

World Scientific Series in Applicable Analysis
Volume 5

Editor

R. P. Agarwal

*Department of Mathematics
National University of Singapore*

RECENT TRENDS IN
OPTIMIZATION THEORY
— AND —
APPLICATIONS

 **World Scientific**
Singapore • New Jersey • London • Hong Kong

RICCATI DIFFERENCE EQUATIONS FOR DISCRETE TIME SPECTRAL FACTORIZATION WITH UNIT CIRCLE ZEROS.¹

JEREMY D. MATSON

BRIAN D.O. ANDERSON

Department of Systems Engineering,
Research School of Information
Sciences and Engineering,
Australian National University,
Canberra ACT 0200, AUSTRALIA

ALAN J. LAUB

Department of Electrical
and Computer Engineering,
University of California,
Santa Barbara, CA 93106-9560
U.S.A.

DAVID J. CLEMENTS

School of Electrical Engineering,
University of New South Wales,
P.O. Box 1,
Kensington 2033, AUSTRALIA

ABSTRACT

Spectral matrices that have unit circle transmission zeros arise in the consideration of H_∞ control, the bounded-real lemma and discrete spectral factorization problems. Spectral matrices with unit circle invariant (but not transmission) zeros arise when considering Kalman filtering for systems with unit circle modes which are not corrupted by process noise. It is well known that if a spectral matrix is generically nonsingular, minimum phase spectral factors can be constructed from a strong solution of an algebraic Riccati equation associated with a state-space realization of the spectral matrix. Closely associated with the algebraic equation is a Riccati difference equation whose iterates are shown to converge to the strong solution under fairly mild conditions on the realization of the spectral matrix. The key observations made in this paper concern the fine structure of the Riccati difference equation iterates from which a convergence rate of $\mathcal{O}(\frac{1}{2})$ is deduced.

1 Introduction.

1.1 Spectral Factorization Review.

In discrete time, a spectral matrix $\Psi(z)$ is a square real rational matrix-valued function of a complex variable z with the properties that $\Psi^T(z^{-1}) = \Psi(z)$ and $\Psi(e^{j\theta}) \geq 0$, for all $\theta \in [0, 2\pi)$. We consider only spectral matrices which are generically nonsingular, in the sense that $\det(\Psi(z)) \neq 0$. Spectral matrices arise naturally in the description of stochastic processes, in the formulation of linear control and filtering problems and in the discrete bounded-real lemma.

¹The authors wish to acknowledge the funding of the activities of the Cooperative Research Centre for Robust and Adaptive Systems by the Australian Commonwealth Government under the Cooperative Research Centres Program.

It is well known that the construction of spectral factors is central to the solution of the abovementioned problems. A *spectral factor* $\Omega(z)$ of $\Psi(z)$ is a real rational matrix-valued function of the complex variable z which satisfies $\Omega^T(z^{-1})\Omega(z) = \Psi(z)$. If, in addition, $\Omega^{-1}(z)$ exists and is analytic when $|z| > 1$, $\Omega(z)$ is called a *minimum phase spectral factor*.

It is also well known^{3,11} that if $\Psi(z)$ is nonsingular, there exists a spectral decomposition of the form

$$\Psi(z) = \Theta^T(z^{-1})N\Theta(z) \quad (1)$$

where N is a positive definite symmetric matrix and $\Theta(z)$ is a square real rational transfer function matrix which is invertible, satisfies $\lim_{z \rightarrow \infty} \Theta(z) = I$ and, along with its inverse, is analytic when $|z| > 1$. Hence a minimum-phase spectral factor of $\Psi(z)$ can be constructed as $\Omega(z) = N^{\frac{1}{2}}\Theta(z)$.

In this paper, we consider a class of nonsingular spectral matrices, having a state-space realization of the form

$$\Psi(z) = U + G^T(z^{-1}I - F^T)^{-1}V(zI - F)^{-1}G \quad (2)$$

where each constant matrix is real and where the following assumptions hold:

A.1 (F, G) is stabilizable.

A.2 $V = V^T$, $U = U^T$ and U is nonsingular.

In fact, by applying appropriate preliminary transformations, most nonsingular spectral factorization problems can be treated via a spectral matrix of the above form.

It should be emphasized that no assumptions regarding the sign-definiteness of either V or U have been made. In H_2 linear-quadratic optimal control and Kalman filtering, spectral matrices arise which are special cases of the above class in which generally $U > 0$ and $V \geq 0$. Clearly these conditions preclude the possibility that the spectral matrix realization has unit circle transmission zeros. Note, however, that unit circle *invariant* zeros can appear if the realization of the spectral matrix is non-minimal: for example if $(F, V^{\frac{1}{2}})$ has unobservable unit-circle modes,^{4,7,8} these become a subset of the invariant zeros of the spectral matrix realization.

Spectral matrices for which unit circle *transmission* zeros can occur arise in discrete time spectral factorization¹ and in the discrete time version of the bounded-real lemma,¹⁰ which is relevant in the H_∞ control problem. Recall that a discrete-time transfer function matrix $L(z)$ is called *bounded real* if all poles of $L(z)$ are inside the unit circle and $\|L\|_\infty \leq 1$. Consider the spectral matrix $\Psi(z)$ defined by $\Psi(z) = I - L^T(z^{-1})L(z)$. Observe that with the state-space realization $L(z) = H_L(zI - F_L)^{-1}G_L$, $\Psi(z)$ is of the standard form given in Eq. 2 with $U = I$, $F = F_L$, $G = G_L$ and $V = -H_L^T H_L$. Should $\sigma_{\max}(L(e^{j\theta^*})) = 1$ for some θ^* , then $\Psi(e^{j\theta^*})$ loses rank at that point, corresponding to a transmission zero.

Real symmetric solutions of the discrete time algebraic Riccati equation (ARE)

$$\Phi = F^T (\Phi - \Phi G (U + G^T \Phi G)^{-1} G^T \Phi) F + V \quad (3)$$

enable the state-space construction of spectral factors of $\Psi(z)$. Such equations have been studied in many contexts including spectral factorization,¹ and infinite-horizon control and filtering problems.^{2,3} It can be demonstrated using Eq. 3 that Eq. 1 is satisfied with the definitions $N = U + G^T \Phi G$ and $\Theta(z) = I + N^{-1} G^T \Phi F (zI - F)^{-1} G$. The resulting spectral factor $\Omega(z) = N^{1/2} \Theta(z)$ has an inverse

$$\Omega^{-1}(z) = (I - N^{-1} G^T \Phi F (zI - \hat{F})^{-1} G) N^{-1/2} \quad (4)$$

where \hat{F} is the closed-loop matrix given by $\hat{F} = (I - GN^{-1}G^T\Phi)F$. A solution Φ of Eq. 3 is said to be *strong* if \hat{F} has all eigenvalues either inside or on the unit circle. Note that the eigenvalues of \hat{F} are also the invariant zeros of $\Omega(z)$ and thus spectral factors constructed from strong solutions of Eq. 3 have the minimum phase property.

Remark: It is not the purpose of the present paper to address the question of when a unique strong solution Φ of Eq. 3 exists. Henceforth we assume that such a solution exists for the realization of the spectral matrix at hand. \square

We now consider the Riccati difference equation (RDE) associated with Eq. 3

$$\Phi_{k+1} = F^T (\Phi_k - \Phi_k G (U + G^T \Phi_k G)^{-1} G^T \Phi_k) F + V \quad (5)$$

where this equation has some real symmetric initial condition Φ_0 . The main result of this paper (which follows immediately) presents conditions under which iterates of the RDE Φ_k ($k \in \{0, 1, 2, \dots\}$) converge to the strong solution of Eq. 3 and describes the convergence rate when the associated spectral matrix has unit circle invariant zeros.

1.2 Main Result.

Firstly, it is demonstrated that RDE convergence results previously established for linear-quadratic control and Kalman filtering problems⁸ and for spectral factorization¹ hold for any spectral matrix of the form in Eq. 2 under assumptions A.1 and A.2. The spectral matrix may have unit circle invariant zeros which arise due to non-minimal modes in its realization, transmission zeros, or any combination of these. Secondly, and most importantly, the fine structure of the Riccati difference equation iterates is investigated. This leads to new results concerning the rates at which the iterates of Eq. 5 converge to the strong solution of Eq. 3.

Theorem 1.1 *Given a realization as in Eq. 2 of a discrete time spectral matrix $\Psi(z)$ satisfying assumptions A.1 and A.2, along with the strong solution Φ of the associated algebraic Riccati equation Eq. 3, then provided $\Phi_0 \geq \Phi$, iterates of Eq. 5 have the following properties:*

1. $\Phi_k \geq \Phi$.
2. $\lim_{k \rightarrow \infty} \Phi_k = \Phi$.

3. If $\Psi(z)$ has an invariant zero on the unit circle, then there exist constants κ_1, κ_2 (depending on the realization of $\Psi(\cdot)$ and on Φ_0) with $\kappa_1 \geq \kappa_2 > 0$ such that:

a) For all $\epsilon > 0$, there exists a k_ϵ such that when $k \geq k_\epsilon$

$$\lambda_{\max}(\Phi_k - \Phi) \leq \frac{\kappa_1 + \epsilon}{k}. \quad (6)$$

b) When $\Phi_0 > \Phi$, there exists a k_0 such that when $k \geq k_0$

$$\lambda_{\max}(\Phi_k - \Phi) \geq \frac{\kappa_2}{k}. \quad (7)$$

A proof of this result is delayed until the final section of the paper. Item 3 a) of this theorem reports a worst-case $\frac{1}{k}$ convergence rate in the case of unit-circle invariant zeros. Item 3 b) says that, with the exclusion of (non-generic) cases where $\Phi_0 - \Phi$ is singular, the convergence rate can be no better than $\frac{1}{k}$.

Remark: In cases where the spectral matrix has no invariant zeros on the unit circle, an exponential convergence rate has been reported:^{1,6} there exist constants A, κ_3 such that $1 > \kappa_3 > 0$ and $\lambda_{\max}(\Phi_k - \Phi) \leq A\kappa_3^k$. \square

2 Preliminary Results.

2.1 Notation.

Let O_m denote the $m \times m$ zero matrix and I_m the $m \times m$ identity matrix. Given a matrix M , let $\{\sigma_i(M)\}$ denote the singular values of M and σ_{\max} be the largest of these; if M is square, denote its eigenvalues as $\{\lambda_i(M)\}$. Suppose M has an even number of rows and columns, consisting of a matrix of 2×2 matrix sub-blocks; for convenience we let $[M]_{ij}$ denote the $(i, j)^{\text{th}}$ 2×2 sub-block of M .

Given $f(l)$ and $g(l)$, both scalar functions of an integer variable l , we say $g(l) = \mathcal{O}(f(l))$ if there exists a constant $\kappa < \infty$ such that $\lim_{l \rightarrow \infty} \frac{|g(l)|}{|f(l)|} = \kappa$. Given $U(l)$, a square matrix-valued function of l , we say that $U(l) = \mathcal{O}(f(l))$ if $\sigma_{\max}(U(l)) = \mathcal{O}(f(l))$. Note that this definition has the following property: If $U(l)$ is such that each $(U(l))_{mn} = \mathcal{O}(f(l))$ or each $[U(l)]_{ij} = \mathcal{O}(f(l))$, then $U(l) = \mathcal{O}(f(l))$.

Real Jordan Form.

The following summarizes standard results concerning the real Jordan decomposition.¹² Any real square matrix B can be expressed as $B = TAT^{-1}$ where

$$A = \text{diag}\{A_1, \dots, A_p\} \quad (8)$$

and p is the number of real Jordan blocks. For each $q \in \{1, \dots, p\}$, A_q has one of the two forms described below.

In the first form, $A_q \in \mathbb{R}^{2n_q \times 2n_q}$, where $n_q \geq 1$ and

$$A_q = \begin{pmatrix} \Lambda_q & & & \\ I_2 & \Lambda_q & & \\ & \ddots & \ddots & \\ & & & \Lambda_q \\ & & & I_2 & \Lambda_q \end{pmatrix} \quad (9)$$

where

$$\Lambda_q = \begin{pmatrix} \sigma_q & \omega_q \\ -\omega_q & \sigma_q \end{pmatrix}. \quad (10)$$

In this case, $\sigma_q \pm j\omega_q$ is a pair of complex conjugate eigenvalues of B . If λ_q is a real eigenvalue of B , then in the second form, $A_q \in \mathbb{R}^{n_q \times n_q}$, where $n_q \geq 1$ and

$$A_q = \begin{pmatrix} \lambda_q & & & \\ 1 & \lambda_q & & \\ & \ddots & \ddots & \\ & & & \lambda_q \\ & & & 1 & \lambda_q \end{pmatrix}. \quad (11)$$

2.2 Asymptotic Behaviour of a Linear Matrix Difference Equation with Jordan Structure.

Due to its significance in describing the convergence behaviour of the Riccati difference equation to its strong solution, we examine the behaviour of the following linear matrix difference equation with initial condition $X_q(0) = 0$:

$$X_q(k+1) = A_q X_q(k) A_q^T + I. \quad (12)$$

In the first instance, we assume that A_q has the first real Jordan form as described in Eq. 9 which corresponds to a complex conjugate pair of eigenvalues. Similar arguments to those which follow for this case can be used in the second case and are thus not treated here. It follows by direct iteration of Eq. 12 that for $k \geq 1$, $X_q(k) = S_q(k)$ where

$$S_q(k) = \sum_{l=0}^{k-1} A_q^l (A_q^T)^l. \quad (13)$$

We focus here on the case where the Jordan block A_q corresponds to a complex conjugate pair of *unit circle* eigenvalues. It will be shown in the next section that it is the behaviour of iterates of this type that are the most important in establishing the convergence rate of RDEs associated with spectral matrices which have unit circle invariant zeros.

Lemma 2.1 Let A_q be a Jordan block of size $2n_q \times 2n_q$ which has the form given in Eq. 9, corresponding to a complex conjugate pair of unit circle eigenvalues.

For sufficiently large l , one has the following identity:

$$A_q^l (A_q^T)^l = C(l) + P(l) = C(l)(I + \mathcal{O}(l^{-1})) \quad (14)$$

where

$$C(l) = \begin{pmatrix} I_2 & l\Lambda_q & \frac{l^2}{2}\Lambda_q^2 & \cdots & \frac{l^{n_q-1}}{(n_q-1)!}\Lambda_q^{n_q-1} \\ l(\Lambda_q^T) & l^2 I_2 & \frac{l^3}{2}\Lambda_q & \cdots & \vdots \\ \frac{l^2}{2}(\Lambda_q^T)^2 & \frac{l^3}{2}(\Lambda_q^T) & \frac{l^4}{4}I_2 & \cdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{l^{n_q-1}}{(n_q-1)!}(\Lambda_q^T)^{n_q-1} & \cdots & \cdots & \cdots & \frac{l^{2n_q-2}}{(n_q-1)!(n_q-1)!}I_2 \end{pmatrix} \quad (15)$$

or equivalently

$$[C(l)]_{ij} = \begin{cases} \frac{l^{i+j-2}}{(i-1)!(j-1)!}\Lambda_q^{j-i} & \text{if } i \leq j \\ [C(l)]_{ji}^T & \text{if } i > j \end{cases}, \quad (16)$$

and

$$[P(l)]_{ij} = \begin{cases} \mathcal{O}(l^{i+j-3}) & \text{if } i+j \geq 3 \\ O_2 & \text{if } i=j=1 \end{cases}. \quad (17)$$

Proof: A straightforward calculation based on approximation of each 2×2 sub-block of $A_q^l (A_q^T)^l$ leads to Eq. 14. \square

Lemma 2.2 Let A_q be a Jordan block of size $2n_q \times 2n_q$ which has the form given in Eq. 9, corresponding to a complex conjugate pair of unit circle eigenvalues.

With $S_q(k)$ defined in Eq. 13, the following identity holds:

$$S_q(k) = D(k) + G(k) = D(k)(I + \mathcal{O}(k^{-1})) \quad (18)$$

where

$$[D(k)]_{ij} = \begin{cases} \frac{k^{i+j-1}}{(i-1)!(j-1)!(i+j-1)!}(\Lambda_q^T)^{i-j} & \text{if } i \geq j \\ [D(k)]_{ji}^T & \text{if } i < j \end{cases} \quad (19)$$

and

$$[G(k)]_{ij} = \begin{cases} \mathcal{O}(k^{i+j-2}) & \text{if } i+j \geq 3 \\ O_2 & \text{if } i=j=1 \end{cases}. \quad (20)$$

Moreover, there exists a symmetric matrix $\Theta \in \mathbb{R}^{2n_q \times 2n_q}$ such that $\Theta > 0$ and

$$D(k) = kH^T(k)\Theta H(k) \quad (21)$$

where

$$H(k) = \text{diag}\{J_2, k\Lambda_q, \dots, k^{(n_q-1)}\Lambda_q^{(n_q-1)}\}. \quad (22)$$

Proof: Follows by application of the results in Lemma 2.1 plus further 2×2 block approximations. \square

Lemma 2.3 Let A_q be a Jordan block of the form given in Eq. 9, having size $2n_q \times 2n_q$, which corresponds to a complex conjugate pair of eigenvalues on the unit circle. Let $S_q(k)$ be defined as in Eq. 13. Then there exists a constant $\zeta > 0$ such that

$$1. \quad \lambda_{\min}(S_q(k)) \leq k \frac{2n_q}{\zeta}, \quad (23)$$

2. for all $\epsilon > 0$, there exists a constant k_ϵ such that $k > k_\epsilon$ implies that

$$\lambda_{\min}(S_q(k)) \geq \frac{k}{\zeta + \epsilon}. \quad (24)$$

Proof: Recall from Lemma 2.2 that $S_q(k) = D(k) + G(k)$ with $D(k)$ given in (21). Since $H(k)$ is invertible, it follows that

$$S_q(k) = kH^T(k) \{ \Theta + W(k) \} H(k) \quad (25)$$

where $W(k) = H^{-T}(k) \frac{G(k)}{k} H^{-1}(k)$ from which it can be verified fairly simply that $W(k) = \mathcal{O}(k^{-1})$.

That $S_q(k)$ is always a positive definite matrix can be seen from its definition in Eq. 13. In order to describe its eigenvalues as a function of k , we investigate those of $S_q^{-1}(k)$. Recall that for any positive definite matrix M , if $\lambda_{\max}(M)$ is the maximum eigenvalue of M , then the minimum eigenvalue of M^{-1} is $\lambda_{\min}(M^{-1}) = \lambda_{\max}^{-1}(M)$.

It follows from Eq. 25 that $\Theta + W(k)$ is positive definite. Observe that $(\Theta + W(k))^{-1} = \Theta^{-\frac{1}{2}}(I + \mathcal{O}(k^{-1}))^{-1}\Theta^{-\frac{1}{2}}$. Now $(I + \mathcal{O}(k^{-1}))^{-1} = I + \mathcal{O}(k^{-1})$ and hence

$$(\Theta + W(k))^{-1} = \Theta^{-\frac{1}{2}}(I + \mathcal{O}(k^{-1}))\Theta^{-\frac{1}{2}} = \Theta^{-1} + \mathcal{O}(k^{-1}). \quad (26)$$

Inverting Eq. 25 and employing Eq. 26 reveals that

$$kS_q^{-1}(k) = H^{-1}(k) \{ \Theta^{-1} + \mathcal{O}(k^{-1}) \} H^{-T}(k) \quad (27)$$

and since $H^{-1}(k) = \mathcal{O}(1)$,

$$kS_q^{-1}(k) = H^{-1}(k)\Theta^{-1}H^{-T}(k) + \mathcal{O}(k^{-1}). \quad (28)$$

With M a nonnegative definite matrix of dimension n_M , recall the standard identity $\lambda_{\max}(M) \leq \text{trace}\{M\} \leq n_M \lambda_{\max}(M)$. Applying this result to Eq. 28 reveals that

$$\lambda_{\max}(kS_q^{-1}(k)) \leq \text{trace}\{H^{-1}(k)\Theta^{-1}H^{-T}(k)\} + \mathcal{O}(k^{-1}) \leq 2n_q \lambda_{\max}(kS_q^{-1}(k)). \quad (29)$$

Observe that

$$[H^{-1}(k)\Theta^{-1}H^{-T}(k)]_{ij} = \frac{1}{k^{2-i-j}}\Lambda_q^{1-i}[\Theta^{-1}]_{ij}(\Lambda_q^T)^{1-j} \quad (30)$$

and as a result that

$$\text{trace}\{H^{-1}(k)\Theta^{-1}H^{-T}(k)\} = \text{trace}\{[\Theta^{-1}]_{11}\} + \mathcal{O}(k^{-2}). \quad (31)$$

With the definition $\zeta = \text{trace}\{[\Theta^{-1}]_{11}\}$, note firstly from Eq. 29 and Eq. 31 that

$$\lambda_{\max}(kS_q^{-1}(k)) \leq \zeta + \mathcal{O}(k^{-1}) \quad (32)$$

and secondly that

$$\lambda_{\max}(kS_q^{-1}(k)) \geq \frac{\zeta}{2n_q}. \quad (33)$$

The stated results follow immediately from Eq. 32 and Eq. 33. \square

3 Riccati Difference Equation Convergence.

3.1 Comparison Theorem for RDE Iterates.

We now state a minor extension of a well-known result which describes the way in which Riccati difference equation iterates behave under perturbations to the initial condition Φ_0 . Having established this result, we will find it has several applications in the proof of RDE convergence.

Lemma 3.1 *Let the sequences $\{\Phi_k^1\}$ and $\{\Phi_k^2\}$ be defined by application of the Riccati Difference Equation, Eq. 5, with initial conditions Φ_0^1 and Φ_0^2 respectively. Define $\bar{\Phi}_k = \Phi_k^2 - \Phi_k^1$. Then*

1. *The following recursions hold for all $k \geq 0$:*

$$\bar{\Phi}_{k+1} = (\hat{F}_k^1)^T \bar{\Phi}_k \hat{F}_k^1 - (\hat{F}_k^1)^T \bar{\Phi}_k G(G^T \bar{\Phi}_k G + G^T \Phi_k^1 G + U)^{-1} G^T \bar{\Phi}_k \hat{F}_k^1 \quad (34)$$

$$\bar{\Phi}_{k+1} = (\hat{F}_k^2)^T \bar{\Phi}_k \hat{F}_k^2 + (\hat{F}_k^2)^T \bar{\Phi}_k G(G^T \Phi_k^1 G + U)^{-1} G^T \bar{\Phi}_k \hat{F}_k^2 \quad (35)$$

$$\text{where} \quad \hat{F}_k^1 = (I - G(G^T \Phi_k^1 G + U)^{-1} G^T \Phi_k^1) F$$

$$\hat{F}_k^2 = (I - G(G^T \Phi_k^2 G + U)^{-1} G^T \Phi_k^2) F.$$

2. *Suppose the RDE, Eq. 5, is associated with a nonsingular spectral matrix $\Psi(z)$ and the ARE, Eq. 3, has a strong solution $\bar{\Phi}$. Suppose also that both $\Phi_0^1 \geq \bar{\Phi}$ and $\Phi_0^2 \geq \bar{\Phi}$. Then if $\bar{\Phi}_0 \geq 0$ it follows that $\bar{\Phi}_k \geq 0$ for all $k \geq 0$.*

Proof: A more general version of the first difference equation in item 1 which also accounts for perturbations in V is well known.⁹ The second difference equation can be obtained from the first simply by first reversing the superscripts and then multiplying the equation by -1 .⁵

Item 2 has been established in the nonnegative definite cost case for LQ control and Kalman filtering.⁵ A generalization of this result to the broader class of spectral matrices we consider here follows from the discussion below.

Observe first that with $\Phi_0^1 = \Phi$, then $\Phi_k^1 = \Phi$ for all subsequent k . By hypothesis, $\Phi_0^2 \geq \Phi$ and hence $\bar{\Phi}_0 \geq 0$. Next observe from Eq. 35 that since $G^T \Phi G + U > 0$ (which follows from the assumed spectral property), it follows that $\Phi_k^2 \geq \Phi$ for all subsequent k .

Suppose now that one is given any $\Phi_0^1 \geq \Phi$. It follows from reversing subscripts in the argument immediately above that $\Phi_k^1 \geq \Phi$ for all subsequent k . Since $G^T \Phi G + U > 0$, it follows that $G^T \Phi_k^1 G + U > 0$ which together with Eq. 35 implies that $\bar{\Phi}_k \geq 0$ for all subsequent k . \square

3.2 A Preliminary Convergence Result.

In this subsection, a weakened version of the main theorem is proven in Lemma 3.2. In the following subsection, we show how the additional assumptions introduced in Lemma 3.2 may be relaxed.

Properties 1 and 2 in the following lemma have been stated in the literature.^{1,7} One of the first observations of the convergence rate stated in item 3 a) was in the context of a Kalman filtering example⁴ in which the plant model has an identity state mapping, with no process noise and observations corrupted by Gaussian white noise. The worst-case convergence rate given in item 3 a) has been stated¹ for a spectral factorization problem. Full proofs which spell out the mechanism and rate of convergence do not seem to be available in the literature, however. We now review the first steps towards a proof of the convergence result^{1,7} and then present a novel and nontrivial completion of the proof which addresses the question of convergence rate.

Lemma 3.2 Consider a realization Eq. 2 of a nonsingular discrete time spectral matrix $\Psi(z)$ which, in addition to assumptions A.1 and A.2, satisfies the following two assumptions:

A.3 (F, G) is controllable.

A.4 F is nonsingular.

Let Φ be the strong solution of the associated algebraic Riccati equation Eq. 3. Then provided $\Phi_0 > \Phi$, iterates of Eq. 5 have the following properties:

1. $\Phi_k > \Phi$.
2. $\lim_{k \rightarrow \infty} \Phi_k = \Phi$.

3. If $\Psi(z)$ has an invariant zero on the unit circle, then there exist constants δ_1, δ_2 (depending on the realization of $\Psi(\cdot)$ and on Φ_0) with $\delta_1 \geq \delta_2 > 0$ such that:

a) For all $\epsilon > 0$, there exists a k_ϵ such that when $k \geq k_\epsilon$

$$\lambda_{\max}(\Phi_k - \Phi) \leq \frac{\delta_1 + \epsilon}{k}. \quad (36)$$

b) There exists a k_0 such that when $k \geq k_0$

$$\lambda_{\max}(\Phi_k - \Phi) \geq \frac{\delta_2}{k}. \quad (37)$$

Proof: With the definitions $\Delta_k = \Phi_k - \Phi$ and $\hat{F} = (I - G(G^T \Phi G + U)^{-1} G^T \Phi) F$, one can apply Lemma 3.1 to obtain

$$\Delta_{k+1} = \hat{F}^T \Delta_k \hat{F} - \hat{F}^T \Delta_k G (G^T \Delta_k G + G^T \Phi G + U)^{-1} G^T \Delta_k \hat{F}. \quad (38)$$

Suppose $\Delta_k > 0$. Then since $N = U + G^T \Phi G > 0$, the so-called matrix inversion lemma may be applied to Eq. 38, revealing that

$$\Delta_{k+1}^{-1} = \hat{F}^{-1} \Delta_k^{-1} \hat{F}^{-T} + \hat{F}^{-1} G N^{-1} G^T \hat{F}^{-T}. \quad (39)$$

Invertibility of \hat{F} is a consequence of the invertibility of F and N ; application of the matrix inversion lemma yields $\hat{F}^{-1} = F^{-1}(I + G U^{-1} G^T \Phi)$. Since $N^{-1} > 0$, Eq. 39 implies that $\Delta_{k+1} > 0$ whenever $\Delta_k > 0$. Thus our assumption that $\Delta_0 > 0$ ensures $\Delta_k > 0$ for all $k \geq 0$. This establishes item 1 in the Lemma statement.

The proof of the convergence of iterates of Eq. 38 to zero is based on the following observation:^{1,7} $\lambda_{\min}(\Delta_k^{-1}) \rightarrow \infty$ implies $\lambda_{\max}(\Delta_k) \rightarrow 0$. An explicit account of the divergent behaviour of $\lambda_{\min}(\Delta_k^{-1}) \rightarrow \infty$ is given which draws upon the preliminary results obtained in Section 2.

Observe that \hat{F} is of the form $\hat{F} = F - GL$ (where $L = (G^T \Phi G + U)^{-1} G^T \Phi F$). It is a well known result that controllability of the pair (F, G) guarantees controllability of (\hat{F}, G) , which in turn implies the controllability of (\hat{F}^{-1}, G) . Since (\hat{F}^{-1}, G) is a controllable pair, so is the pair (\hat{F}^{-1}, GN^{-1}) . Hence the controllability Gramian for the latter pair satisfies

$$\hat{W} = \sum_{j=0}^{n-1} \hat{F}^{-j} G N^{-1} G^T (\hat{F}^T)^{-j} > 0,$$

with n the dimension of the state space.

Iteration of the identity Eq. 39 reveals that

$$\Delta_{k+n}^{-1} = \hat{F}^{-n} \Delta_k^{-1} (\hat{F}^T)^{-n} + \hat{W}. \quad (40)$$

With $\hat{A} = \hat{F}^{-n}$ and $\hat{X}_j = \Delta_{j+n}^{-1}$, Eq. 40 reads

$$\hat{X}_{j+1} = \hat{A} \hat{X}_j \hat{A}^T + \hat{W}. \quad (41)$$

Let \hat{A} have the real Jordan decomposition

$$\hat{A} = TAT^{-1} \quad (42)$$

where A has the structure described in Eq. 8.

Since Φ is a strong solution of the algebraic Riccati equation, we know that $|\lambda_i(\hat{A})| \leq 1$. It can be easily checked that $|\lambda_i(\hat{A})| \geq 1$ is a consequence of this.

With T the transformation in Eq. 42 and $W = T^{-1}\hat{W}T^{-T} > 0$, observe that one can express \hat{X}_j as $\hat{X}_j = TX_jT^T$, where X_j are iterates defined by the equation

$$X_{j+1} = AX_jA^T + W, \quad X_0 = T^{-1}\hat{X}_0T^{-T}. \quad (43)$$

We next define a sequence of matrices $\{Y_j\}$ which under-bounds $\{X_j\}$.

$$Y_{j+1} = AY_jA^T + \lambda_{\min}(W)I, \quad Y_0 = 0. \quad (44)$$

It is trivial to show by induction that $Y_j \leq X_j$ for all j . Thus if we can show that $\{Y_j\}$ diverges, divergence of $\{X_j\}$ and $\{\hat{X}_j\}$ follow.

A closed-form expression for Y_j can be found immediately:

$$Y_j = \lambda_{\min}(W) \left\{ \sum_{l=0}^{j-1} A^l (A^T)^l \right\} \quad (45)$$

$$= \lambda_{\min}(W) \text{diag}\{S_1(j), \dots, S_p(j)\} \quad (46)$$

where $S_q(j)$ is defined in Eq. 13.

Thus the set of eigenvalues of Y_j is simply the union of all the eigenvalues of $S_q(k)$ for all q . A straightforward but lengthy argument employing item 2 of Lemma 2.3 then establishes item 3 a) of the lemma.

Provided $\mu \geq \lambda_{\max}(X_1)$, the sequence of matrices $\{Z_j\}$ defined below over-bounds \hat{X}_j :

$$Z_{j+1} = AZ_jA^T + \mu I, \quad Z_0 = 0. \quad (47)$$

It is trivial to show by induction that $Z_j \geq X_j$ for all $j \geq 1$.

In an identical manner to that employed in investigating Y_j , one can deduce the following expression for Z_j : $Z_j = \mu \text{diag}\{S_1(j), \dots, S_p(j)\}$. A straightforward but lengthy argument employing item 1 of Lemma 2.3 yields part 3 b) of the lemma. \square

3.3 Proof of the Main Theorem.

Having established convergence and the associated rate under the preliminary assumptions A.3, A.4 and $\Phi_0 > 0$ of Lemma 3.2, we now successively relax each of these assumptions to give Theorem 1.1. It has been shown⁸ (albeit by different means to those proposed here) that in the case of Kalman filtering problems, these assumptions can be relaxed to extend previously established convergence results.⁷ Here we

consider the more general class of spectral matrices given in Eq. 2 and present a proof of convergence which as well as relaxing these assumptions, also enables statements to be made concerning the convergence rate of the RDE.

Relaxing assumption A.3 (that (F, G) is controllable.)

This assumption has previously been relaxed⁸ via a sequence of perturbations on the original problem, each of which has (F, G) controllable. The emphasis in the present paper is to investigate the structure of RDE iterates associated with the stable and uncontrollable modes of (F, G) . These observations give rise to statements concerning the convergence rate.

We assume now that (F, G) is stabilizable and that, without loss of generality, $F = \begin{pmatrix} F_{11} & F_{12} \\ 0 & F_{22} \end{pmatrix}$ and $G = \begin{pmatrix} G_1 \\ 0 \end{pmatrix}$, where $|\lambda_i(F_{22})| < 1$ and (F_{11}, G_1) is a controllable pair. We partition Φ_k and V conformally: $\Phi_k = \begin{pmatrix} \Phi_{11}^k & \Phi_{12}^k \\ \Phi_{12}^k & \Phi_{22}^k \end{pmatrix}$ and $V = \begin{pmatrix} V_{11} & V_{12} \\ V_{12}^T & V_{22} \end{pmatrix}$. Expression of Eq. 5 in terms of this partitioning reveals that Φ_{11}^k satisfies the Riccati difference equation

$$\Phi_{11}^{k+1} = F_{11}^T (\Phi_{11}^k - \Phi_{11}^k G_1 (U + G_1^T \Phi_{11}^k G_1)^{-1} G_1^T \Phi_{11}^k) F_{11} + V_{11}. \quad (48)$$

With conformal partitioning of Φ (the strong solution of Eq. 3), it can be readily shown that Φ_{11} is a strong solution of the algebraic equation

$$\Phi_{11} = F_{11}^T (\Phi_{11} - \Phi_{11} G_1 (U + G_1^T \Phi_{11} G_1)^{-1} G_1^T \Phi_{11}) F_{11} + V_{11} \quad (49)$$

in the sense that the following matrix only has eigenvalues with magnitude less than or equal to unity: $\hat{F}_{11} = (I - G_1 N_{11}^{-1} G_1^T \Phi_{11}) F_{11}$ where $N_{11} = N = U + G_1^T \Phi_{11} G_1$.

Observe that in fact $\Psi(z) = U + G_1^T (z^{-1} I - F_1^T)^{-1} V_{11} (z I - F_1)^{-1} G_1$. Recall also that for the moment, we maintain the assumption that F is invertible, from which it follows that F_1 is also invertible. Since we also assume that $\Phi_0 > \Phi$ and therefore that $\Delta_0 > 0$, it follows that $\Delta_{11}^0 > 0$. Since (F_{11}, G_1) is controllable, we can apply Lemma 3.2 to deduce that iterates of the reduced-order RDE Eq. 48 satisfy

$$\Delta_{11}^k = \mathcal{O}\left(\frac{1}{k}\right). \quad (50)$$

It also follows from Eq. 5 that the partitions Φ_{12}^k of the iterates Φ_k satisfy

$$\Phi_{12}^{k+1} = (\hat{F}_{11}^k)^T \Phi_{12}^k F_{22} + W_{12}^k \quad (51)$$

where

$$\hat{F}_{11}^k = (I - G_1 (N_{11}^k)^{-1} G_1^T \Phi_{11}^k) F_{11} \quad (52)$$

$$N_{11}^k = U + G_1^T \Phi_{11}^k G_1 \quad (53)$$

$$W_{12}^k = \hat{F}_{11}^k \Phi_{11}^k F_{12} + V_{12}. \quad (54)$$

Lemma 3.3 Let $\{\Upsilon_k\}$ be a bounded sequence of matrices defined for $k \geq 0$. Consider the linear matrix difference equation with (possibly non-square) iterates Ξ_k having a finite initial condition Ξ_0 :

$$\Xi_{k+1} = A_k \Xi_k B_k + \Upsilon_k. \quad (55)$$

Suppose the (square) matrix sequences $\{A_k\}$ and $\{B_k\}$ are such that $A_k \rightarrow A$ and $B_k \rightarrow B$ where $\|\lambda_i(A)\lambda_j(B)\| < 1$ for all i and j . Then if $\Upsilon_k = \mathcal{O}(\frac{1}{k})$, it follows that

$$\Xi_k = \mathcal{O}(\frac{1}{k}). \quad (56)$$

Proof: A number of standard stability results for difference equations can be applied to show this result. \square

Recall from Eq. 50 that $\Phi_{11}^k = \Phi_{11} + \mathcal{O}(\frac{1}{k})$. It follows from Eq. 53 that $N_{11}^k = N_{11} + \mathcal{O}(\frac{1}{k})$, from Eq. 52 that $\hat{F}_{11}^k = \hat{F}_{11} + \mathcal{O}(\frac{1}{k})$ and hence that $W_{12}^k = W_{12} + \mathcal{O}(\frac{1}{k})$ where $W_{12} = \hat{F}_{11}^T \Phi_{11} F_{12} + V_{12}$.

Observe that by hypothesis there exists a solution Φ_{12} of the algebraic equation

$$\Phi_{12} = \hat{F}_{11}^T \Phi_{12} F_{22} + W_{12}. \quad (57)$$

Subtracting this equation from Eq. 51 and simultaneously adding and subtracting the term $(\hat{F}_{11}^k)^T \Phi_{12} F_{22}$ yields the equation

$$\Delta_{12}^{k+1} = (\hat{F}_{11}^k)^T \Delta_{12}^k F_{22} + (\hat{F}_{11}^k - \hat{F}_{11})^T \Phi_{12} F_{22} + W_{12}^k - W_{12}. \quad (58)$$

Recall that \hat{F}_{11} has all eigenvalues in the closed unit circle. Observe that F_{22} is stable. It follows that $\|\lambda_i(\hat{F}_{11})\lambda_j(F_{22})\| < 1$. We now identify Ξ_k with Δ_{12}^k , A_k with $(\hat{F}_{11}^k)^T$, B_k with F_{22} and Υ_k with the remaining terms in Eq. 58, which can be easily shown to be $\mathcal{O}(\frac{1}{k})$. We now apply Lemma 3.3 to Eq. 58 to conclude that

$$\Delta_{12}^k = \mathcal{O}(\frac{1}{k}). \quad (59)$$

Note that examination of the (2, 2) partition of Eq. 5 reveals the following iteration:

$$\Phi_{22}^{k+1} = F_{22}^T \Phi_{22}^k F_{22} + S_{22}^k \quad (60)$$

where

$$S_{22}^k = F_{12}^T \Phi_{11}^k \hat{F}_{12}^k + (\hat{F}_{12}^k)^T \Phi_{12}^k F_{22} + F_{22}^T (\Phi_{12}^k)^T (\hat{F}_{12}^k - G_1 N_{11}^{-k} G_1^T \Phi_{12}^k F_{22}) + V_{22} \quad (61)$$

$$\hat{F}_{12}^k = (I - G_1 (N_{11}^k)^{-1} G_1^T \Phi_{11}^k) F_{12}. \quad (62)$$

Recall that by hypothesis, there exists a solution Φ_{22} of the equation

$$\Phi_{22} = F_{22}^T \Phi_{22} F_{22} + S_{22} \quad (63)$$

where S_{22} is given by taking the limit of Eq. 61. Subtracting this equation from Eq. 60 yields the equation

$$\Delta_{22}^{k+1} = F_{22}^T \Delta_{22}^k F_{22} + S_{22}^k - S_{22}. \quad (64)$$

From Eq. 50 and Eq. 59 it follows that $S_{22}^k = S_{22} + \mathcal{O}(\frac{1}{k})$. Since F_{22} is stable, we can apply Lemma 3.3 with $A_k = B_k = F_{22}$ and $T_k = S_{22}^k - S_{22}$ to conclude that $\Delta_{22}^k = \mathcal{O}(\frac{1}{k})$.

Since $\Delta_{ij}^k = \mathcal{O}(\frac{1}{k})$ for each partition of Φ_k , it follows that $\Delta_k = \mathcal{O}(\frac{1}{k})$. This establishes the worst-case convergence result in item 3 a) of Theorem 1.1.

We now establish the best-case result in item 3 b). Choose any Φ_0 such that $\Delta_0 > 0$ and note therefore that $\Delta_{11}^0 > 0$. Recall that the invariant zeros of the minimum phase spectral factor $\Omega(z)$ are the eigenvalues of \hat{F} . It is easy to check that in the new basis, \hat{F} has diagonal blocks \hat{F}_{11} and F_{22} . Since F_{22} is stable, all unit circle invariant zeros of $\Psi(z)$ are eigenvalues of \hat{F}_{11} . Note that we can apply item 3 b) of Lemma 3.2 to deduce that there exists a k_0 such that when $k \geq k_0$, $\lambda_{\max}(\Phi_{11}^k - \Phi_{11}) \geq \frac{\epsilon_2}{k}$. Note now that the positive definite matrix Δ_{11}^k is a partition of the larger positive definite matrix Δ_k and hence that $\lambda_{\max}(\Phi_k - \Phi) \geq \lambda_{\max}(\Phi_{11}^k - \Phi_{11}) \geq \frac{\epsilon_2}{k}$, which establishes item 3 b).

First strengthening of Lemma 3.2:

With the additional assumptions A.4 and $\Phi_0 > \Phi$, each item of Theorem 1.1 holds.

Relaxing the assumption: $\Phi_0 > \Phi$.

Suppose now that $\Phi_0 \geq \Phi$ but not $\Phi_0 > \Phi$. It is well known that Eq. 38, the difference equation for Δ_k , holds also when Δ_k is singular. In particular, from item 2 of Lemma 3.1, it follows that $\Delta_k \geq 0$ for all $k \geq 0$ which establishes item 1 of the Lemma statement.

Suppose we have any $\bar{\Phi}_0$ such that $\bar{\Phi}_0 \geq \Phi_0 \geq \Phi$ and $\bar{\Phi}_0 > \Phi$. From item 2 of Lemma 3.1 it follows that $\bar{\Phi}_k \geq \Phi_k \geq \Phi$ for all $k \geq 0$ (where $\bar{\Phi}_k$ are iterates of the RDE with initial condition $\bar{\Phi}_0$). Since $\bar{\Phi}_k \geq \Phi_k$ and the convergence of $\{\bar{\Phi}_k\}$ is guaranteed by item 2 of the first strengthening of Lemma 3.2, item 2 in the theorem statement is established.

Item 3 a) in the first strengthening of Lemma 3.2 establishes a worst-case bound for the convergence rate of $\{\bar{\Phi}_k\}$ which, by virtue of the above observations, guarantees the same convergence rate for $\{\Phi_k\}$ which is stated in item 3 a) of the theorem.

Since the restriction $\Phi_0 > \Phi$ is maintained in item 3 b) of the theorem, clearly this worst-case convergence result of item 3 b) in the first strengthening of Lemma 3.2 remains.

Second strengthening of Lemma 3.2:

With the additional assumption A.4, the statements in Theorem 1.1 hold.

Relaxing assumption A.4 (that F is nonsingular).

If F is singular then $\hat{F} = (I - GN^{-1}G^T\Phi)F$ will be also. Suppose it has a Jordan

canonical form $\tilde{F} = M^{-1}\bar{F}M$ where

$$\tilde{F} = \text{diag} \{ \bar{F} \quad F_z \} \quad (65)$$

and \bar{F} and F_z are block diagonal and contain Jordan blocks corresponding to the non-zero and zero eigenvalues of \tilde{F} , respectively.

Note that the difference equation for Δ_k given in Eq. 38 still holds under the new assumptions (i.e., those stated in Theorem 1.1). Now express this equation in the coordinate basis introduced above and define $\tilde{\Delta}_k = M^T \Delta_k M$. Let \bar{F}_q be any Jordan block of \bar{F} of the form in Eq. 11, corresponding to a zero eigenvalue. We now investigate the RDE evolution in the subspace corresponding to this Jordan block. It can be shown fairly easily via Eq. 38 that

$$(\tilde{\Delta}_k)_q = \text{diag} \{ D_k \quad O_{n_q} \} \quad (66)$$

where $(\tilde{\Delta}_k)_q$ is the $n_q \times n_q$ diagonal sub-block of $\tilde{\Delta}_k$ and $D_k \in \mathbb{R}^{(n_q-k) \times (n_q-k)}$ is a nonzero matrix in general. Observe that $(\tilde{\Delta}_k)_q = O_{n_q}$ for all iterations $k \geq n_q$. This reasoning can be applied to each Jordan block which has a zero eigenvalue. It follows that there exists an integer n_z (the size of the largest zero-eigenvalue Jordan block) such that when $k \geq n_z$,

$$\tilde{\Delta}_k = \text{diag} \{ \tilde{\Delta}_k \quad O_{n_z} \} \quad (67)$$

where $\tilde{\Delta}_k \geq 0$ and n_z is the size of the whole invariant subspace corresponding to an eigenvalue of zero. Now define a lower dimensional problem by considering iterates of $\tilde{\Delta}_k$ only. Let G in this basis be partitioned as follows: $G^T = (\tilde{G}^T \quad G_z^T)$.

Note that for $n \geq n_z$, it follows from Eq. 38 that $\tilde{\Delta}_k$ satisfies the following recursion:

$$\tilde{\Delta}_{k+1} = \tilde{F}^T \tilde{\Delta}_k \tilde{F} - \tilde{F}^T \tilde{\Delta}_k \tilde{G} (\tilde{G}^T \tilde{\Delta}_k \tilde{G} + G_z^T \Phi G_z^T + U)^{-1} \tilde{G}^T \tilde{\Delta}_k \tilde{F}. \quad (68)$$

Now consider the above RDE as being associated with the factorization of a spectral matrix $\tilde{\Psi}(z) = U + G^T \Phi G$ of the form in Eq. 2 with " F " replaced by \tilde{F} , " G " replaced by \tilde{G} , " U " replaced by $U + G_z^T \Phi G_z$ and " V " replaced by zero. Clearly the factorization of this matrix from the original state-space realization is trivial and the strong solution of the algebraic equation associated with Eq. 68 is $\tilde{\Delta} = 0$.

It can be easily checked that if (F, G) is stabilizable then (\tilde{F}, \tilde{G}) is also. Since we assume that $\Delta_0 \geq 0$, it follows that $\tilde{\Delta}_0 \geq 0$. By construction, \tilde{F} is invertible. Convergence of $\tilde{\Delta}_k$ then follows by application of the second strengthening of Lemma 3.2 to $\tilde{\Psi}(z)$.

Recall that the parts of the RDE iterates associated with invariant subspaces corresponding to zero eigenvalues converge in a finite number of iterations. The best and worst case convergence behaviour of $\tilde{\Delta}_k$ are therefore inherited by Δ_k , as stated in items 3 a) and 3 b) of the theorem. \square

Acknowledgements.

The authors wish to gratefully acknowledge helpful correspondence and discussions with Professor F.M. Callier, Dr. Robert Bitmead, Dr. Michael Green, Dr. Michel Gevers and Dr. Leonid Gurvits.

References

- [1] B.D.O. Anderson, K.L. Hitz, N.D. Diem, Recursive Algorithm for Spectral Factorization, *IEEE Trans. Circuits Syst.*, vol. 21, no. 6, 1974, pp. 742-750.
- [2] B.D.O. Anderson, J.B. Moore, *Optimal Control - Linear Quadratic Methods*, Prentice-Hall, Englewood Cliffs, NJ, 1990.
- [3] B.D.O. Anderson, J.B. Moore, *Optimal Filtering*, Prentice-Hall, Englewood Cliffs, NJ, 1979.
- [4] B.D.O. Anderson, Stability Properties of Kalman-Bucy Filters, *Journal of the Franklin Institute*, vol. 291, No. 2, 1971, pp. 137-144.
- [5] R.R. Bitmead, M. Gevers, Riccati Difference and Differential Equations: Convergence, Monotonicity and Stability, *The Riccati Equation*, S. Bittanti, A.J. Laub, J.C. Willems (eds.), Springer-Verlag, Berlin, 1991.
- [6] P.E. Caines, *Linear Stochastic Systems*, Wiley, New York, 1988.
- [7] S.W. Chan, G.C. Goodwin, K.S. Sin, Convergence Properties of the Riccati Difference Equation in Optimal Filtering of Nonstabilizable Systems, *IEEE Trans. Automat. Contr.*, vol. AC-29, February 1984, pp. 110-118.
- [8] C.E. de Souza, M.R. Gevers, G.C. Goodwin, Riccati Equations in Optimal Filtering of Nonstabilizable Systems Having Singular State Transition Matrices, *IEEE Trans. Automat. Contr.*, vol. AC-31, September 1986, pp. 831-838.
- [9] C.E. de Souza, On Stabilizing Properties of Solutions of the Riccati Difference Equation, *IEEE Trans. Automat. Contr.*, vol. 34, September 1989, pp. 1313-1316.
- [10] C.E. de Souza and L. Xie, On the Discrete-time Bounded Real Lemma with application in the characterization of static state feedback H_∞ controllers, *Systems & Control Letters*, vol. 18, 1992, pp. 61-71.
- [11] B.P. Molinari, The Stabilizing Solution of the Discrete Algebraic Riccati Equation, *IEEE Trans. Automat. Contr.*, vol. AC-20, June 1975, pp. 396-399.
- [12] W.M. Wonham, *Linear Multivariable Control. A Geometric Approach*, Springer-Verlag, New York, 1974.