

Note on Finding an Optimal Deflation for Quadratic Matrix Polynomials

Xin Liang*

Yau Mathematical Sciences Center, Tsinghua University, Beijing 100084, China.

Received 14 May 2020; Accepted 6 October 2020

Abstract. This paper is concerned with the way to find an optimal deflation for the eigenvalue problem associated with quadratic matrix polynomials. This work is a response of the work by Tisseur et al., *Linear Algebra Appl.*, 435:464-479, 2011, and solves one of open problems raised by them. We build an equivalent unconstrained optimization problem on eigenvalues of a hyperbolic quadratic matrix polynomial of order 2, and develop a technique that transforms the quadratic matrix polynomial to an equivalent one that is easy to solve. Numerical tests are given to illustrate several properties of the problem.

AMS subject classifications: 65F35, 65F15

Key words: deflation, quadratic matrix polynomials, hyperbolic, eigenvalue optimization.

1 Introduction

Given a quadratic matrix polynomial

$$Q(\lambda) = \lambda^2 M + \lambda C + K,$$

where $M, C, K \in \mathbb{R}^{n \times n}$ with M nonsingular. Its associated quadratic eigenvalue problem is

$$Q(\lambda)x = 0, \quad y^H Q(\lambda) = 0,$$

where λ is an eigenvalue and x, y are its corresponding (right) eigenvector and left eigenvector respectively. An eigenvalue is of positive type, if $y^H Q'(\lambda)x = y^H(2\lambda M + C)x > 0$; An eigenvalue is of negative type, if $y^H Q'(\lambda)x = y^H(2\lambda M + C)x < 0$.

Suppose that $\lambda(Q)$, the spectra of $Q(\lambda)$, is $\{\lambda_1, \dots, \lambda_{2n}\}$. Deflating two distinct eigenvalues λ_1, λ_2 is to construct a new quadratic matrix polynomial

$$\tilde{Q}(\lambda) = \begin{bmatrix} Q_d(\lambda) & \\ & q(\lambda) \end{bmatrix} = \lambda^2 \begin{bmatrix} M_d & \\ & m \end{bmatrix} + \lambda \begin{bmatrix} C_d & \\ & c \end{bmatrix} + \begin{bmatrix} K_d & \\ & k \end{bmatrix}$$

*Corresponding author. *Email address:* liangxinslm@tsinghua.edu.cn (X. Liang)

such that $\lambda(q) = \{\lambda_1, \lambda_2\}$, $\lambda(Q_d) = \{\lambda_3, \dots, \lambda_{2n}\}$. Usually, in the two eigenvalues, one is of positive type, and the other is of negative type. If M, C, K are symmetric, and $\lambda_1 \in \lambda(Q)$ but λ_1 is nonreal, then $\overline{\lambda_1} \in \lambda(Q)$, and in this case, it is usually required to deflate this conjugate pair together.

The deflation technique is very useful and popular in computing the eigenvalues of a matrix, so that it is hoped to be used for computing the eigenvalues of quadratic matrix polynomials. However, as far as we know, not many works discussed on this topic. Meini [5] discussed a deflation method coupled in her so-called "shift-and-deflate" technique. Tisseur et al. [6] presented a general way to deflate two distinct eigenvalues of quadratic matrix polynomials if the corresponding eigenvectors are given.

Here we briefly describe the idea of the method introduced by Tisseur et al. First they developed a method to deflate a quadratic polynomial for two given eigenvalues whose eigenvectors are parallel. Then they invented a way to transform a quadratic polynomial for two given eigenvalues whose eigenvectors are nonparallel into a new one that has two eigenvalues whose eigenvectors are parallel, i.e., transform this case into the solved case.

The new quadratic matrix polynomial produced by the deflation may have a significantly large condition number compared to the original quadratic matrix polynomial. For the special case that M, C, K are symmetric, Tisseur et al. gave an optimal choice to minimize the condition number for the parallel case, but for the nonparallel case, in [6, Section 3], they reported:

Identifying which solution minimizes the condition number $\kappa_2(T) = \|T\|_2 \|T^{-1}\|_2$ remains an open problem.

Here T is a related transformation matrix, of which the detailed form will be given below.

The aim of this paper is to solve this problem. First this problem is formulated and simplified in Section 2, which induces a constrained optimization problem on the eigenvalues of a hyperbolic quadratic matrix polynomial of order 2. Next we parameterize (or equivalently nondimensionalize) it and obtain an unconstrained optimization problem in Section 3. Then we calculate the gradient and the Hessian matrix of the objective function in Section 4. Then we make several numerical tests to show the properties of this problem and suggest a technique that transforms it to an equivalent problem whose objective function is easy to solve, as is shown in Section 5. Finally some concluding remarks is given in Section 6.

Notation. Throughout this paper, I_n (or simply I if its dimension is clear from the context) is the $n \times n$ identity matrix. For any scalar, vector, or matrix X , $\Re X$ and $\Im X$ are its real part and imaginary part respectively; while $\|X\|_2$ and $\|X\|_\infty$ are its spectral norm and sum-of-row norm. For any matrix X , $\lambda(X)$ represents its spectra, and $\lambda_*(X)$ represents the set consisting of all its nonzero eigenvalues. For any real symmetric matrix X , $X \succ 0$ ($X \succeq 0$) means that X is positive (semi-)definite, and $X \prec 0$ ($X \preceq 0$) if $-X \succ 0$ ($-X \succeq 0$).

By $X \otimes Y$ we denote the Kronecker product of two matrices X, Y . For any set S , S^\perp is its orthogonal complement.

2 Simplify the problem

First we state the problem.

Given a symmetric quadratic matrix polynomial $Q(\lambda) = \lambda^2 M + \lambda C + K$ with M non-singular, where (λ_1, x_1) and (λ_2, x_2) are its two eigenpairs of opposite types. Write

$$(\Lambda, X) = \begin{cases} \left(\begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}, \begin{bmatrix} x_1 & x_2 \end{bmatrix} \right), & \text{if } \lambda_1 \text{ and } \lambda_2 \text{ are real,} \\ \left(\begin{bmatrix} \Re\lambda_1 & \Im\lambda_1 \\ -\Im\lambda_1 & \Re\lambda_1 \end{bmatrix}, \begin{bmatrix} \Re x_1 & \Im x_1 \end{bmatrix} \right), & \text{if } \lambda_1 = \overline{\lambda_2} \text{ are nonreal and } x_1 = \overline{x_2}. \end{cases}$$

If x_1, x_2 are nonparallel vectors, then the transformation T is

$$T = I_{2n} + \begin{bmatrix} ab^T & ad^T \\ af^T & ah^T \end{bmatrix},$$

where

$$a = \frac{Xp}{\|Xp\|_2}, \quad [b \ f \ d \ h] = \left(I_n - \frac{zz^T}{z^T z} \right) B A^+ + U(I_4 - A A^+) + \frac{zw^T}{z^T z} =: V \in \mathbb{R}^{n \times 4}.$$

Here U can be any matrix with $z^T U = 0$, and

$$A = \frac{1}{2} \begin{bmatrix} 2\alpha_M & \alpha_C & 0 \\ \alpha_C & 2\alpha_K & 0 \\ 0 & 2\alpha_M & \alpha_C \\ 0 & \alpha_C & 2\alpha_K \end{bmatrix}, \quad B = -[Ma \ Ca \ Ka],$$

where $\alpha_M, \alpha_C, \alpha_K$ are $a^T M a, a^T C a, a^T K a$ respectively. Also $z \in \mathbb{R}^n, w \in \mathbb{R}^4$ are given by

$$z = X \Lambda p - \frac{e_\ell^T X \Lambda p - 1}{e_\ell^T a} a, \quad w = \frac{1}{e_\ell^T a} \begin{bmatrix} e_\ell^T X \Lambda p - 1 \\ e_\ell^T a \|Xp\|_2 \\ e_\ell^T X \Lambda q \\ e_\ell^T X q - 1 \end{bmatrix},$$

where e_ℓ is a column of I_n that makes $|e_\ell^T a| = \|a\|_\infty$, and

$$p = \frac{\gamma}{\lambda_1 - \lambda_2} \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad \gamma = 1 \text{ or } i \text{ to make } p \text{ real,} \quad q = \Lambda p - (\lambda_1 + \lambda_2)p.$$

The optimization problem to be solved is

$$\min_{U \in \mathbb{R}^{n \times 4}; z^T U = 0} \kappa_2(T). \tag{2.1}$$

Besides, a, V, A, B, z, w have the properties below, which has been shown in [6]:

$$z^T V = w^T, \quad VA = B, \quad w^T A = z^T B; \tag{2.2}$$

$$a^T B + a_0^T A = 0 \quad \text{where} \quad a_0 = [1 \ 0 \ 0 \ 1]^T; \tag{2.3}$$

$$A \text{ has full column rank.} \tag{2.4}$$

In the following, we will simplify the problem.

Define a linear mapping

$$\begin{aligned} \text{rs:} \quad \mathbb{R}^{n \times 2m} &\rightarrow \mathbb{R}^{2n \times m} \\ Z = [Z_1 \ Z_2] &\mapsto \text{rs}(Z) = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix}, \end{aligned}$$

which satisfies

$$\text{rs}(XY) = (I_2 \otimes X)\text{rs}(Y), \quad \forall X, Y. \tag{2.5}$$

Let

$$A_0 = I_2 \otimes a = \begin{bmatrix} a & 0 \\ 0 & a \end{bmatrix}, \quad B_0 = \text{rs}(V) = \begin{bmatrix} b & f \\ d & h \end{bmatrix},$$

and then

$$T = I_{2n} + A_0 B_0^T.$$

To solve the optimization problem equation (2.1), we need to compute $\|T\|_2$ and $\|T^{-1}\|_2$ first, or equivalently, the minimal and maximal eigenvalues of TT^T . Note that

$$\begin{aligned} TT^T &= I_{2n} + A_0 B_0^T + B_0 A_0^T + A_0 B_0^T B_0 A_0^T \\ &= I_{2n} + [A_0 \ B_0] \begin{bmatrix} B_0^T B_0 & I_2 \\ I_2 & 0 \end{bmatrix} [A_0 \ B_0]^T. \end{aligned}$$

Since $A_0^T A_0 = (a^T a) I_2 = I_2$, noticing $\lambda_*(XY) = \lambda_*(YX), \forall X, Y$,

$$\begin{aligned} \lambda_*(TT^T - I_{2n}) &= \lambda_* \left([A_0 \ B_0] \begin{bmatrix} B_0^T B_0 & I_2 \\ I_2 & 0 \end{bmatrix} [A_0 \ B_0]^T \right) \\ &= \lambda_* \left(\begin{bmatrix} B_0^T B_0 & I_2 \\ I_2 & 0 \end{bmatrix} [A_0 \ B_0]^T [A_0 \ B_0] \right) \end{aligned}$$

$$\begin{aligned}
 &= \lambda_* \left(\begin{bmatrix} B_0^T B_0 & I_2 \\ I_2 & 0 \end{bmatrix} \begin{bmatrix} I_2 & A_0^T B_0 \\ B_0^T A_0 & B_0^T B_0 \end{bmatrix} \right) \\
 &= \lambda_* \left(\begin{bmatrix} B_0^T B_0 & I_2 \\ I_2 & 0 \end{bmatrix} \begin{bmatrix} I_2 & 0 \\ B_0^T A_0 & I_2 \end{bmatrix} \begin{bmatrix} I_2 & 0 \\ 0 & B_0^T B_0 - B_0^T A_0 A_0^T B_0 \end{bmatrix} \begin{bmatrix} I_2 & A_0^T B_0 \\ 0 & I_2 \end{bmatrix} \right) \\
 &= \lambda_* \left(\begin{bmatrix} I_2 & A_0^T B_0 \\ 0 & I_2 \end{bmatrix} \begin{bmatrix} B_0^T B_0 & I_2 \\ I_2 & 0 \end{bmatrix} \begin{bmatrix} I_2 & 0 \\ B_0^T A_0 & I_2 \end{bmatrix} \begin{bmatrix} I_2 & 0 \\ 0 & B_0^T B_0 - B_0^T A_0 A_0^T B_0 \end{bmatrix} \right) \\
 &= \lambda_* \left(\begin{bmatrix} B_0^T B_0 + B_0^T A_0 + A_0^T B_0 & B_0^T B_0 - B_0^T A_0 A_0^T B_0 \\ I_2 & 0 \end{bmatrix} \right).
 \end{aligned}$$

Due to the relationship between the eigenvalues of a quadratic matrix polynomial and its linearization, $\lambda_*(TT^T - I_{2n}) = \lambda(\tilde{H})$, where

$$\tilde{H}(\lambda) = -\lambda^2 I_2 + \lambda(B_0^T B_0 + B_0^T A_0 + A_0^T B_0) + (B_0^T B_0 - B_0^T A_0 A_0^T B_0).$$

Note that

$$\begin{aligned}
 H(\lambda) &:= -\tilde{H}(\lambda - 1) \\
 &= (\lambda - 1)^2 I_2 - (\lambda - 1)(B_0^T B_0 + B_0^T A_0 + A_0^T B_0) - (B_0^T B_0 - B_0^T A_0 A_0^T B_0) \\
 &= \lambda^2 I_2 - \lambda(B_0^T B_0 + B_0^T A_0 + A_0^T B_0 + 2I_2) + (I_2 + B_0^T A_0 + A_0^T B_0 + B_0^T A_0 A_0^T B_0) \\
 &= \lambda^2 I_2 - \lambda(B_0^T B_0 + B_0^T A_0 + A_0^T B_0 + 2I_2) + (I_2 + A_0^T B_0)^T (I_2 + A_0^T B_0). \tag{2.6}
 \end{aligned}$$

We can see the eigenvalues of H are also the eigenvalues of TT^T , and the other eigenvalues of TT^T are 1. Since $A_0^T A_0 = I_2$,

$$H(1) = -B_0^T B_0 + B_0^T A_0 A_0^T B_0 = -B_0^T (I_{2n} - A_0 A_0^T) B_0 \leq 0. \tag{2.7}$$

Hence $H(\lambda)$ is a semi-hyperbolic quadratic matrix polynomial (see, e.g., [2]), which implies

$$\lambda(H) = \{\lambda_i, i = 1, 2, 3, 4\} \quad \text{satisfying} \quad 0 \leq \lambda_1 \leq \lambda_2 \leq 1 \leq \lambda_3 \leq \lambda_4.$$

Here $\lambda_1 \geq 0$ is guaranteed by the fact λ_1 is an eigenvalue of TT^T . Thus,

$$\kappa_2(T) = \sqrt{\frac{\lambda_4}{\lambda_1}}.$$

Now the problem is to compute the maximal and minimal eigenvalues of H . Note that by Eqs. (2.6) and (2.7),

$$H(\lambda) = \lambda^2 I_2 - \lambda \left[(I_2 + A_0^T B_0)^T (I_2 + A_0^T B_0) + I_2 - H(1) \right] + (I_2 + A_0^T B_0)^T (I_2 + A_0^T B_0). \tag{2.8}$$

The next step is to calculate $A_0^T B_0$ and $B_0^T B_0$.

First, we simplify V . Write

$$a_{\perp} = \frac{1}{2} [\alpha_C \quad -2\alpha_M \quad 2\alpha_K \quad -\alpha_C]^T.$$

Then it is easy to check $a_{\perp}^T A = 0$. By Eq. (2.4), namely the fact A has full column rank, $A^+ = (A^T A)^{-1} A^T$, and $I_4 - AA^+ = \frac{a_{\perp} a_{\perp}^T}{a_{\perp}^T a_{\perp}}$. Thus, by Eq. (2.2),

$$\begin{aligned} V &= BA^+ - \frac{zz^T}{z^T z} BA^+ + U \frac{a_{\perp} a_{\perp}^T}{a_{\perp}^T a_{\perp}} + \frac{zw^T}{z^T z} \\ &= BA^+ + \frac{zw^T}{z^T z} (I_4 - AA^+) + U \frac{a_{\perp} a_{\perp}^T}{a_{\perp}^T a_{\perp}} \\ &= BA^+ + \left(\frac{zw^T}{z^T z} + U \right) \frac{a_{\perp} a_{\perp}^T}{a_{\perp}^T a_{\perp}}. \end{aligned}$$

Note that

$$\{Ua_{\perp} : z^T U = 0\} = \{z\}^{\perp} =: \mathcal{U}.$$

Then

$$V = BA^+ + va_{\perp}^T, \tag{2.9}$$

where

$$v = \frac{w^T a_{\perp}}{a_{\perp}^T a_{\perp} z^T z} z + u, \quad u \in \mathcal{U},$$

or equivalently,

$$v \in \{v \in \mathbb{R}^n : v^T z a_{\perp}^T a_{\perp} = w^T a_{\perp}\}. \tag{2.10}$$

Then we simplify $A_0^T B_0$. By Eq. (2.3),

$$a^T BA^+ = -a_0^T AA^+ = -a_0^T \left(I_4 - \frac{a_{\perp} a_{\perp}^T}{a_{\perp}^T a_{\perp}} \right) = -a_0^T. \tag{2.11}$$

Thus,

$$\begin{aligned} A_0^T B_0 &= (I_2 \otimes a)^T \text{rs}(V) && \text{by Eq. (2.5)} \\ &= \text{rs}(a^T V) && \text{by Eq. (2.9)} \\ &= \text{rs}(a^T BA^+ + a^T va_{\perp}^T) && \text{by Eq. (2.11)} \\ &= \text{rs}(-a_0^T) + a^T v \text{rs}(a_{\perp}^T). \end{aligned}$$

Write

$$S = \frac{1}{2} \begin{bmatrix} 2\alpha_M & \alpha_C \\ \alpha_C & 2\alpha_K \end{bmatrix}, \quad J = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}.$$

It is easy to see that $\text{rs}(a_{\perp}^T) = SJ$. Therefore,

$$A_0^T B_0 = -I_2 + a^T v S J. \tag{2.12}$$

Next we simplify $B_0^T B_0$. By Eq. (2.5),

$$B_0 = \text{rs}(V) = \text{rs}(B(A^T A)^{-1} A^T) + \text{rs}(v a_{\perp}^T) = (I_2 \otimes B(A^T A)^{-1}) \text{rs}(A^T) + (I_2 \otimes v) \text{rs}(a_{\perp}^T).$$

It is easy to see that $\text{rs}(A^T) = (p_0 \otimes I_2) S$ where $p_0 = [1 \ 0 \ 1]^T = e_1 + e_3$. Thus

$$B_0 = \text{rs}(V) = (I_2 \otimes B(A^T A)^{-1})(p_0 \otimes I_2) S + (I_2 \otimes v) S J =: (I_2 \otimes v) S J + D S,$$

where

$$D = (I_2 \otimes B(A^T A)^{-1})(p_0 \otimes I_2). \tag{2.13}$$

Besides, by Eq. (2.11),

$$(I_2 \otimes a^T) D S = (I_2 \otimes a^T) \text{rs}(B A^+) = \text{rs}(a^T B A^+) = -I_2. \tag{2.14}$$

Then

$$\begin{aligned} B_0^T B_0 &= [(I_2 \otimes v) S J + D S]^T [(I_2 \otimes v) S J + D S] \\ &= -J S (I_2 \otimes v^T v) S J + S D^T (I_2 \otimes v) S J - J S (I_2 \otimes v^T) D S + S D^T D S \\ &= -v^T v J S^2 J + S D^T (I_2 \otimes v) S J - J S (I_2 \otimes v^T) D S + S D^T D S, \end{aligned}$$

and by Eqs. (2.7) and (2.12),

$$\begin{aligned} H_1 &:= -J H(1) J \\ &= -[v^T v - (a^T v)^2] S^2 - J S D^T (I_2 \otimes v) S + S (I_2 \otimes v^T) D S J + J S D^T D S J + I_2 - a^T v (J S - S J). \end{aligned}$$

Turn back to Eq. (2.8). Substituting Eq. (2.12),

$$\begin{aligned} H(\lambda) &= \lambda^2 I_2 - \lambda [I_2 - (a^T v)^2 J S^2 J - H(1)] - (a^T v)^2 J S^2 J \\ &= -J \left(\lambda^2 I_2 - \lambda [I_2 + (a^T v)^2 S^2 - H_1] + (a^T v)^2 S^2 \right) J. \end{aligned}$$

To sum up, the optimization problem equation (2.1) is equivalent to:

$$\min_{v \in \mathbb{R}^n: v^T z a_{\perp}^T a_{\perp} = w^T a_{\perp}} \kappa_2(T) = \min_{v \in \mathbb{R}^n: v^T z a_{\perp}^T a_{\perp} = w^T a_{\perp}} \sqrt{\frac{\lambda_4}{\lambda_1}}, \tag{2.15}$$

where $0 \leq \lambda_1 \leq \lambda_2 \leq 1 \leq \lambda_3 \leq \lambda_4$ are four eigenvalues of the hyperbolic quadratic polynomial

$$-J H(\lambda) J = \lambda^2 I_2 - \lambda [I_2 + (a^T v)^2 S^2 - H_1] + (a^T v)^2 S^2$$

and

$$H_1 := I_2 - [v^T v - (a^T v)^2] S^2 - J S D^T (I_2 \otimes v) S + S (I_2 \otimes v^T) D S J + J S D^T D S J - a^T v (J S - S J). \tag{2.16}$$

3 Parameterize the problem

First we give a lemma to show the relationship between hyperbolic quadratic matrix polynomial and its eigenvalues.

Lemma 3.1. *For any two hyperbolic quadratic matrix polynomial $H^{(i)}(\lambda) = \lambda^2 I + \lambda B^{(i)} + C$ satisfying $B^{(i)} \preceq 0, C \succeq 0$, let $\lambda_{\min}^{(i)}$ and $\lambda_{\max}^{(i)}$ be their minimal and maximal eigenvalues respectively. If $B^{(1)} \succeq B^{(2)}$, then $\lambda_{\min}^{(1)} \geq \lambda_{\min}^{(2)}, \lambda_{\max}^{(1)} \leq \lambda_{\max}^{(2)}$.*

Proof. Since the polynomial is hyperbolic, for any vector x satisfying $x^T x = 1$, we can well define

$$\rho_{\pm}(x) := \frac{1}{2} \left(-x^T Bx \pm \sqrt{(x^T Bx)^2 - 4x^T Cx} \right).$$

Then

$$\frac{\partial \rho_+(x)}{\partial (x^T Bx)} = \frac{1}{2} \left(-1 + \frac{2x^T Bx}{2\sqrt{(x^T Bx)^2 - 4x^T Cx}} \right) = \frac{-\rho_+(x)}{\sqrt{(x^T Bx)^2 - 4x^T Cx}} \leq 0.$$

Thus, $B^{(1)} \succeq B^{(2)}$ implies $x^T B^{(1)}x \geq x^T B^{(2)}x$ for all x . Then, $\rho_+^{(1)}(x) \leq \rho_+^{(2)}(x)$. As a result, by the Courant-Fischer min-max theorem of hyperbolic quadratic matrix polynomials [1], $\lambda_{\max}^{(1)} = \max \rho_+^{(1)}(x) \leq \max \rho_+^{(2)}(x) = \lambda_{\max}^{(2)}$.

Similarly,

$$\frac{\partial \rho_-(x)}{\partial (x^T Bx)} = \frac{1}{2} \left(-1 - \frac{2x^T Bx}{2\sqrt{(x^T Bx)^2 - 4x^T Cx}} \right) = \frac{\rho_-(x)}{\sqrt{(x^T Bx)^2 - 4x^T Cx}} \geq 0.$$

Then, $\rho_-^{(1)}(x) \geq \rho_-^{(2)}(x)$, and $\lambda_{\min}^{(1)} = \min \rho_-^{(1)}(x) \geq \min \rho_-^{(2)}(x) = \lambda_{\min}^{(2)}$. □

We may continue simplifying the optimization problem.

First, we know the minimal v must lie in $\{a, Ma, Ca, Ka\}$. Otherwise, let

$$v = v_0 + v_{\perp} \quad \text{where} \quad v_0 \in \{a, Ma, Ca, Ka\}, \quad v_{\perp} \in \{a, Ma, Ca, Ka\}^{\perp}.$$

Then $a^T v = a^T v_0, B^T v = B^T v_0$, which implies $(I_2 \otimes v^T)D = (I_2 \otimes v_0^T)D$ by Eq. (2.13). However, $v^T v = v_0^T v_0 + v_{\perp}^T v_{\perp} \geq v_0^T v_0$. By Lemma 3.1, $\kappa_2(T; v) \geq \kappa_2(T; v_0)$.

Let

$$B_a = B + a \begin{bmatrix} \alpha_M & \alpha_C & \alpha_K \end{bmatrix}, \quad \tilde{B}_a = B_a (A^T A)^{-1}. \tag{3.1}$$

It is easy to check $a^T B_a = 0$ and $B^T B_a = B_a^T B_a$. Since the minimal v must lie in $\{a, Ma, Ca, Ka\}$, we can write

$$v = \xi a + \tilde{B}_a \tilde{y}_B, \quad \text{where} \quad \xi \in \mathbb{R}, \quad \tilde{y}_B \in \mathbb{R}^3. \tag{3.2}$$

Thus, by Eq. (2.14),

$$a^T v = \xi, \quad v^T v = \xi^2 + \tilde{y}_B^T \tilde{B}_a^T \tilde{B}_a \tilde{y}_B, \quad (I_2 \otimes v^T)DS = -\xi I_2 + (I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T)DS. \tag{3.3}$$

Substituting Eq. (3.3) into Eq. (2.16),

$$H_1 = I_2 - \tilde{y}_B^T \tilde{B}_a^T \tilde{B}_a \tilde{y}_B S^2 - JSD^T(I_2 \otimes \tilde{B}_a \tilde{y}_B)S + S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T)DSJ + JSD^TDSJ.$$

Clearly,

$$I_n = aa^T + \tilde{B}_a \tilde{B}_a^+ + B_\perp B_\perp^+,$$

where B_\perp is a basis of $\{a, Ma, Ca, Ka\}^\perp$. Thus,

$$\begin{aligned} DSJ &= (I_2 \otimes [aa^T + \tilde{B}_a \tilde{B}_a^+ + B_\perp B_\perp^+])DSJ \\ &\quad \text{by Eq. (2.13)} \\ &= (I_2 \otimes aa^T)DSJ + (I_2 \otimes [\tilde{B}_a \tilde{B}_a^+ + B_\perp B_\perp^+])(I_2 \otimes B(A^T A)^{-1})(p_0 \otimes I_2)SJ \\ &\quad \text{by Eq. (2.14)} \\ &= -(I_2 \otimes a)J + (I_2 \otimes \tilde{B}_a \tilde{B}_a^+ B(A^T A)^{-1})(p_0 \otimes I_2)SJ \\ &\quad \text{by Eq. (3.1)} \\ &= -(I_2 \otimes a)J + \left(I_2 \otimes \tilde{B}_a \tilde{B}_a^+ (\tilde{B}_a - a [\alpha_M \quad \alpha_C \quad \alpha_K] (A^T A)^{-1}) \right) (p_0 \otimes I_2)SJ \\ &\quad \text{by } a^T B_a = a^T \tilde{B}_a = 0 \\ &= -(I_2 \otimes a)J + (I_2 \otimes \tilde{B}_a)(p_0 \otimes I_2)SJ. \end{aligned}$$

Then

$$\begin{aligned} -JSD^TDSJ &= -J(I_2 \otimes a^T a)J - JS(p_0^T \otimes I_2)(I_2 \otimes \tilde{B}_a^T \tilde{B}_a)(p_0 \otimes I_2)SJ \\ &= I_2 - JS(p_0^T \otimes I_2)(I_2 \otimes \tilde{B}_a^T \tilde{B}_a)(p_0 \otimes I_2)SJ \\ &=: I_2 + SD_a^T D_a S, \end{aligned}$$

where $D_a = (I_2 \otimes \tilde{B}_a)(p_0 \otimes I_2)SJS^{-1}$, and by $a^T \tilde{B}_a = 0$,

$$\begin{aligned} -S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T)DSJ &= S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T a)J - S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T \tilde{B}_a)(p_0 \otimes I_2)SJ \\ &= -S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T \tilde{B}_a)(p_0 \otimes I_2)SJ \\ &= -S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T)D_a S, \end{aligned}$$

which gives

$$\begin{aligned} H_1 &= I_2 - \tilde{y}_B^T \tilde{B}_a^T \tilde{B}_a \tilde{y}_B S^2 + S(I_2 \otimes \tilde{y}_B^T \tilde{B}_a^T)D_a S + SD_a^T(I_2 \otimes \tilde{B}_a \tilde{y}_B)S - I_2 - SD_a^T D_a S \\ &= -S(D_a - (I_2 \otimes \tilde{B}_a \tilde{y}_B))^T (D_a - (I_2 \otimes \tilde{B}_a \tilde{y}_B))S \\ &= -SD_y^T D_y S, \end{aligned} \tag{3.4}$$

where $D_y = D_a - (I_2 \otimes \tilde{B}_a \tilde{y}_B)$. Note that

$$\begin{aligned} D_a &= (I_2 \otimes \tilde{B}_a)(p_0 \otimes I_2)SJS^{-1} \\ &= (I_2 \otimes B_a(A^T A)^{-1})(p_0 \otimes I_2)SJS^{-1} = (I_2 \otimes B_a)R, \end{aligned}$$

where calculation tells

$$R = (I_2 \otimes (A^T A)^{-1})(p_0 \otimes I_2)SJS^{-1} = \delta \begin{bmatrix} r & -e_1 \\ e_3 & r - e_2 \end{bmatrix},$$

in which $I_3 = [e_1 \ e_2 \ e_3]$, and

$$\delta = \frac{4}{4\alpha_K\alpha_M - \alpha_C^2}, \quad r = \frac{1}{2(2\alpha_M^2 + 2\alpha_K^2 + \alpha_C^2)} \begin{bmatrix} 2\alpha_C(\alpha_M + \alpha_K) \\ 4\alpha_K^2 + \alpha_C^2 \\ -2\alpha_C(\alpha_M + \alpha_K) \end{bmatrix}.$$

Thus, writing $y_B = (A^T A)^{-1}\tilde{y}_B$ and $y = r - \frac{1}{\delta}y_B$,

$$D_y = (I_2 \otimes B_a)R - (I_2 \otimes B_a)y_B = \delta(I_2 \otimes B_a) \begin{bmatrix} y & -e_1 \\ e_3 & y - e_2 \end{bmatrix}.$$

Then, letting $B_W = B_a^T B_a = B^T B - [\alpha_M \ \alpha_C \ \alpha_K]^T [\alpha_M \ \alpha_C \ \alpha_K]$,

$$\begin{aligned} D_y^T D_y &= \delta^2 \begin{bmatrix} y^T B_W y & (e_3 - e_1)^T B_W y \\ (e_3 - e_1)^T B_W y & y^T B_W y - 2e_2^T B_W y \end{bmatrix} + \delta^2 \begin{bmatrix} e_3^T B_W e_3 & -e_2^T B_W e_3 \\ -e_2^T B_W e_3 & e_2^T B_W e_2 + e_1^T B_W e_1 \end{bmatrix} \\ &=: \delta^2(W_y + W_e). \end{aligned} \tag{3.5}$$

Note that $\begin{bmatrix} y & -e_1 \\ e_3 & y - e_2 \end{bmatrix}$ has full column rank, because otherwise $y \parallel e_1, (y - e_2) \parallel e_3$, which is a contradiction. Thus $D_y^T D_y \succ 0$ and then by Eq. (3.4) $H_1 = -SD_y^T D_y S \prec 0$, which implies $\lambda_2 < 1 < \lambda_3$.

On the other hand, by Eq. (3.2),

$$v = \zeta a + B_a y_B = \zeta a + \delta B_a(r - y), \quad \zeta \in \mathbb{R}, \quad y \in \mathbb{R}^3.$$

Then the only constraint on v , namely Eq. (2.10), is

$$\frac{w^T a_\perp}{a_\perp^T a_\perp} = v^T z = \zeta a^T z + \delta z^T B_a r - \delta z^T B_a y,$$

or equivalently,

$$z_a^T y = \zeta + \beta,$$

where

$$z_a = \frac{\delta B_a^T z}{a^T z}, \quad \beta = z_a^T r - \frac{w^T a_\perp}{a^T z a_\perp^T a_\perp}.$$

Also,

$$-JH(\lambda)J = \lambda^2 I_2 - \lambda(I_2 + \delta^2 S(W_y + W_e)S + \zeta^2 S^2) + \zeta^2 S^2.$$

To sum up, the constrained optimization problem Eq. (2.1), or Eq. (2.15), is equivalent to this unconstrained optimization problem:

$$\min_{y \in \mathbb{R}^3} \kappa_2(T) = \min_{y \in \mathbb{R}^3} \sqrt{\frac{\lambda_4}{\lambda_1}} \tag{3.6}$$

where $0 \leq \lambda_1 \leq \lambda_2 < 1 < \lambda_3 \leq \lambda_4$ are four eigenvalues of the hyperbolic quadratic polynomial

$$-JH(\lambda)J = \lambda^2 I_2 - \lambda [I_2 + \delta^2 S W_e S + \delta^2 S W_y S + (z_a^T y - \beta)^2 S^2] + (z_a^T y - \beta)^2 S^2, \tag{3.7}$$

and

$$W_y = (y^T B_W y) I_2 + \begin{bmatrix} 0 & (e_3 - e_1)^T B_W y \\ (e_3 - e_1)^T B_W y & -2e_2^T B_W y \end{bmatrix}.$$

4 Analyze the problem

It might be possible to have an explicit solution of the unconstrained optimization problem Eq. (3.6), but it must be complicated and difficult to use in practice. Note that $H(\lambda)$ is of order 2, we can directly use some numerical optimization method to reach the minimal point y_{\min} and the minimum $\kappa_2(T; y_{\min})$.

If $z_a^T y = \beta$, then $-JH(\lambda)J = \lambda^2 I_2 - \lambda(I_2 - H_1)$, which implies that $\lambda_1 = \lambda_2 = 0$ and $\kappa_2(T) = +\infty$. Hence the whole domain is split into two connected domains: $\{y : z_a^T y < \beta\}$ and $\{y : z_a^T y > \beta\}$. Then we may use optimization methods to try to attain the local minimum in both of the two domains.

Most optimization methods require the gradient and the Hessian of the objective function, and we calculate it theoretically.

First try to obtain those for the eigenvalues. For any eigenpair of Eq. (3.7), it holds that $-JH(\lambda)Jx = 0$, and without loss of generality, $x^T x = 1$. Let $W = \delta^2(W_e + W_y)$, $\xi = z_a^T y - \beta$, and then

$$-JH(\lambda)J = \lambda^2 I_2 - \lambda(I_2 + SWS + \xi^2 S^2) + \xi^2 S^2. \tag{4.1}$$

Taking differential on both sides gives

$$-Jd(H(\lambda))Jx - JH(\lambda)Jdx = 0, \tag{4.2}$$

and

$$0 = x^T JH(\lambda)Jdx = -x^T Jd(H(\lambda))Jx.$$

Thus,

$$0 = -x^T \left(2\lambda d\lambda I_2 - d\lambda [I_2 + SWS + \xi^2 S^2] - \lambda [SdWS + 2\xi d\xi S^2] + 2\xi d\xi S^2 \right) x.$$

Write $t := Sx = [\tau_1 \quad \tau_2]^T$, and noticing $x^T x = 1$,

$$0 = 2\lambda d\lambda - d\lambda [1 + t^T W t + \xi^2 t^T t] - \lambda [t^T dW t + 2\xi d\xi t^T t] + 2\xi d\xi t^T t,$$

and then

$$d\lambda = \frac{\lambda t^T dWt + 2(\lambda - 1)\zeta d\zeta^T t}{2\lambda - [1 + t^T Wt + \zeta^2 t^T t]}.$$

Note that $d\zeta = z_a^T dy$ and by Eq. (3.5),

$$\begin{aligned} dWt &= \delta^2 d \left((y^T B_W y) I_2 + \begin{bmatrix} 0 & (e_3 - e_1)^T B_W y \\ (e_3 - e_1)^T B_W y & -2e_2^T B_W y \end{bmatrix} \right) t \\ &= \delta^2 t d \left(y^T B_W y \right) + \delta^2 \begin{bmatrix} 0 & (e_3 - e_1)^T B_W dy \\ (e_3 - e_1)^T B_W dy & -2e_2^T B_W dy \end{bmatrix} t \\ &= 2\delta^2 t y^T B_W dy + \delta^2 \begin{bmatrix} \tau_2 (e_3 - e_1)^T \\ \tau_1 (e_3 - e_1)^T - 2\tau_2 e_2^T \end{bmatrix} B_W dy \\ &= \delta^2 (2ty^T + E_t^T) B_W dy, \end{aligned}$$

where

$$E_t := [\tau_2(e_3 - e_1) \quad \tau_1(e_3 - e_1) - 2\tau_2 e_2].$$

Then

$$t^T dWt = 2\delta^2 \left[(t^T t)y + \frac{1}{2} E_t t \right]^T B_W dy = 2\delta^2 [(t^T t)y + \tau_1 \tau_2 (e_3 - e_1) - \tau_2^2 e_2]^T B_W dy,$$

and we have

$$d\lambda = \frac{2\lambda\delta^2 [(t^T t)y + \tau_1 \tau_2 (e_3 - e_1) - \tau_2^2 e_2]^T B_W + 2(\lambda - 1)t^T t \zeta z_a^T}{2\lambda - [1 + t^T Wt + \zeta^2 t^T t]} dy,$$

which implies

$$\nabla\lambda = 2 \frac{\lambda\delta^2 B_W [(t^T t)y + \tau_1 \tau_2 (e_3 - e_1) - \tau_2^2 e_2] + (\lambda - 1)t^T t \zeta z_a}{2\lambda - [1 + t^T Wt + \zeta^2 t^T t]}.$$

Writing $\tau_W := t^T Wt$, then

$$\nabla\lambda = 2 \frac{\lambda\delta^2 B_W [(t^T t)y + \tau_1 \tau_2 (e_3 - e_1) - \tau_2^2 e_2] + (\lambda - 1)t^T t \zeta z_a}{2\lambda - [1 + \tau_W + \zeta^2 t^T t]},$$

and

$$d(\nabla\lambda) = \frac{\partial \nabla\lambda}{\partial \lambda} d\lambda + \frac{\partial \nabla\lambda}{\partial \zeta} d\zeta + \frac{\partial \nabla\lambda}{\partial \tau_1} d\tau_1 + \frac{\partial \nabla\lambda}{\partial \tau_2} d\tau_2 + \frac{\partial \nabla\lambda}{\partial \tau_W} d\tau_W. \tag{4.3}$$

Let

$$\sigma := 2\lambda - [1 + t^T Wt + \zeta^2 t^T t], \quad g := \lambda\delta^2 B_W \left[(t^T t)y + \frac{1}{2} E_t t \right] + (\lambda - 1)t^T t \zeta z_a, \tag{4.4}$$

then

$$\nabla\lambda = 2 \frac{g}{\sigma}. \tag{4.5}$$

Putting

$$\begin{aligned} \frac{\partial \nabla \lambda}{\partial \lambda} &= \frac{2}{\sigma^2} \left((\delta^2 B_W \left[(t^T t)y + \frac{1}{2} E_t t \right] + t^T t \zeta z_a) \sigma - 2g \right), \\ \frac{\partial \nabla \lambda}{\partial \xi} &= \frac{2}{\sigma^2} \left((\lambda - 1) t^T t z_a \sigma + 2 t^T t \zeta g \right), \\ \frac{\partial \nabla \lambda}{\partial \tau_1} &= \frac{2}{\sigma^2} \left((\lambda \delta^2 B_W [2\tau_1 y + \tau_2 (e_3 - e_1)] + 2(\lambda - 1) \tau_1 \zeta z_a) \sigma + \zeta^2 2\tau_1 g \right), \\ \frac{\partial \nabla \lambda}{\partial \tau_2} &= \frac{2}{\sigma^2} \left((\lambda \delta^2 B_W [2\tau_2 y + \tau_1 (e_3 - e_1) - 2\tau_2 e_2] + 2(\lambda - 1) \tau_2 \zeta z_a) \sigma + \zeta^2 2\tau_2 g \right), \\ \frac{\partial \nabla \lambda}{\partial \tau_W} &= \frac{2}{\sigma^2} g, \\ d\tau_W &= 2t^T W dt + t^T dWt, \end{aligned}$$

into Eq. (4.3) gives

$$\begin{aligned} d(\nabla \lambda) &= \frac{2}{\sigma^2} \left(\left((\delta^2 B_W \left[(t^T t)y + \frac{1}{2} E_t t \right] + t^T t \zeta z_a) \sigma - 2g \right) 2 \frac{g^T}{\sigma} dy \right. \\ &\quad + \left((\lambda - 1) t^T t z_a \sigma + 2 t^T t \zeta g \right) z_a^T dy \\ &\quad + 2\sigma \lambda \delta^2 B_W y t^T dt + \sigma \lambda \delta^2 B_W E_t dt + 2\sigma(\lambda - 1) \zeta z_a t^T dt + 2\zeta^2 g t^T dt \\ &\quad \left. + g \left(2t^T W dt + 2\delta^2 \left[(t^T t)y + \frac{1}{2} E_t t \right]^T B_W dy \right) \right) \\ &=: \frac{2}{\sigma^2} \left(G_y dy + G_t dt \right), \end{aligned} \tag{4.6}$$

where

$$\begin{aligned} G_y &= \left((\delta^2 B_W \left[(t^T t)y + \frac{1}{2} E_t t \right] + t^T t \zeta z_a) \sigma - 2g \right) 2 \frac{g^T}{\sigma} \\ &\quad + \left((\lambda - 1) t^T t z_a \sigma + 2 t^T t \zeta g \right) z_a^T + 2g \delta^2 \left[(t^T t)y + \frac{1}{2} E_t t \right]^T B_W \\ &\quad \text{by } \delta^2 B_W \left[(t^T t)y + \frac{1}{2} E_t t \right] = \frac{1}{\lambda} [g - (\lambda - 1) t^T t \zeta z_a] \\ &= \left(\frac{4}{\lambda} - \frac{4}{\sigma} \right) g g^T + \frac{2}{\lambda} t^T t \zeta (z_a g^T + g z_a^T) + \sigma(\lambda - 1) t^T t z_a z_a^T \\ &\quad \text{by } -x^T JH(\lambda) Jx = \lambda^2 - \lambda(2\lambda - \sigma) + \zeta^2 t^T t = 0 \Rightarrow \sigma - \lambda = -\frac{\zeta^2 t^T t}{\lambda} \end{aligned}$$

$$\begin{aligned}
 &= \frac{4}{\sigma\lambda} \frac{\zeta^2 t^T t}{-\lambda} g g^T + \frac{2}{\lambda} t^T t \zeta (z_a g^T + g z_a^T) + \sigma(\lambda - 1) t^T t z_a z_a^T \\
 &= \sigma t^T t \left(- \left(z_a - \frac{2\zeta}{\sigma\lambda} g \right) \left(z_a - \frac{2\zeta}{\sigma\lambda} g \right)^T + \lambda z_a z_a^T \right),
 \end{aligned}$$

and

$$G_t = 2gt^T W + 2(\sigma\lambda\delta^2 B_W y + \sigma(\lambda - 1)\zeta z_a + \zeta^2 g)t^T + \sigma\lambda\delta^2 B_W E_t. \tag{4.7}$$

Then we calculate dt . By Eq. (4.2), dx satisfies $-JH(\lambda)Jdx = Jd(H(\lambda))Jx$, which gives

$$\begin{aligned}
 -JH(\lambda)Jdx &= Jd(H(\lambda))Jx \\
 &= -2\lambda d\lambda x + d\lambda [I_2 + SWS + \zeta^2 S^2]x + \lambda SdWSx + 2\lambda\zeta d\zeta S^2 x - 2\zeta d\zeta S^2 x \\
 &= (1 - 2\lambda)d\lambda x + d\lambda [SWS + \zeta^2 S^2]x + \lambda SdWSx + 2(\lambda - 1)\zeta d\zeta S^2 x.
 \end{aligned}$$

Note that $x^T Jx = 0$, $x^T dx = 0$ and $x \in \mathbb{R}^2$. Hence $dx = \eta Jx$ holds for some η . Moreover, since x is an eigenvector of the matrix $-JH(\lambda)J$ corresponding to 0, Jx has to be its eigenvector corresponding to a positive eigenvalue ω , and

$$\begin{aligned}
 \omega &= (Jx)^T (-JH(\lambda)J)(Jx) = (Jx)^T [\lambda^2 I_2 - \lambda(I_2 + SWS + \zeta^2 S^2) + \zeta^2 S^2](Jx) \\
 &= \lambda^2 - \lambda(1 + s^T Ws + \zeta^2 s^T s) + \zeta^2 s^T s,
 \end{aligned}$$

where $s := SJx = SJS^{-1}t$. Thus

$$(Jx)^T (-JH(\lambda)J)(\eta Jx) = -d\lambda x^T J[SWS + \zeta^2 S^2]x - \lambda x^T JSdWSx - 2(\lambda - 1)\zeta d\zeta x^T JS^2 x,$$

and

$$\begin{aligned}
 \eta &= \frac{-d\lambda x^T J[SW + \zeta^2 S]t - \lambda x^T JSdWt - 2(\lambda - 1)\zeta d\zeta x^T JS t}{\omega} \\
 &= \frac{-2\frac{g^T dy}{\sigma} x^T J[SW + \zeta^2 S]t - \lambda x^T JS(2\delta^2 ty^T B_W dy + \delta^2 E_t^T B_W dy) - 2(\lambda - 1)\zeta z_a^T dy x^T JS t}{\omega} \\
 &= \frac{-2\frac{x^T J[SW + \zeta^2 S]t}{\sigma} g^T - 2\lambda\delta^2 (x^T JS t) y^T B_W - \lambda\delta^2 x^T JSE_t^T B_W - 2(\lambda - 1)\zeta (x^T JS t) z_a^T}{\omega} dy \\
 &= \frac{2\frac{s^T Wt + \zeta^2 s^T t}{\sigma} g^T + 2\lambda\delta^2 s^T ty^T B_W + \lambda\delta^2 s^T E_t^T B_W + 2(\lambda - 1)\zeta s^T t z_a^T}{\omega} dy.
 \end{aligned}$$

Since by Eq. (4.7),

$$\begin{aligned}
 g_t &:= G_t s = 2gt^T Ws + 2(\sigma\lambda\delta^2 B_W y + \sigma(\lambda - 1)\zeta z_a + \zeta^2 g)t^T s + \sigma\lambda\delta^2 B_W E_t s \\
 &= 2(t^T Ws + \zeta^2 t^T s)g + \sigma[2\lambda\delta^2 t^T s B_W y + 2(\lambda - 1)t^T s \zeta z_a + \lambda\delta^2 B_W E_t s],
 \end{aligned}$$

we have

$$\eta = \frac{g_t^T}{\sigma\omega} dy. \tag{4.8}$$

Moreover,

$$\begin{aligned}
 g_t &= 2(t^T Ws + \zeta^2 t^T s)g + \sigma[2\lambda\delta^2 t^T s B_W y + 2(\lambda - 1)t^T s \zeta z_a + \lambda\delta^2 B_W E_t s] \quad \text{by Eq. (4.4)} \\
 &= 2(t^T Ws + \zeta^2 t^T s)g + \sigma\left[2\frac{t^T s}{t^T t}\left(g - \lambda\delta^2 B_W \frac{1}{2} E_t t\right) + \lambda\delta^2 B_W E_t s\right] \\
 &= 2\left(t^T Ws + \zeta^2 t^T s + \frac{\sigma t^T s}{t^T t}\right)g + \frac{\sigma}{t^T t}\lambda\delta^2 B_W E_t [(t^T t)s - (t^T s)t]. \tag{4.9}
 \end{aligned}$$

Note that

$$(t^T t)s - (t^T s)t = [(t^T t)I_2 - tt^T]s = -Jtt^T J s = -Jt(x^T SJSJx) = -Jtx^T \left(-\frac{4}{\delta}I_2\right)x = \frac{4}{\delta}Jt,$$

and

$$\begin{aligned}
 0 &= x^T(-J^T H(\lambda)J)Jx = x^T(\lambda^2 I_2 - \lambda[I_2 + SWS + \zeta^2 S^2] + \zeta^2 S^2)Jx \\
 &= -\lambda[t^T Ws + \zeta^2 t^T s] + \zeta^2 t^T s,
 \end{aligned}$$

which gives

$$t^T Ws = \frac{(1 - \lambda)\zeta^2 t^T s}{\lambda}. \tag{4.10}$$

Note that by Eq. (4.4)

$$\zeta^2 t^T t + \sigma\lambda = 2\lambda^2 - \lambda[1 + t^T Wt + \zeta^2 t^T t] + \zeta^2 t^T t = \lambda^2. \tag{4.11}$$

Substituting Eqs. (4.10) and (4.11) into Eq. (4.9) gives

$$\begin{aligned}
 g_t = G_t s &= 2\left(\frac{\zeta^2}{\lambda}t^T s + \frac{\sigma t^T s}{t^T t}\right)g + \frac{4\sigma}{t^T t}\lambda\delta B_W E_t Jt \\
 &= 2t^T s\left(\frac{\zeta^2}{\lambda} + \frac{\sigma}{t^T t}\right)g + \frac{4\sigma}{t^T t}\lambda\delta B_W E_t Jt \\
 &= 2\frac{\lambda}{t^T t}[t^T s g + 2\sigma\delta B_W E_t Jt].
 \end{aligned}$$

Noticing by Eq. (4.8)

$$dt = Sdx = \eta S J x = \eta s = \left(\frac{g_t^T}{\sigma\omega} dy\right) s,$$

Eq. (4.6) becomes

$$d(\nabla\lambda) = \frac{2}{\sigma^2}\left(\sigma t^T t\left(\lambda z_a z_a^T - \left(z_a - \frac{2\zeta}{\sigma\lambda}g\right)\left(z_a - \frac{2\zeta}{\sigma\lambda}g\right)^T\right) dy + \frac{g_t g_t^T}{\sigma\omega} dy\right),$$

and

$$\nabla^2\lambda = \frac{2}{\sigma^2}\left(\sigma t^T t\left(\lambda z_a z_a^T - \left(z_a - \frac{2\zeta}{\sigma\lambda}g\right)\left(z_a - \frac{2\zeta}{\sigma\lambda}g\right)^T\right) + \frac{g_t g_t^T}{\sigma\omega}\right). \tag{4.12}$$

Following this, we may easily have $\nabla\lambda_4, \nabla\lambda_1$ and then

$$\nabla\kappa = \frac{\lambda_1\nabla\lambda_4 - \lambda_4\nabla\lambda_1}{2\lambda_1^2\kappa} = \frac{\kappa}{2} \left(\frac{\nabla\lambda_4}{\lambda_4} - \frac{\nabla\lambda_1}{\lambda_1} \right), \tag{4.13}$$

with $\kappa = \kappa_2(T)$, which can be used in the optimization methods. Moreover,

$$\begin{aligned} \nabla^2\kappa &= \left(\frac{\nabla\lambda_4}{\lambda_4} - \frac{\nabla\lambda_1}{\lambda_1} \right) \frac{\nabla\kappa^T}{2} + \frac{\kappa}{2} \left(\frac{\nabla^2\lambda_4}{\lambda_4} - \frac{\nabla\lambda_4\nabla\lambda_4^T}{\lambda_4^2} - \frac{\nabla^2\lambda_1}{\lambda_1} + \frac{\nabla\lambda_1\nabla\lambda_1^T}{\lambda_1^2} \right) \\ &= \frac{\kappa}{2} \left(\frac{2\nabla\kappa\nabla\kappa^T}{\kappa^2} + \frac{\nabla^2\lambda_4}{\lambda_4} - \frac{\nabla\lambda_4\nabla\lambda_4^T}{\lambda_4^2} - \frac{\nabla^2\lambda_1}{\lambda_1} + \frac{\nabla\lambda_1\nabla\lambda_1^T}{\lambda_1^2} \right) \\ &= \frac{\kappa}{2} \left(\frac{\nabla^2\lambda_4}{\lambda_4} - \frac{\nabla^2\lambda_1}{\lambda_1} - \frac{2\nabla\lambda_1\nabla\kappa^T}{\lambda_1\kappa} - \frac{2\nabla\kappa\nabla\lambda_1^T}{\kappa\lambda_1} - \frac{2\nabla\kappa\nabla\kappa^T}{\kappa^2} \right), \end{aligned} \tag{4.14}$$

of which the last equality holds for $\frac{\nabla\lambda_4}{\lambda_4} = 2\frac{\nabla\kappa}{\kappa} + \frac{\nabla\lambda_1}{\lambda_1}$ by Eq. (4.13).

5 Solve the problem

Usually optimization methods converge to the local minimal point or stationary point, which would not be the global minimal point. It is well known that optimization methods performing on convex functions always converge to the unique global minimal point. Hence we would like to know whether the objective function is convex in its connected domain. However, numerical tests do not support that.

Example 5.1 (Convex Test). All the numerical tests in the paper are implemented in MATLAB R2017a. We randomly generate three symmetric matrices M, C, K (not necessarily positive definite), where $M = U^T \Lambda U$ in which U is generated by the MATLAB code `rand(10)`, and Λ is diagonal and each diagonal entry follows the uniform distribution on $[-0.5, 0.5]$. Then we pick its two eigenpairs of different types, and calculate the related unconstrained optimization problem equation (3.6).

Then we randomly generate 1000 points where any coordinate follows the uniform distribution $[-100, 100]$. For each point we compute its Hessian matrix $\nabla^2\kappa$, and then compute its minimal eigenvalue λ_{\min} . Here the gradient $\nabla\kappa$ and the Hessian $\nabla^2\kappa$ are generated according to Eq. (4.13) and Eq. (4.14) respectively. Specially, if the point (nearly) locates in the hyperplane $z_a^T y = \beta$, then we simply discard the pair, or equivalently let the minimal eigenvalue be $+\infty$. Then we find the minimal λ_{\min} of the 1000 pairs.

The procedure above will be repeated 1000 times. The distribution of the 1000 minimal λ_{\min} is shown below.

$(-10^8, -10^6)$	$(-10^6, -10^4)$	$(-10^4, -10^3)$	$(-10^3, -10^2)$	$(-10^2, -10^0)$
27	294	351	255	73

The numerical test suggests us that the objective function is hardly convex.

Although the objective function may be not convex, the global minimal point can still be converged if it is pseudoconvex. A function $f: S \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is called pseudoconvex, if $f(y_1) < f(y_2) \Rightarrow (y_1 - y_2)^T \nabla f(y_2) < 0$ for any $y_1, y_2 \in S$. A function $f(y)$ pseudoconvex on S has this property: if $y_0 \in S$ is a stationary point, then y_0 is a global minimal point; if y_0 is a local minimal point, then y_0 is a global minimal point. If $\nabla^2 f(y)$ is positive definite on the subspace $\{\nabla f(y)\}^\perp$, then f is pseudoconvex. In detail, we will judge whether $Z^H \nabla^2 \kappa Z \succ 0$, where Z is an orthonormal basis matrix of the orthogonal complement of $\nabla \kappa$, and $Z^T \nabla \kappa = 0$. Specially $Z = I_3$ if $\nabla \kappa = 0$. More information can be found in classic textbooks (e.g. [3]), or this early paper [4].

Numerical tests imply the objective function $\kappa_2(T)$ with respect to y is pseudoconvex in its connected domain for a few instances.

Example 5.2 (Pseudoconvex Test). The whole process is nearly the same as in Example 5.1, but rather than the minimal eigenvalue of $\nabla^2 \kappa$, we calculate the minimal eigenvalue $\tilde{\lambda}_{\min}$ of $Z^H \nabla^2 \kappa Z$.

The distribution of the 1000 minimal $\tilde{\lambda}_{\min}$ is shown below.

$(-10^6, -10^2)$	$(-10^2, -10^{-2})$	$(-10^{-2}, -10^{-6})$	$(-10^{-6}, -10^{-11})$	$(10^{-13}, 10^{-8})$
97	595	275	26	7

The numerical test suggests us that the objective function is pseudoconvex for a few instances only.

Although in most instances the objective function is likely not to be pseudoconvex, for any given problem, namely given the quadratic matrix polynomial $Q(\lambda)$ and its two eigenpairs to be deflated, we are able to transform it into an equivalent problem whose associated tiny-scale unconstrained optimization problem with a pseudoconvex objective function.

Recall the problem setting in the beginning of Section 2. Note that a is determined by the eigenvectors only. Thus, performing an affine transformation on the eigenvalue does not change either a or the eigenvectors to be deflated. In details, the problem

$$\tilde{Q}(\lambda) = \frac{1}{v_1^2} Q(v_1 \lambda + v_2) = \lambda^2 M + \lambda \frac{2v_2 M + C}{v_1} + \frac{v_2^2 M + v_2 C + K}{v_1^2} =: \lambda^2 M + \lambda \tilde{C} + \tilde{K} \quad (5.1)$$

will be treated instead. However,

$$\tilde{\alpha}_C = a^T \tilde{C} a = \frac{2v_2 \alpha_M + \alpha_C}{v_1}, \quad \tilde{\alpha}_K = a^T \tilde{K} a = \frac{v_2^2 \alpha_M + v_2 \alpha_C + \alpha_K}{v_1^2}.$$

It is not difficult to make $\tilde{\alpha}_C = 0, |\alpha_M| = |\tilde{\alpha}_K|$. In fact,

$$v_2 = -\frac{\alpha_C}{2\alpha_M}, \quad v_1 = \frac{\sqrt{|4\alpha_M \alpha_K - \alpha_C^2|}}{2\alpha_M} \neq 0.$$

Moreover, write $\chi = \text{sign}(4\alpha_M\alpha_K - \alpha_C^2)$, and then $\tilde{\alpha}_K = \chi\alpha_M$.

Since the deflation works on eigenvectors only, the tiny-scale unconstrained optimization problem must have the same form, except the terms involving α_C, α_K . Hence in the following, we will directly use α_C, α_K rather than $\tilde{\alpha}_C, \tilde{\alpha}_K$, and all the notation will be kept. Now

$$\delta = \frac{1}{\alpha_M^2 \chi}, \quad S = \alpha_M \begin{bmatrix} 1 & \\ & \chi \end{bmatrix}.$$

So $|\alpha_M| = \frac{1}{\sqrt{|\delta|}}$. The involving hyperbolic quadratic polynomial, namely Eq. (4.1), is

$$-JH(\lambda)J = \lambda^2 I_2 - \lambda [I_2 + SWS + \zeta^2 S^2] + \zeta^2 S^2 = \lambda^2 I_2 - \lambda \left[\left(1 + \frac{\zeta^2}{\chi\delta}\right) I_2 + SWS \right] + \frac{\zeta^2}{\chi\delta} I_2.$$

Since I_2 and SWS can be simultaneously diagonalizable by an orthonormal matrix, that hyperbolic quadratic polynomial is diagonalizable by an orthonormal matrix, which implies: 1) an eigenvector corresponding to a positive-type eigenvalue has to be that corresponding to a negative-type eigenvalue; 2) the eigenvectors corresponding to two positive-type eigenvalues are orthogonal. There are two cases to be considered.

Case 1: λ_1, λ_4 cannot share a same eigenvector, which implies that λ_1, λ_3 share a same eigenvector x_1 and λ_2, λ_4 share a same eigenvector x_4 satisfying $x_4^T x_1 = 0$.

Since $x_1^T x_1 = x_4^T x_4 = 1$, without loss of generality, assume $x_1 = Jx_4$. Then $t_4^T t_4 = x_4^T S^2 x_4 = \frac{1}{\chi\delta}$, and similarly $t_1^T t_1 = \frac{1}{\chi\delta}$. Thus, $\lambda_1 \lambda_3 = \frac{\zeta^2}{\chi\delta} = \lambda_2 \lambda_4$ and then $\frac{\lambda_1}{\lambda_2} = \frac{\lambda_4}{\lambda_3}$, which implies $\lambda_1 = \lambda_2, \lambda_3 = \lambda_4$. This tells λ_1, λ_4 share a same eigenvector vector, a contradiction.

Case 2: λ_1, λ_4 share a same eigenvector. Since $x_1^T x_1 = x_4^T x_4 = 1$, without loss of generality, assume $x_1 = x_4$, then $t_1 = t_4, s_1 = s_4$. Then $t_4^T s_4 = x_4^T S^2 Jx_4 = 0, t_4^T t_4 = x_4^T S^2 x_4 = \frac{1}{\chi\delta}$. Also, $\sigma_4 = -\sigma_1 = \lambda_4 - \lambda_1, g_{t_4} = g_{t_1} = 4\chi\delta^2 \sigma_4 \lambda_4 B_W E_{t_4} J t_4$ by Eq. (4.9).

By Eqs. (4.5) and (4.13),

$$\begin{aligned} \frac{\nabla \kappa}{\kappa} &= \frac{1}{2} \left(\frac{\nabla \lambda_4}{\lambda_4} - \frac{\nabla \lambda_1}{\lambda_1} \right) = \frac{g_4}{\sigma_4 \lambda_4} - \frac{g_1}{\sigma_1 \lambda_1} \\ &= \delta^2 B_W \left[\frac{t_4^T t_4}{\sigma_4} y + \frac{1}{2\sigma_4} E_{t_4} t_4 \right] + \left(1 - \frac{1}{\lambda_4}\right) \frac{t_4^T t_4}{\sigma_4} \zeta z_a \\ &\quad - \delta^2 B_W \left[\frac{t_1^T t_1}{\sigma_1} y + \frac{1}{2\sigma_1} E_{t_1} t_1 \right] - \left(1 - \frac{1}{\lambda_1}\right) \frac{t_1^T t_1}{\sigma_1} \zeta z_a \\ &= \delta^2 B_W \left[\left(\frac{1}{\chi\delta\sigma_4} - \frac{1}{\chi\delta\sigma_1} \right) y + \frac{1}{2\sigma_4} E_{t_4} t_4 - \frac{1}{2\sigma_1} E_{t_1} t_1 \right] \\ &\quad + \left[\frac{1}{\chi\delta\sigma_4} - \frac{1}{\chi\delta\sigma_1} - \frac{1}{\chi\delta\lambda_4\sigma_4} + \frac{1}{\chi\delta\lambda_1\sigma_1} \right] \zeta z_a \\ &= \frac{1}{\chi\delta\sigma_4} \left(\delta^2 B_W [2y + \chi\delta E_{t_4} t_4] + \left[2 - \left(\frac{1}{\lambda_4} + \frac{1}{\lambda_1} \right) \right] \zeta z_a \right). \end{aligned}$$

By $Z^T \frac{\nabla \kappa}{\kappa} = 0$,

$$\delta^2 Z^T B_W [2y + \chi \delta E_{t_4} t_4] = \left[\frac{1}{\lambda_4} + \frac{1}{\lambda_1} - 2 \right] \xi Z^T z_a,$$

and then

$$\begin{aligned} \frac{2\xi}{\sigma_4 \lambda_4} Z^T g_4 &= \frac{2\xi}{\sigma_4 \lambda_4} \left(\lambda_4 \delta^2 Z^T B_W \left[\frac{1}{\chi \delta} y + \frac{1}{2} E_{t_4} t_4 \right] + \frac{\lambda_4 - 1}{\chi \delta} \xi Z^T z_a \right) \\ &= \frac{2\xi}{\sigma_4 \lambda_4} \left(\frac{\lambda_4}{2\chi \delta} \left[\frac{1}{\lambda_4} + \frac{1}{\lambda_1} - 2 \right] \xi Z^T z_a + \frac{\lambda_4 - 1}{\chi \delta} \xi Z^T z_a \right) \\ &= \frac{\xi^2}{\chi \delta \sigma_4 \lambda_4} \left(\frac{\lambda_4}{\lambda_1} - 1 \right) Z^T z_a \\ &= \frac{\lambda_1 \lambda_4}{(\lambda_4 - \lambda_1) \lambda_4} \left(\frac{\lambda_4}{\lambda_1} - 1 \right) Z^T z_a \quad \text{by } \frac{\xi^2}{\chi \delta} = \lambda_1 \lambda_4 \\ &= Z^T z_a. \end{aligned}$$

Similarly, $\frac{2\xi}{\sigma_1 \lambda_1} Z^T g_1 = Z^T z_a$. Thus, by Eqs. (4.12) and (4.14),

$$Z^T \nabla^2 \kappa Z = \frac{\kappa}{2} Z^T \left(\frac{\nabla^2 \lambda_4}{\lambda_4} - \frac{\nabla^2 \lambda_1}{\lambda_1} \right) Z = \kappa Z^T \left(\frac{2}{\chi \delta \sigma_4} z_a z_a^T + \frac{1}{\sigma_4^3} \left[\frac{1}{\lambda_4 \omega_4} + \frac{1}{\lambda_1 \omega_1} \right] g_{t_4} g_{t_4}^T \right) Z \succeq 0.$$

Moreover, in the generic case, except several critical points

$$\text{rank}(Z^T \nabla^2 \kappa Z) = \text{rank}(Z^T [z_a \quad g_{t_4}]) = 2,$$

which implies $Z^T \nabla^2 \kappa Z \succ 0$.

To sum up, the objective function of this transformed problem equation (5.1) is pseudoconvex. Then, optimization methods would effectively attain the global minimum in either of the two connected domains. When using optimization method to solve the problem, we may start at one initial y with $z_a^T y > \beta$ and another initial y with $z_a^T y < \beta$. As long as the problem is pseudoconvex, the two numerical solutions come out fast. After that, we may pick the better one in the two. Please note that this transforming technique is implemented for the pseudoconvexity, so it would produce a structure-preserving transformation for the deflation whose condition is worse than that of the untransformed one.

Example 5.3 (Convergence Test). The unconstrained optimization problem equation (3.6) is built as in Example 5.1.

Then we randomly generate 100 points where any coordinate follows the uniform distribution $[-100, 100]$. Then we use each point as the initial point and perform both the quasi-Newton method and the trust-region method to observe the convergence behavior. All the options of the methods are set in default by MATLAB. Among those, the tolerance of the difference of the function value is 10^{-6} , and the iteration will stop if the number of iterations reaches 300.

The procedure above will be repeated 1000 times. Two rounds of the procedure are made for the untransformed problem and the transformed problem. The data of the numerical tests are shown below.

problem	method	domain	min attained	iterations			func. evaluations		
				max	min	mean	max	min	mean
untransformed	trust-region	$z_a^T y > \beta$	502	300	10	28.2	301	11	29.2
		$z_a^T y < \beta$	498	300	11	28.5	301	12	29.5
	quasi-Newton	$z_a^T y > \beta$	489	68	1	13.1	280	14	43.3
		$z_a^T y < \beta$	511	70	1	13.0	187	13	43.2
transformed	trust-region	$z_a^T y > \beta$	468	300	4	28.0	301	5	29.5
		$z_a^T y < \beta$	532	300	5	28.4	301	6	29.4
	quasi-Newton	$z_a^T y > \beta$	496	65	1	13.4	300	8	39.2
		$z_a^T y < \beta$	504	70	1	13.5	300	10	39.1

From the data, we know that generally the quasi-Newton method uses fewer iterations but more function evaluations in the sense of the average case and the worst case. However, the trust-region method needs the (numerical) Hessian. Since the number of function evaluations needed by the quasi-Newton method does not exceed twice of that need by the trust-region method, the quasi-Newton method is faster than the trust-region methods according to the data.

It seems that in many cases the minimums are attained in different connected domain by the two methods. Actually in fact, the minimums are nearly the same, which can be illustrated by the distribution of the relative difference of the 1000 pairs of minimums obtained by trust-region and quasi-Newton, namely $\frac{\kappa_{\text{trust-region}} - \kappa_{\text{quasi-Newton}}}{\kappa_{\text{trust-region}}}$ given below.

problem	$(-2, -1)$	$(-1, -10^{-2})$	$(-10^{-2}, 0)$	$(0, 10^{-2})$	$(10^{-2}, 1)$
untransformed	0	153	44	115	688
transformed	1	33	39	365	562

From the data, we can see the minimums found by the trust-region method is between half and twice those found by the quasi-Newton method, which would produce little affect on the condition number; in most cases, the quasi-Newton method is a little better than the trust-region.

On the other hand, from the data above, we know performing on the transformed problem costs slightly less than performing on the untransformed problem. However, if we compare the minimums found for the untransformed and transformed problems, namely $\frac{\kappa_{\text{untransformed}}}{\kappa_{\text{transformed}}}$ give below,

$(10^{-5}, 0.1)$	$(0.1, 1)$	$(1, 10)$	$(10, 100)$	$(100, 1000)$	$(1000, 10^5)$
30	123	299	396	135	17

we may discover that in most cases performing on the transformed problems is likely to find out a much better structure-preserving transformation than performing on the untransformed problems.

To sum up, a good choice is to use the quasi-Newton method to solve the unconstrained optimization problem Eq. (3.6) induced by the transformed quadratic matrix polynomial.

Besides, from the data, we can see the minimum has no preference on the connected domain, which makes us have to compute on them both.

6 Concluding remarks

We have built a tiny-scale unconstrained optimization problem and suggested using optimization methods to solve it, which make it successful to find an optimal deflation for symmetric quadratic matrix polynomials. This would be a good answer to the question asked by Tisseur et al. [6]. However, it is quite natural to ask whether this method can be used for asymmetric problems. It is likely that some analogues would hold, but to ensure this is left for future work. In addition, investigating the theoretical solution to the unconstrained optimization problem would be also valuable to consider in future.

Acknowledgments

Part of the work was done when the author worked in Max Planck Institute for Dynamics of Complex Technical Systems. The author would like to thank Prof. Peter Benner and Prof. Sara Grundel for the discussion, and thank the anonymous referees for the comments that spur the author to survey further.

References

- [1] X. Liang and R.-C. Li. The hyperbolic quadratic eigenvalue problem. *Forum of Mathematics, Sigma*, 3(e13), 2015.
- [2] K. Veselić. Note on interlacing for hyperbolic quadratic pencils. In J. Behrndt, K.-H. Förster, and C. Trunk (Eds.), *Recent Advances in Operator Theory in Hilbert and Krein Spaces*, 198:305-307, 2010.
- [3] S. Boyd and L. Vandenberghe. *Convex Optimization*, Cambridge University Press, 2004.
- [4] O. L. Mangasarian. Pseudo-convex functions. *SIAM J. Control*, 3(2):281-290, 1965.
- [5] B. Meini. A "shift-and-deflate" technique for quadratic matrix polynomials. *Linear Algebra Appl.*, 438:1946-1961, 2013.
- [6] F. Tisseur, S. D. Garvey, and C. Munro. Deflating quadratic matrix polynomials with structure preserving transformations. *Linear Algebra Appl.*, 435:464-479, 2011.