quantization
$\mathbf{D} = 0$	28.6187		
Separate	28.6187		
Scalar Separate [9]	20.1235	20.1810	26.0527
Scalar Joint	16.2006	16.2370	18.2063
Diagonal Separate [9]	16.4104	16.4547	17.4039
Diagonal Joint	12.7097	12.7722	13.7535
General Separate [9]	11.6352	11.7054	16.5814
General Joint	4.8823	23.4873	293.0187

braic techniques. An efficient quasi-Newton algorithm has been employed to solve the unconstrained optimization problem. The proposed technique has been applied to the cases where EF matrix is a general, diagonal, or scalar matrix. The effectiveness for the cases of a scalar EF matrix and a diagonal EF matrix compared with the existing method [9] has been illustrated by a numerical example.

REFERENCES

- [1] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, pp. 551-562, Sept. 1976.
- [2] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust. Speech, Signal Processing*, vol. ASSP-25, pp. 273-281, Aug. 1977.
- [3] W. E. Higgins and D. C. Munson, "Optimal and suboptimal error-spectrum shaping for cascade-form digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp. 429-437, May 1984.
- [4] T. I. Laakso and I. O. Hartimo, "Noise reduction in recursive digital filters using high-order error feedback," *IEEE Trans. Signal Processing*, vol. 40, pp. 1096-1107, May 1992.
- [5] P. P. Vaidyanathan, "On error-spectrum shaping in state-space digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-32, pp. 88-92, Jan. 1985.
- [6] D. Williamson, "Roundoff noise minimization and pole-zero sensitivity in fixed-point digital filters using residue feedback," *IEEE Trans. Acoust. Speech, Signal Processing*, vol. ASSP-34, pp. 1210-1220, Oct. 1986.
- [7] W.-S. Lu and T. Hinamoto, "Jointly optimized error-feedback and realization for roundoff noise minimization in state-space digital filters," *IEEE Trans. Signal Processing*, vol. 53, pp. 2135-2145, June 2005.
- [8] G. Li and M. Gevers, "Optimal finite precision implementation of a state-estimate feedback controller," *IEEE Trans. Circuits Syst.*, vol. CAS-37, pp. 1487-1498, Dec. 1990.
- [9] T. Hinamoto and S. Yamamoto, "Error spectrum shaping in closed-loop systems with state-estimate feedback controller," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'02)*, May 2002, vol. 1, pp. 289-292.
- [10] R. Fletcher, *Practical Methods of Optimization*, 2nd ed. New York: Wiley, 1987.