東北大学機関リポジトリ
**TOUR**
Tohoku University Repository

# On Hankel Singular Values and Reflected Zeros of Linear Dynamical Systems

[6] P. Gahinet and P. Apkarian, "A linear matrix inequality approach to $H_\infty$ control," *Int. J. Robust Nonlin. Control*, vol. 4, no. 4, pp. 421–448, 1994.

[7] A. Packard, "Gain scheduling via linear fractional transformations," *Syst. Control Lett.*, vol. 22, no. 2, pp. 79–92, 1994.

[8] G. E. Dullerud and S. Lall, "A new approach for analysis and synthesis of time-varying systems," *IEEE Trans. Automat. Control*, vol. 44, no. 8, pp. 1486–1497, Aug. 1999.

[9] H. S. Witsenhausen, "Separation of estimation and control for discrete time systems," *Proc. IEEE*, vol. 59, no. 11, pp. 1557–1566, Nov. 1971.

[10] J.-W. Lee and G. E. Dullerud, "Uniform stabilization of discrete-time switched and Markovian jump linear systems," *Automatica*, vol. 42, no. 2, pp. 205–218, 2006.

[11] J.-W. Lee and G. E. Dullerud, "Optimal disturbance attenuation for discrete-time switched and Markovian jump linear systems," *SIAM J. Control Optim.*, vol. 45, no. 4, pp. 1329–1358, 2006.

[12] J.-W. Lee and P. P. Khargonekar, "Optimal output regulation for discrete-time switched and markovian jump linear systems," *SIAM J. Control Optim.*, vol. 47, no. 1, pp. 40–72, 2008.

[13] C. Scherer, P. Gahinet, and M. Chilali, "Multiobjective output-feedback control via LMI optimization," *IEEE Trans. Automat. Control*, vol. 42, no. 7, pp. 896–911, Jul. 1997.

[14] I. Masubuchi, A. Ohara, and N. Suda, "LMI-based controller synthesis: A unified formulation and solution," *Int. J. Robust Nonlin. Control*, vol. 8, no. 8, pp. 669–686, 1998.

# On Hankel Singular Values and Reflected Zeros of Linear Dynamical Systems

Shunsuke Koshita, *Member, IEEE*, Masahide Abe, *Member, IEEE*, Masayuki Kawamata, *Senior Member, IEEE*, and Athanasios C. Antoulas, *Fellow, IEEE*

*Abstract*—This note discusses a relationship between the Hankel singular values and reflected zeros of linear systems. Our main result proves that the Hankel singular values of a linear continuous-time system increase (decrease) pointwise when one or more zeros of the transfer function are reflected with respect to the imaginary axis, that is, move from the left-(right-)half to the right-(left-)half of the complex plane. We also derive a similar result for linear discrete-time systems.

*Index Terms*—Hankel singular values, linear continuous-time system, linear discrete-time system, reflected zeros.

## I. Introduction

The study of the Hankel singular values of linear dynamical systems is an important subject since they play crucial roles in many fields of linear system theory. One of the well-known examples is approximation of dynamical systems such as balanced model reduction and Hankel norm approximation [1]–[4], where the Hankel singular values give a priori theoretical upper bound of the infinity norm of approxi-

mation error. Other practically important issues can be seen in the field of signal processing theory, where the Hankel singular values are referred to as the second-order modes. They provide the optimal dynamic range of analog filters [5], [6], i.e. the highest ratio of the maximal and minimal signal levels that can be processed in the filters. Also, in the literature on digital signal processing, it is well known that the Hankel singular values characterize the minimum attainable value of roundoff noise [7], [8] and statistical coefficient sensitivity [9], [10] of digital filters.

Our main result is to derive a pointwise inequality that relates the Hankel singular values to reflection of the zeros of the transfer function. For linear continuous-time systems we establish the fact that, when one or more zeros of the transfer function are reflected with respect to the imaginary axis, the Hankel singular values increase or decrease pointwise. We also derive a similar result for linear discrete-time systems: we show that, the Hankel singular values of a discrete-time system increase or decrease pointwise when one or more zeros of the transfer function are relfected with respect to the unit circle. Although a part of these topics is also mentioned in [11], our result to be presented in this note will offer further significant insights into the linear system theory. Details on the contribution of this note with respect to [11] will be discussed in Section IV.

Throughout this note, we will use the following notations. $\mathbb{R}, \mathbb{C}, \mathbb{Z}$ denote the sets of real numbers, complex numbers and integers, respectively. $\mathbb{R}^{m \times n}$ and $\mathbb{C}^{m \times n}$ respectively denote the sets of $m \times n$ real matrices and $m \times n$ complex matrices. $\boldsymbol{A}^T$ and $\boldsymbol{A}^*$ respectively stand for the transpose and the complex conjugate transpose of a matrix $\boldsymbol{A}$. The symbol $\lambda_i(\boldsymbol{A})$ for $i = 1, 2, \cdots, n$ denotes the eigenvalues of $\boldsymbol{A} \in \mathbb{R}^{n \times n}$. When the eigenvalues are all real, they are always arranged in decreasing order, i.e. $\lambda_1(\boldsymbol{A}) \geq \lambda_2(\boldsymbol{A}) \geq \cdots \geq \lambda_n(\boldsymbol{A})$.

## II. Hankel Singular Values

We consider an $N$-th order single-input/single-output linear continuous-time or discrete-time system described by

$$\dot{\boldsymbol{x}}(t) = \boldsymbol{A}\boldsymbol{x}(t) + \boldsymbol{b}u(t), \quad y(t) = \boldsymbol{c}\boldsymbol{x}(t) + du(t), \quad t \in \mathbb{R}, \quad (1)$$

or $\boldsymbol{x}(t+1) = \boldsymbol{A}\boldsymbol{x}(t) + \boldsymbol{b}u(t), \quad y(t) = \boldsymbol{c}\boldsymbol{x}(t) + du(t), \quad t \in \mathbb{Z}.$ (2)

In both cases, $u(t) \in \mathbb{R}^{1 \times 1}$ is the input, $y(t) \in \mathbb{R}^{1 \times 1}$ is the output, $\boldsymbol{x}(t) \in \mathbb{R}^{N \times 1}$ is the state, and $\boldsymbol{A} \in \mathbb{R}^{N \times N}, \boldsymbol{b} \in \mathbb{R}^{N \times 1}, \boldsymbol{c} \in \mathbb{R}^{1 \times N}$ and $d \in \mathbb{R}^{1 \times 1}$ are real coefficients. The transfer function is represented as

$$G(\xi) = d + \boldsymbol{c}(\xi \boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b} \quad (3)$$

where $\xi = s$ (Laplace transform) for continuous-time systems and $\xi = z$ ($\mathcal{Z}$-transform) for discrete-time systems. Throughout this note, the system $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{c}, d)$ is assumed to be asymptotically stable, controllable, and observable.

For continuous-time systems, the solutions $\boldsymbol{P}$ and $\boldsymbol{Q}$ to the following Lyapunov equations are called the controllability Gramian and the observability Gramian, respectively:

$$\boldsymbol{A}\boldsymbol{P} + \boldsymbol{P}\boldsymbol{A}^T = -\boldsymbol{b}\boldsymbol{b}^T, \quad \boldsymbol{A}^T\boldsymbol{Q} + \boldsymbol{Q}\boldsymbol{A} = -\boldsymbol{c}^T\boldsymbol{c}. \quad (4)$$

In the discrete-time case, the controllability and observability Gramians are given from the following Lyapunov equations:

$$\boldsymbol{P} - \boldsymbol{A}\boldsymbol{P}\boldsymbol{A}^T = \boldsymbol{b}\boldsymbol{b}^T, \quad \boldsymbol{Q} - \boldsymbol{A}^T\boldsymbol{Q}\boldsymbol{A} = \boldsymbol{c}^T\boldsymbol{c}. \quad (5)$$

In both cases, $\boldsymbol{P}$ and $\boldsymbol{Q}$ are symmetric and positive definite, i.e. $\boldsymbol{P} = \boldsymbol{P}^T > 0$ and $\boldsymbol{Q} = \boldsymbol{Q}^T > 0$, because the system $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{c}, d)$ is assumed to be asymptotically stable, controllable and observable.

The Hankel singular values of $G(\xi)$, which are denoted by $\sigma_1, \sigma_2, \cdots, \sigma_N$, are defined as the positive square roots of the eigenvalues of the matrix product $\boldsymbol{PQ}$, i.e.

$$\sigma_i = \sqrt{\lambda_i(\boldsymbol{PQ})}, \quad 1 \leq i \leq N. \tag{6}$$

Throughout this note, we assume that the Hankel singular values are arranged in decreasing order, i.e. $\sigma_1 \geq \sigma_2 \geq \cdots \sigma_N$. It is well known that the Hankel singular values are invariant under similarity transformations of the state. This means that the Hankel singular values depend only on the transfer function.

## III. PROBLEM STATEMENT

The purpose of this note is to reveal a relationship between the Hankel singular values and reflected zeros of the transfer function in an explicit form. To this end, we first consider the following minimum phase transfer function

$$G_{\min \mathrm{p}}(\xi) = K \frac{(\xi - z_1)(\xi - z_2)\cdots(\xi - z_M)}{(\xi - p_1)(\xi - p_2)\cdots(\xi - p_N)}, \quad M \leq N \tag{7}$$

where $K$ is a real nonzero constant, $p_k$ and $z_l$ for $1 \leq k \leq N$ and $1 \leq l \leq M$ are respectively the poles and zeros of $G(\xi)$. Since $G_{\min \mathrm{p}}(\xi)$ is a stable and minumum phase system, $\mathrm{Re}(p_k), \mathrm{Re}(z_l) < 0$ for continuous-time case and $|p_k|, |z_l| < 1$ for discrete-time case. From (7), we define non-minimum phase transfer functions as

$$G_{\mathrm{np}}(s) = G_{\min \mathrm{p}}(s)|_{(s-z_l)\leftarrow(s+z_l^*)}$$
$$\text{(continuous} - \text{time case)} \tag{8}$$
$$G_{\mathrm{np}}(z) = G_{\min \mathrm{p}}(z)|_{(z-z_l)\leftarrow(1-z_l^*z)}$$
$$\text{(discrete} - \text{time case)} \tag{9}$$

for some $l$. In the continuous-time case, the replacement $(s - z_l) \leftarrow (s + z_l^*)$ means that $z_l$'s of $G_{\min \mathrm{p}}(s)$ are reflected with respect to the imaginary axis. In the discrete-time case, the replacement $(z - z_l) \leftarrow (1 - z_l^*z)$ means that $z_l$'s of $G_{\min \mathrm{p}}(z)$ are reflected with respect to the unit circle.[1] If this replacement is carried out for all $l$, the resultant system becomes the maximum phase system, which is denoted by $G_{\max \mathrm{p}}(\xi)$. Throughout this note, it is assumed that no pole-zero cancellations occur in the family of $G_{\min \mathrm{p}}(s)$, $G_{\max \mathrm{p}}(s)$ and $G_{\mathrm{np}}(s)$ and the family of $G_{\min \mathrm{p}}(z)$, $G_{\max \mathrm{p}}(z)$ and $G_{\mathrm{np}}(z)$. Also, it is easy to see that these two families have identical magnitude responses: $|G_{\min \mathrm{p}}(j\Omega)| = |G_{\max \mathrm{p}}(j\Omega)| = |G_{\mathrm{np}}(j\Omega)|$ holds in the continuous-time case, and $|G_{\min \mathrm{p}}(e^{j\omega})| = |G_{\max \mathrm{p}}(e^{j\omega})| = |G_{\mathrm{np}}(e^{j\omega})|$ holds in the discrete-time case.

Our contribution in this note is to derive a pointwise inequality of the Hankel singular values of $G_{\min \mathrm{p}}(\xi)$, $G_{\max \mathrm{p}}(\xi)$ and $G_{\mathrm{np}}(\xi)$. This result will be presented in the next two sections.

## IV. HANKEL SINGULAR VALUES AND REFLECTED ZEROS

Our main result will be derived from a description of the Gramians of spectral factors using the bounded-real Riccati equations. Although this approach can be applied to limited classes of transfer functions, the same conclusion can be derived for the other classes of transfer functions by making use of bilinear transformation or frequency transformation.

Before showing our main result, we first need to give some mathematical preliminaries on these concepts.

---

[1]Letting $z_l = r_l e^{j\omega_l}$ with $0 < r_l < 1$, we know that the roots of $1 - z_l^*z$ are given as $e^{j\omega_l}/r_l$, i.e. the reciprocal conjugate of $z_l$. Hence the roots of $1 - z_l^*z$ can be interpreted as the zeros that are reflected from $z_l$ with respect to the unit circle.

### A. Preliminaries

We first discuss state-space description of spectral factors using the bounded-real Riccati equation for continuous-time systems. This is summarized in the following two lemmas.

*Lemma 1 ([12]–[14]):* Let a continuous-time system $G(s) = d + \boldsymbol{c}(s\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b}$ be bounded-real, i.e. let $|G(j\Omega)|^2 \leq \gamma^2$ for all $\Omega$ and some real constant $\gamma$. Also, define $w^2 = \gamma^2 - d^2$ and assume $w^2 > 0$. Then, there exists a positive definite symmetric matrix $\boldsymbol{Y}$ satisfying the following bounded-real Riccati equation:

$$\boldsymbol{A}^T\boldsymbol{Y} + \boldsymbol{Y}\boldsymbol{A} + \boldsymbol{c}^T\boldsymbol{c} + (\boldsymbol{Y}\boldsymbol{b} + \boldsymbol{c}^Td)w^{-2}(\boldsymbol{b}^T\boldsymbol{Y} + \boldsymbol{c}d) = \boldsymbol{0}. \tag{10}$$

Any solution $\boldsymbol{Y}$ to (10) lies between two external solutions, i.e. $0 < \boldsymbol{Y}_{\min} \leq \boldsymbol{Y} \leq \boldsymbol{Y}_{\max}$. The matrix $\boldsymbol{Y}_{\min}$ is the unique solution to (10) such that the eigenvalues of $\boldsymbol{A} + \boldsymbol{b}w^{-2}(\boldsymbol{b}^T\boldsymbol{Y} + \boldsymbol{c}d)$ are all in the left-half plane. The matrix $\boldsymbol{Y}_{\max}$ is the unique solution to (10) such that the eigenvalues of $\boldsymbol{A} + \boldsymbol{b}w^{-2}(\boldsymbol{b}^T\boldsymbol{Y} + \boldsymbol{c}d)$ are all in the right-half plane.

*Lemma 2 ([13], [14]):* Let $G(s) = d + \boldsymbol{c}(s\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b}$ be bounded-real and let $\boldsymbol{Y}$ be a solution to (10). Also, define $\boldsymbol{l} = -(\boldsymbol{b}^T\boldsymbol{Y} + \boldsymbol{c}d)w^{-1}$. Then, the system $\overline{G}(s) = w + \boldsymbol{l}(s\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b}$ is a spectral factor of $G(s)$, i.e. $|G(j\Omega)|^2 + |\overline{G}(j\Omega)|^2 = \gamma^2$ holds for all $\Omega$. The system $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{l}, w)$ in $\overline{G}(s)$ is asymptotically stable and controllable because $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{c}, d)$ in $G(s)$ is assumed to be asymptotically stable and controllable.

Spectral factors for discrete-time systems can be also described in state-space form in a similar manner to the continuous-time case. Such a description is provided by the following lemma, which is obtained as a consequence of [15], [16].

*Lemma 3:* Let $G(z) = d + \boldsymbol{c}(z\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b}$ be bounded-real, i.e. $|G(e^{j\omega})|^2 \leq \gamma^2$ for all $\omega$. Then, there exists a positive definite symmetric matrix $\boldsymbol{Y}$ satisfying the following discrete-time bounded-real Riccati equation:

$$\boldsymbol{Y} - \boldsymbol{A}^T\boldsymbol{Y}\boldsymbol{A} - \boldsymbol{c}^T\boldsymbol{c} - (\boldsymbol{A}^T\boldsymbol{Y}\boldsymbol{b} + \boldsymbol{c}^Td)w^{-2}(\boldsymbol{b}^T\boldsymbol{Y}\boldsymbol{A} + \boldsymbol{c}d) = \boldsymbol{0} \tag{11}$$

where $w^{-2} = \gamma^2 - d^2 - \boldsymbol{b}^T\boldsymbol{Y}\boldsymbol{b}$. Any solution $\boldsymbol{Y}$ to (11) satisfies $0 < \boldsymbol{Y}_{\min} \leq \boldsymbol{Y} \leq \boldsymbol{Y}_{\max}$, where $\boldsymbol{Y}_{\min}$ and $\boldsymbol{Y}_{\max}$ are the external solutions such that the eigenvalues of $\boldsymbol{A} + \boldsymbol{b}w^{-2}(\boldsymbol{b}^T\boldsymbol{Y}\boldsymbol{A} + \boldsymbol{c}d)$ are all inside the unit circle and all outside the unit circle, respectively. Also, let $\boldsymbol{l} = -(\boldsymbol{b}^T\boldsymbol{Y}\boldsymbol{A} + \boldsymbol{c}d)w^{-1}$ and consider the system $\overline{G}(z) = w + \boldsymbol{l}(z\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b}$. Then, $\overline{G}(z)$ is a spectral factor of $G(z)$, i.e. $|G(e^{j\omega})|^2 + |\overline{G}(e^{j\omega})|^2 = \gamma^2$ holds for all $\omega$, and $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{l}, w)$ is asymptotically stable and controllable.

In addition to the above well-known theory, we give the following lemma that offers a simple description of the Gramians of spectral factors.

*Lemma 4:* Let $(\boldsymbol{P}, \boldsymbol{Q})$ be the controllability and observability Gramians of a bounded-real system $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{c}, d)$ with the transfer function $G(\xi)$. Also, let $(\overline{\boldsymbol{P}}, \overline{\boldsymbol{Q}})$ be the controllability and observability Gramians of $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{l}, w)$ with the transfer function $\overline{G}(\xi)$. Then, the following equations hold:

$$\overline{\boldsymbol{P}} = \boldsymbol{P} \tag{12}$$
$$\overline{\boldsymbol{Q}} + \boldsymbol{Q} = \boldsymbol{Y} \tag{13}$$

where $\boldsymbol{Y}$ is a solution to (10) or (11).

*Proof:* Here we give the proof for continuous-time systems. The proof in the discrete-time case can be derived in a similar way and is omitted for brevity.

Eq. (12) is trivial from the state-space representations of $G(s)$ and $\overline{G}(s)$.

Eq. (13) is proved as follows. Since $(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{c}, d)$ is assumed to be asymptotically stable, a solution $\boldsymbol{Y} = \boldsymbol{Y}^T > 0$ to (10) can be represented as

$$\boldsymbol{Y} = \boldsymbol{U} + \boldsymbol{V} \tag{14}$$

where $U = U^T > 0$ and $V = V^T \geq 0$ are given as the solutions to the following Lyapunov equations:

$$A^T U + UA = -c^T c \qquad (15)$$
$$A^T V + VA = -(Yb + c^T d)w^{-2}(b^T Y + cd). \qquad (16)$$

From (4) and (15), it is obvious that $U = Q$. Moreover, from (16) and $l = -(b^T Y + cd)w^{-1}$, it follows that $V = \overline{Q}$. These relationships show (13).   ∎

Before leaving this subsection, we introduce some well-known properties of bilinear transformation and frequency transformation.

*Lemma 5:* Consider a continuous-time transfer funtion $G(s)$ and the discrete-time system $H(z)$ given by the bilinear transformation

$$H(z) = G\left(\frac{z-1}{z+1}\right). \qquad (17)$$

Then, the followings hold:
1) The poles and zeros at $\mathrm{Re}(s) < 0$, $\mathrm{Re}(s) = 0$, and $\mathrm{Re}(s) > 0$ are respectively mapped into $|z| < 1$, $|z| = 1$, and $|z| > 1$.
2) Let a pair of reflected zeros of $G(s)$ be $z_G$ and $-z_G^*$, and let $z_H$ and $z_H'$ be respectively the zeros mapped from $z_G$ and $-z_G^*$ by the bilinear transformation. Then, $z_H' = 1/z_H^*$ holds, i.e. $z_H'$ is given by reflecting $z_H$ with respect to the unit circle.
3) The Hankel singular values of $H(z)$ are the same as those of $G(s)$.

*Lemma 6:* Consider a discrete-time transfer funtion $G(z)$ and another discrete-time transfer function $H(z)$ that is given by the following first-order frequency transformation

$$H(z) = G\left(\frac{z-\alpha}{1-\alpha z}\right) \qquad (18)$$

where $\alpha \in \mathbb{R}$ is an arbitrary constant satisfying $0 < |\alpha| < 1$. Then, the followings hold:
1) The poles and zeros of $G(z)$ at $|z| < 1$, $|z| = 1$, and $|z| > 1$ are respectively mapped into $|z| < 1$, $|z| = 1$, and $|z| > 1$.
2) Let a pair of reflected zeros of $G(z)$ be $z_G$ and $1/z_G^*$, and let $z_H$ and $z_H'$ be respectively the zeros mapped from $z_G$ and $1/z_G^*$ by the frequency transformation. Then, $z_H' = 1/z_H^*$ holds. In particular, $z = 0$ and $z = \infty$ are respectively mapped into $z = \alpha$ and $z = 1/\alpha$.
3) The Hankel singular values of $H(z)$ are the same as those of $G(z)$.

*Remark 1:* The proof of the statement 3) of Lemma 6 is given in [17].

### B. Main Result

We now reveal the relationship between the Hankel singular values and reflected zeros of the transfer function.

First, we consider the case where $M = N$ holds in $G(\xi)$, i.e. $G(\xi)$ is non-strictly proper. For this transfer function, the Hankel singular values and reflected zeros are related by the following proposition.

*Proposition 1:* Let $G_{\min p}(\xi)$, $G_{\max p}(\xi)$ and $G_{np}(\xi)$ be the linear continuous-time or discrete-time transfer functions that are respectively defined in Section III, and assume that $d \neq 0$. In the discrete-time case, it is also assumed that these transfer functions have no zeors at $z = 0$. Now, let $\sigma_{\min p,k}$, $\sigma_{\max p,k}$ and $\sigma_{np,k}$ for $1 \leq k \leq N$ be the Hankel singular values of $G_{\min p}(\xi)$, $G_{\max p}(\xi)$ and $G_{np}(\xi)$, respectively. Then, the following pointwise inequality holds for all $k$:

$$\sigma_{\min p,k} \leq \sigma_{np,k} \leq \sigma_{\max p,k}. \qquad (19)$$

*Proof:* Here we give the proof for linear continuous-time systems. By Lemma 2, there exists an asymptotically stable system $F(s) = d_F + c_F(sI - A_F)^{-1}b_F$ such that $|F(j\Omega)|^2 + |G_{\min p}(j\Omega)|^2 = |F(j\Omega)|^2 + |G_{\max p}(j\Omega)|^2 = |F(j\Omega)|^2 + |G_{np}(j\Omega)|^2 = \gamma^2$ for some

$\gamma$ and all $\Omega$. Using this relationship allows us to describe $G_{\min p}(s)$, $G_{\max p}(s)$ and $G_{np}(s)$ of the form

$$G_{\min p}(s) = w_F + l_{F\min p}(sI - A_F)^{-1}b_F$$
$$G_{\max p}(s) = w_F + l_{F\max p}(sI - A_F)^{-1}b_F$$
$$G_{np}(s) = w_F + l_{Fnp}(sI - A_F)^{-1}b_F \qquad (20)$$

where

$$w_F = \pm\sqrt{\gamma^2 - d_F^2}$$
$$l_{F\min p} = -\left(b_F^T Y_{F\min} + c_F d_F\right)w_F^{-1}$$
$$l_{F\max p} = -\left(b_F^T Y_{F\max} + c_F d_F\right)w_F^{-1}$$
$$l_{Fnp} = -\left(b_F^T Y_F + c_F d_F\right)w_F^{-1} \qquad (21)$$

and $Y_{F\min}$, $Y_{F\max}$, and $Y_F$ are the solutions to (10) for $(A_F, b_F, c_F, d_F)$ such that $Y_{F\min} \leq Y_F \leq Y_{F\max}$. Now, let $(P_{\min p}, Q_{\min p})$, $(P_{\max p}, Q_{\max p})$ and $(P_{np}, Q_{np})$ be the controllability and observability Gramians of $(A_F, b_F, l_{F\min p}, w_F)$, $(A_F, b_F, l_{F\max p}, w_F)$ and $(A_F, b_F, l_{Fnp}, w_F)$, respectively. Then, from Lemma 4 and the fact that $Y_{F\min} \leq Y_F \leq Y_{F\max}$, it immediately follows that

$$P_{\min p} = P_{np} = P_{\max p}, \quad Q_{\min p} \leq Q_{np} \leq Q_{\max p}. \qquad (22)$$

Hence $\lambda_k(P_{\min p}Q_{\min p}) \leq \lambda_k(P_{np}Q_{np}) \leq \lambda_k(P_{\max p}Q_{\max p})$ is derived for $1 \leq k \leq N$ and this shows $\sigma_{\min p,k} \leq \sigma_{np,k} \leq \sigma_{\max p,k}$.

The proof for discrete-time systems can be derived in a similar way with the help of Lemma 3 and is omitted here.   ∎

*Remark 2:* In the above proof, we have used the fact that for $K, L, M > 0$, the relationship $L \leq M$ implies $\lambda_i(KL) \leq \lambda_i(KM)$, $i = 1, \cdots, n$. This fact can be easily shown by deriving $K^{1/2}LK^{1/2} \leq K^{1/2}MK^{1/2}$ and using the facts that $\lambda_i(L) \leq \lambda_i(M)$ and $\lambda_i(RS) = \lambda_i(SR)$ for $R, S > 0$.

*Remark 3:* In the above proof, the assumption of $d \neq 0$ is required because the family of transfer functions $G_{\min p}(\xi)$, $G_{\max p}(\xi)$ and $G_{np}(\xi)$ are given as spectral factors of $F(\xi)$. If $d = 0$, these spectral factors cannot be obtained by the bounded-real Riccati equations. In the discrete-time case, we have imposed another restriction that $G(z)$ has no zeros at $z = 0$ because reflecting a zero at $z = 0$ with respect to the unit circle yields a zero at $z = \infty$, which results in $d = 0$.

As stated in Remark 3, Proposition 1 is applicable to only non-strictly proper transfer functions (with no zeros at $z = 0$ in the discrete-time case). However, for strictly proper transfer functions (and non-strictly proper transfer functions with some zeros at $z = 0$ in the discrete-time case), we can derive the same pointwise inequality as in Proposition 1. This fact can be easily shown by using bilinear transformation or discrete-time frequency transformation. A detailed discussion now follows.

First, we consider the discrete-time case. From Lemma 6 it is obvious that, given a strictly proper transfer function $G(z)$, the function $H(z)$ generated by (18) becomes non-strictly proper. In addition, we can generate $H(z)$ with no zeros at $z = 0$ by taking an appropriate $\alpha$. Furthermore, Lemma 6 shows that $H(z)$ has the same Hankel singular values and the same relationship with respect to reflected zeros as those of $G(z)$. Consequently, even when $G(z)$ is strictly proper or $G(z)$ has some zeros at $z = 0$, the pointwise inequality (19) can be derived for $G(z)$ by applying Proposition 1 to $H(z)$.

We next discuss the continuous-time case. Let $G(s)$ be strictly proper and $H(z)$ be obtained by the bilinear transformation (17). Then, as is well-known, $H(z)$ becomes non-strictly proper in most cases. From this fact and Lemma 5, it immediately follows that (19) holds for most of the strictly proper continuous-time transfer functions. Even if $H(z)$ becomes strictly proper or $H(z)$ has some

zeros at $z = 0$, we can easily arrive at (19) by applying an appropriate frequency transformation to such $H(z)$.

As a result of the above discussion, we finally present the following theorem.

*Theorem 1:* Let $G_{\min p}(\xi)$, $G_{\max p}(\xi)$ and $G_{np}(\xi)$ be the linear dynamical systems that are respectively defined in Section III. Also, let $\sigma_{\min p,k}$, $\sigma_{\max p,k}$ and $\sigma_{np,k}$ for $1 \leq k \leq N$ be the Hankel singular values of $G_{\min p}(\xi)$, $G_{\max p}(\xi)$ and $G_{np}(\xi)$, respectively. Then, the following pointwise inequality holds for all $k$:

$$\sigma_{\min p,k} \leq \sigma_{np,k} \leq \sigma_{\max p,k}. \qquad (23)$$

*Remark 4:* In [11], it is stated that the minimum phase factors in spectral factorization have the smallest Hankel singular values. Although [11] does not give the proof of this fact, it is mentioned in [11] that the use of Nehari's theorem [18] will provide the proof of this fact. Taking this into account, it seems that our result is considered as an alternative proof of the fact mentioned in [11]. However, our result provides another significant insight into the linear system theory in that the relationship between reflected zeros and the Gramians of systems has been described in simple and explicit form. Also, in [11] only the strictly proper spectral factors are discussed, whereas our result given by Theorem 1 clearly applies to both strictly proper and non-strictly proper transfer functions. In particular, it should be stressed that Theorem 1 includes the case of reflected zeros at $z = 0$ and $z = \infty$ for discrete-time systems, which is not discussed in [11]. Furthermore, the discussion to be presented in the next section will give further insights into the topic of this note, in that it explicitly describes a direct relationship of reflected zeros to the Hankel singular values.

## V. HANKEL SINGULAR VALUES AND REFLECTED ZEROS: DETAILED ANALYSIS FOR STRICTLY PROPER CONTINUOUS-TIME SYSTEMS

In this section, we restrict ourselves to strictly proper continuous-time transfer functions with simple poles, and derive the pointwise inequality on the Hankel singular values and reflected zeros by a different approach. Although the conclusion is the same as in the previous section, the analysis to be presented here gives further insights into the topic of this note, in that the analysis provides explicit formulation of the Gramians of systems in terms of reflected zeros.

### A. Preliminaries

Given a strictly proper transfer function $G(s) = \boldsymbol{c}(s\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b}$, let one of its zeros be denoted by $\zeta$. Our main result will show that if this zero is reflected with respect to the imaginary axis, i.e. becomes $-\zeta^*$, depending on whether $\mathrm{Re}(\zeta) < 0$ or $\mathrm{Re}(\zeta) > 0$, the $k$-th Hankel singular value of the resulting system is bigger or smaller than or equal, respectively, to the $k$-th singular value of the original system, for all $k$.

Towards this goal we will make use of the partial fraction expansion of the transfer function

$$G(s) = \boldsymbol{c}(s\boldsymbol{I} - \boldsymbol{A})^{-1}\boldsymbol{b} = \sum_{k=1}^{N} \frac{\beta_k}{s + \alpha_k}, \quad \mathrm{Re}(\zeta_k) > 0, \qquad (24)$$

where we assume for simplicity that all poles are simple.

A state-space representation of this system is given by[2]

$$\boldsymbol{A} = \mathrm{diag}(-\alpha_1, \cdots, -\alpha_N), \quad \boldsymbol{b} = [1, \cdots, 1]^T$$
$$\boldsymbol{c} = [\beta_1, \cdots, \beta_N]. \qquad (25)$$

[2]For analysis purpose, in this section we assume that state-space coefficients may be complex, i.e. $\boldsymbol{A} \in \mathbb{C}^{N \times N}$, $\boldsymbol{b} \in \mathbb{C}^{N \times 1}$ and $\boldsymbol{c} \in \mathbb{C}^{1 \times N}$. Accordingly, the Gramians $\boldsymbol{P}$ and $\boldsymbol{Q}$ are assumed to be complex, and the associated Lyapunov equations are respectively described as $\boldsymbol{A}\boldsymbol{P} + \boldsymbol{P}\boldsymbol{A}^* = -\boldsymbol{b}\boldsymbol{b}^*$ and $\boldsymbol{A}^*\boldsymbol{Q} + \boldsymbol{Q}\boldsymbol{A} = -\boldsymbol{c}^*\boldsymbol{c}$.

From (25), it follows that the controllability and observability Gramians of this system are

$$\boldsymbol{P} = \begin{bmatrix} \frac{1}{\alpha_1 + \alpha_1^*} & \cdots & \frac{1}{\alpha_1 + \alpha_N^*} \\ \vdots & \ddots & \vdots \\ \frac{1}{\alpha_N + \alpha_1^*} & \cdots & \frac{1}{\alpha_N + \alpha_N^*} \end{bmatrix}, \quad \boldsymbol{Q} = \boldsymbol{\mathcal{B}}^*\boldsymbol{\mathcal{S}}\boldsymbol{\mathcal{B}} \qquad (26)$$

where

$$\boldsymbol{\mathcal{B}} = \mathrm{diag}(\beta_1, \cdots, \beta_N), \quad \boldsymbol{\mathcal{S}} = \begin{bmatrix} \frac{1}{\alpha_1^* + \alpha_1} & \cdots & \frac{1}{\alpha_1^* + \alpha_N} \\ \vdots & \ddots & \vdots \\ \frac{1}{\alpha_N^* + \alpha_1} & \cdots & \frac{1}{\alpha_N^* + \alpha_N} \end{bmatrix}. \qquad (27)$$

### B. Main Result

Here we discuss the relationship between the Hankel singular values and reflected zeros. We first consider the case where a reflected zero is real, and present the following proposition.

*Proposition 2:* Let $G(s)$ be a strictly proper transfer function with simple poles, and let one of its real zeros be denoted by $\zeta_r$. Also, consider the new transfer function $\widehat{G}(s)$ that is obtained from $G(s)$ by reflecting $\zeta_r$, i.e. by moving $\zeta_r$ to $-\zeta_r$. Now, let the Hankel singular values of $G(s)$ and $\widehat{G}(s)$ be respectively denoted by $\sigma_k$ and $\widehat{\sigma}_k$ for $1 \leq k \leq N$. Then, $\widehat{\sigma}_k \geq \sigma_k$ holds for $\zeta_r < 0$, and $\widehat{\sigma}_k \leq \sigma_k$ holds for $\zeta_r > 0$.

*Proof:* We first consider state-space formulation of $\widehat{G}(s)$. From (25), it is easy to see that $\widehat{G}(s)$ has a realization $(\widehat{\boldsymbol{A}}, \widehat{\boldsymbol{b}}, \widehat{\boldsymbol{c}})$ such that

$$\widehat{\boldsymbol{A}} = \boldsymbol{A}, \quad \widehat{\boldsymbol{b}} = \boldsymbol{b}, \quad \widehat{\boldsymbol{c}} = [\widehat{\beta}_1, \cdots, \widehat{\beta}_N] \qquad (28)$$

where

$$\widehat{\beta}_k = \underbrace{\frac{-\alpha_k + \zeta_r}{-\alpha_k - \zeta_r}}_{=: \gamma_k} \beta_k = \gamma_k \beta_k, \quad 1 \leq k \leq N. \qquad (29)$$

The controllability and observability Gramians $(\widehat{\boldsymbol{P}}, \widehat{\boldsymbol{Q}})$ of this system are related to $(\boldsymbol{P}, \boldsymbol{Q})$ as follows:

$$\widehat{\boldsymbol{P}} = \boldsymbol{P}, \quad \widehat{\boldsymbol{Q}} = \boldsymbol{\Gamma}^*\boldsymbol{Q}\boldsymbol{\Gamma} \qquad (30)$$

where we let $\boldsymbol{\Gamma} = \mathrm{diag}(\gamma_1, \cdots, \gamma_N)$.

Now, consider $\widehat{\boldsymbol{Q}} - \boldsymbol{Q}$. This matrix is simplified by the above relationships as follows:

$$\begin{aligned} \widehat{\boldsymbol{Q}} - \boldsymbol{Q} &= \boldsymbol{\Gamma}^*(\boldsymbol{\mathcal{B}}^*\boldsymbol{\mathcal{S}}\boldsymbol{\mathcal{B}})\boldsymbol{\Gamma} - \boldsymbol{\mathcal{B}}^*\boldsymbol{\mathcal{S}}\boldsymbol{\mathcal{B}} \\ &= \boldsymbol{\mathcal{B}}^*(\boldsymbol{\Gamma}^*\boldsymbol{\mathcal{S}}\boldsymbol{\Gamma} - \boldsymbol{\mathcal{S}})\boldsymbol{\mathcal{B}} \\ &= -2\zeta_r\boldsymbol{\mathcal{B}}^*\boldsymbol{v}^*\boldsymbol{v}\boldsymbol{\mathcal{B}} \end{aligned} \qquad (31)$$

where $\boldsymbol{v} = [1/(\alpha_1 + \zeta_r), \cdots, 1/(\alpha_N + \zeta_r)]$. Eq. (31) shows that $\widehat{\boldsymbol{Q}} - \boldsymbol{Q}$ is positive (negative) semidefinite (of rank one) depending on whether $\zeta_r < 0 (\zeta_r > 0)$. Hence $\lambda_k(\widehat{\boldsymbol{P}}\widehat{\boldsymbol{Q}}) \geq \lambda_k(\boldsymbol{P}\boldsymbol{Q})$ holds for $\zeta_r < 0$, and $\lambda_k(\widehat{\boldsymbol{P}}\widehat{\boldsymbol{Q}}) \leq \lambda_k(\boldsymbol{P}\boldsymbol{Q})$ holds for $\zeta_r > 0$. This result shows the desired inequality for the Hankel singular values. ∎

Next we turn our attention to the case of reflecting a pair of complex conjugate zeros. The result is given as the following proposition.

*Proposition 3:* Let $G(s)$ be a strictly proper transfer function with simple poles, and let a pair of its complex conjugate zeros be denoted by $\zeta_c$ and $\zeta_c^*$, respectively. Also, let $\widehat{G}(s)$ be the transfer fucntion that is obtained from $G(s)$ by reflecting these zeros, i.e. by moving $\zeta_c$ and $\zeta_c^*$ to $-\zeta_c^*$ and $-\zeta_c$, respectively. Now, let the Hankel singular values of $G(s)$ and $\widehat{G}(s)$ be respectively $\sigma_k$ and $\widehat{\sigma}_k$. Then, $\widehat{\sigma}_k \geq \sigma_k$ holds if $\mathrm{Re}(\zeta_c) < 0$, and $\widehat{\sigma}_k \leq \sigma_k$ holds if $\mathrm{Re}(\zeta_c) > 0$.

*Proof:* In this case, $\gamma_k$ in (29) becomes

$$\gamma_k = \frac{\alpha_k - \zeta_c}{\alpha_k + \zeta_c} \cdot \frac{\alpha_k - \zeta_c^*}{\alpha_k + \zeta_c^*}. \tag{32}$$

Consequently, the $(i,j)$-th entry of $\boldsymbol{\Gamma}^* \boldsymbol{S} \boldsymbol{\Gamma} - \boldsymbol{S}$ becomes

$$
\begin{aligned}
& [\boldsymbol{\Gamma}^* \boldsymbol{S} \boldsymbol{\Gamma} - \boldsymbol{S}]_{i,j} \\
&= \frac{\gamma_i^* \gamma_j - 1}{\alpha_i^* + \alpha_j} \\
&= -2\left(\zeta_c + \zeta_c^*\right) \\
&\quad \times \frac{\zeta_c^* \zeta_c + \alpha_i^* \alpha_j}{\left(\alpha_i^* + \zeta_c\right)\left(\alpha_i^* + \zeta_c^*\right)\left(\alpha_j + \zeta_c\right)\left(\alpha_j + \zeta_c^*\right)}
\end{aligned} \tag{33}
$$

which leads to

$$\widehat{\boldsymbol{Q}} - \boldsymbol{Q} = -2\left(\zeta_c + \zeta_c^*\right) \boldsymbol{\mathcal{B}}^* \boldsymbol{\Delta}^* [\boldsymbol{Z}^* \boldsymbol{Z} + \boldsymbol{\alpha}^* \boldsymbol{\alpha}] \boldsymbol{\Delta} \boldsymbol{\mathcal{B}} \tag{34}$$

where

$$
\begin{aligned}
\boldsymbol{\Delta} &= \operatorname{diag}\left[\frac{1}{(\alpha_1 + \zeta_c)(\alpha_1 + \zeta_c^*)}, \cdots, \frac{1}{(\alpha_N + \zeta_c)(\alpha_N + \zeta_c^*)}\right] \\
\boldsymbol{Z} &= [\zeta_c, \cdots, \zeta_c], \quad \boldsymbol{\alpha} = [\alpha_1, \cdots, \alpha_N].
\end{aligned} \tag{35}
$$

Thus, since $\boldsymbol{\Delta}^*[\boldsymbol{Z}^* \boldsymbol{Z} + \boldsymbol{\alpha}^* \boldsymbol{\alpha}] \boldsymbol{\Delta} \geq 0$, it follows that $\widehat{\boldsymbol{Q}} - \boldsymbol{Q}$ is positive (negative) semidefinite provided that $\operatorname{Re}(\zeta_c) < 0 (\operatorname{Re}(\zeta_c) > 0)$. This shows $\lambda_k(\widehat{\boldsymbol{P}}\widehat{\boldsymbol{Q}}) \geq \lambda_k(\boldsymbol{P}\boldsymbol{Q})$ for $\operatorname{Re}(\zeta_c) < 0$ and $\lambda_k(\widehat{\boldsymbol{P}}\widehat{\boldsymbol{Q}}) \leq \lambda_k(\boldsymbol{P}\boldsymbol{Q})$ for $\operatorname{Re}(\zeta_c) > 0$, which completes the proof. ∎

Propositions 2 and 3 show that if *one* (real or complex) zero is moved to a location which is symmertic with respect to the imaginary axis, the Hankel singular increase (decrease) pointwise. Clearly this will be even more the case if a greater number of zeros is moved (all in the same direction). We conclude that if a minimum phase system becomes maximum phase (i.e. all stable zeros are reflected with respect to the imaginary axis), the Hankel singular values will increase pointwise. Hence the pointwise inequality (23) is derived by this approach.

The significance of the above analysis is that description of the Gramians are provided in terms of reflected zeros of the transfer functions. Therefore, this description clearly tells us the influence of reflecting zeros upon the Hankel singular values. It appears that this analysis leads to characterization of the values of the shift of the Hankel singular values in terms of the values of reflected zeros, although this topic is currently an open question.

## VI. NUMERICAL EXAMPLE

This section gives a numerical example to demonstrate our main result. Consider the following family of 6th-order continuous-time transfer functions:

$$G_{\min p}(s) = \frac{(s-z_1)(s-z_2)(s-z_3)(s-z_4)(s-z_5)}{(s-p_1)^2(s-p_2)(s-p_3)(s-p_4)(s-p_5)} \tag{36}$$

$$G_{\mathrm{np}1}(s) = \frac{(s+z_1^*)(s-z_2)(s-z_3)(s-z_4)(s-z_5)}{(s-p_1)^2(s-p_2)(s-p_3)(s-p_4)(s-p_5)} \tag{37}$$

$$G_{\mathrm{np}2}(s) = \frac{(s+z_1^*)(s+z_2^*)(s+z_3^*)(s-z_4)(s-z_5)}{(s-p_1)^2(s-p_2)(s-p_3)(s-p_4)(s-p_5)} \tag{38}$$

$$G_{\max p}(s) = \frac{(s+z_1^*)(s+z_2^*)(s+z_3^*)(s+z_4^*)(s+z_5^*)}{(s-p_1)^2(s-p_2)(s-p_3)(s-p_4)(s-p_5)} \tag{39}$$

where $(z_1, z_2, z_3, z_4, z_5) = (-1, -0.5 + j0.5, -0.5 - j0.5, -0.8 + j1.2, -0.8 - j1.2)$ and $(p_1, p_2, p_3, p_4, p_5) = (-1.5, -0.2 + j0.6, -0.2 - j0.6, -2 + j, -2 - j)$. Note that $G_{\min p}(s)$ is the minimum phase system, and the other systems $G_{\mathrm{np}1}(s)$, $G_{\mathrm{np}2}(s)$ and $G_{\max p}(s)$ are obtained from $G_{\min p}(s)$ by reflecting $\{z_1\}$, $\{z_1, z_2, z_3\}$ and $\{z_1, z_2, z_3, z_4, z_5\}$, respectively. The Hankel singular

### TABLE I
RELATIONSHIP BETWEEN HANKEL SINGULAR VALUES AND REFLECTED ZEROS

| Transfer function | Hankel singular values |
|---|---|
| $G_{\min p}(s)$ | $(0.1662, 0.1041, 0.0700, 0.0191, 0.0026, 0.0001)$ |
| $G_{\mathrm{np}1}(s)$ | $(0.2874, 0.1305, 0.0789, 0.0369, 0.0030, 0.0023)$ |
| $G_{\mathrm{np}2}(s)$ | $(0.3325, 0.3189, 0.2059, 0.0927, 0.0220, 0.0062)$ |
| $G_{\max p}(s)$ | $(0.3540, 0.3289, 0.2467, 0.1979, 0.1729, 0.0811)$ |

values of these systems are given in Table I, which shows that these sets of Hankel singular values satisfy the pointwise inequality.

## VII. CONCLUSION

This note has derived a pointwise inequality that is concerned with the Hankel singular values and reflected zeros of the transfer function. It has been shown that the Hankel singular values increase (decrease) pointwise when a zero is reflected with respect to the imaginary axis or the unit circle. This leads to the fact that the minimum phase system has the smallest Hankel singular values while the maximum phase system has the largest ones, and that the Hankel singular values of the other non-minimum phase systems lie in between. Our numerical example has demonstrated this property.

Practical implications of our main result can be derived in the field of signal processing theory. As stated in Section I, the Hankel singular values characterize optimal values with respect to the dynamic range of analog filters and the quantization effects of digital filters. Therefore, from this fact and our main result, we immediately know that analog or digital filters of minimum phase attain higher performance with respect to the dynamic range and quantization effects than any other non-minimum phase filter. This property also holds in the field of balanced model reduction: it follows that the minimum phase transfer function has the smallest value of the upper bound of the approximation error.[3]

## REFERENCES

[1] B. C. Moore, "Principal component analysis in linear systems: Controllability, observability, and model reduction," *IEEE Trans. Automat. Control*, vol. AC-26, no. 1, pp. 17–32, Feb. 1981.

[2] L. Pernebo and L. M. Silverman, "Model reduction via balanced state space representations," *IEEE Trans. Automat. Control*, vol. AC-27, no. 2, pp. 382–387, Apr. 1982.

[3] K. Glover, "All optimal Hankel-norm approximations of linear multivariable systems and their $L_\infty$-error bounds," *Int. J. Control*, vol. 39, pp. 1115–1193, 1984.

[4] A. C. Antoulas, "Approximation of large-scale dynamical systems," in *Advances in Design and Control*. Philadelphia, PA: SIAM, 2005, vol. DC-06.

[5] G. Groenewold, "The design of high dynamic range continuous-time integratable bandpass filters," *IEEE Trans. Circuits Syst.*, vol. CAS-38, no. 8, pp. 838–852, Aug. 1991.

[6] W. M. Snelgrove and A. S. Sedra, "Synthesis and analysis of state-space active filters using intermediate transfer functions," *IEEE Trans. Circuits Syst.*, vol. CAS-33, no. 3, pp. 287–301, Mar. 1986.

[7] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-25, no. 4, pp. 273–281, Aug. 1977.

[8] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, no. 9, pp. 551–562, Sep. 1976.

[9] M. Kawamata and T. Higuchi, "A unified approach to the optimal synthesis of fixed-point state-space digital filters," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-33, no. 4, pp. 911–920, Aug. 1985.

[10] 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. CAS-37, no. 1, pp. 72–80, Jan. 1990.

[3]This fact is originally proved in [19] by describing the decay rates of the Hankel singular values.

[11] B. Hanzon, "The area enclosed by the (oriented) Nyquist diagram and the Hilbert–Schmidt–Hankel norm of a linear system," *IEEE Trans. Automat. Control*, vol. 37, no. 6, pp. 835–839, Jun. 1992.

[12] S. Gugercin and A. C. Antoulas, "A survey of model reduction by balanced truncation and some new results," *Int. J. Control*, vol. 77, no. 8, pp. 748–766, May 2004.

[13] P. C. Opdenacker and E. A. Jonckheere, "A contraction mapping preserving balanced reduction scheme and its infinity norm error bounds," *IEEE Trans. Circuits Syst.*, vol. CAS-35, no. 2, pp. 184–189, Feb. 1988.

[14] R. Ober, "Balanced parametrization of classes of linear systems," *SIAM J. Control Optim.*, vol. 29, no. 6, pp. 1251–1287, Nov. 1991.

[15] P. P. Vaidyanathan, "The discrete-time bounded-real lemma in digital filtering," *IEEE Trans. Circuits Syst.*, vol. CAS-32, no. 9, pp. 918–924, Sep. 1985.

[16] C. E. de Souza and L. Xie, "On the discrete-time bounded real lemma with application in the characterization of static state feedback $H_\infty$ controllers," *Syst. Control Lett.*, vol. 18, no. 1, pp. 61–71, Jan. 1992.

[17] C. T. Mullis and R. A. Roberts, "Roundoff noise in digital filters: Frequency transformations and invariants," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-24, no. 6, pp. 538–550, Dec. 1976.

[18] J. R. Partington, *An Introduction to Hankel Operators*.   London, U.K.: Cambridge University Press, 1988.

[19] A. C. Antoulas, D. C. Sorensen, and Y. Zhou, "On the decay rate of Hankel singular values and related issues," *Syst. Control Lett.*, vol. 46, no. 5, pp. 323–342, Aug. 2002.

# Likelihood Gradient Evaluation Using Square-Root Covariance Filters

M. V. Kulikova

*Abstract*—Using the array form of numerically stable square-root implementation methods for Kalman filtering formulas, we construct a new square-root algorithm for the log-likelihood gradient (score) evaluation. This avoids the use of the conventional Kalman filter with its inherent numerical instabilities and improves the robustness of computations against roundoff errors. The new algorithm is developed in terms of covariance quantities and based on the "condensed form" of the array square-root filter.

*Index Terms*—Gradient methods, identification, Kalman filtering, maximum likelihood estimation, numerical stability.

## I. INTRODUCTION

Consider the discrete-time linear stochastic system

$$x_k = F_k x_{k-1} + G_k w_k \qquad (1)$$
$$z_k = H_k x_k + v_k, \quad k = 1, \ldots, N \qquad (2)$$

where $x_k \in \mathbb{R}^n$ and $z_k \in \mathbb{R}^m$ are, respectively, the state and the measurement vectors; $k$ is a discrete time, i.e. $x_k$ means $x(t_k)$. The noises $w_k \in \mathbb{R}^q$, $v_k \in \mathbb{R}^m$ and the initial state $x_0 \sim \mathcal{N}(\bar{x}_0, \Pi_0)$ are taken from mutually independent Gaussian distributions with zero mean and covariance matrices $Q_k$ and $R_k$, respectively, i.e. $w_k \sim \mathcal{N}(0, Q_k)$, $v_k \sim \mathcal{N}(0, R_k)$. Additionally, system (1), (2) is parameterized by a vector of unknown system parameters $\theta \in \mathbb{R}^p$, which needs to be estimated. This means that the entries of the matrices $F_k$, $G_k$, $H_k$, $Q_k$, $R_k$

and $\Pi_0$ are functions of $\theta \in \mathbb{R}^p$. However, for the sake of simplicity we will suppress the corresponding notations below, i.e instead of $F_k(\theta)$, $G_k(\theta)$, $H_k(\theta)$, $Q_k(\theta)$, $R_k(\theta)$ and $\Pi_0(\theta)$ we will write $F_k$, $G_k$, $H_k$, $Q_k$, $R_k$ and $\Pi_0$.

Solving the parameter estimation problem by the method of maximum likelihood requires the maximization of the likelihood function (LF) with respect to unknown system parameters. It is often done by using a gradient approach where the computation of the likelihood gradient (LG) is necessary. For the state-space system (1), (2) the negative Log LF is given as [1]

$$L_\theta\left(Z_1^N\right) = \frac{1}{2}\sum_{k=1}^{N}\left\{\frac{m}{2}\ln(2\pi) + \ln(\det R_{e,k}) + e_k^T R_{e,k}^{-1} e_k\right\}$$

where $Z_1^N = [z_1, \ldots, z_N]$ is $N$-step measurement history and $e_k$ are the innovations, generated by the discrete-time Kalman filter (KF), with zero mean and covariance matrix $R_{e,k}$. They are $e_k = z_k - H_k \hat{x}_{k|k-1}$ and $R_{e,k} = H_k P_{k|k-1} H_k^T + R_k$, respectively. The KF defines the one-step ahead predicted state estimate $\hat{x}_{k|k-1}$ and the one-step predicted error covariance matrix $P_{k|k-1}$.

Straight forward differentiation of the KF equations is a direct approach to the Log LG evaluation, known as a "score". This leads to a set of $p$ vector equations, known as the *filter sensitivity equations*, for computing $\partial \hat{x}_{k|k-1}/\partial\theta$, and a set of $p$ matrix equations, known as the *Riccati-type sensitivity equations*, for computing $\partial P_{k|k-1}/\partial\theta$.

Consequently, the main disadvantage of the standard approach is the problem of numerical instability of the conventional KF, i.e divergence due to the lack of reliability of the numerical algorithm. Solution of the matrix Riccati equation is a major cause of numerical difficulties in the conventional KF implementation, from the standpoint of computational load as well as from the standpoint of computational errors [2].

The alternative approach can be found in, so-called, square-root filtering algorithms. It is well known that numerical solution of the Riccati equation tends to be more robust against roundoff errors if Cholesky factors or modified Cholesky factors (such as the $U^T D U$-algorithms [3]) of the covariance matrix are used as the dependent variables. The resulting KF implementation methods are called square-root filters (SRF). They are now generally preferred for practical use [2], [4], [5]. For more insights about numerical properties of different KF implementation methods we refer to the celebrated paper of Verhaegen and Van Dooren [6].

Increasingly, the preferred form for algorithms in many fields is now the array form [7]. Several useful SRF algorithms for KF formulas formulated in the array form have been recently proposed by Park and Kailath [8]. For this implementations the reliability of the filter estimates is expected to be better because of the use of numerically stable orthogonal transformations for each recursion step. Apart from numerical advantages, array SRF algorithms appear to be better suited to parallel and to very large scale integration (VLSI) implementations [8], [9].

The development of numerically stable implementation methods for KF formulas has led to the hope that the Log LG (with respect to unknown system parameters) might be computed more accurately. For this problem, a number of questions arise:

- Is it possible to extend reliable array SRF algorithms to the case of the Log LG evaluation?
- If such methods exist, will they inherit the advantages from the source filtering implementations? In particular, will they improve the robustness of the computations against roundoff errors compared to the conventional KF technique? The question about suitability for parallel implementation is beyond the scope of this technical note.