

**A Krylov Subspace Method for Large  
Scale Second Order Cone Linear  
Complementarity Problem**

**Lei-Hong Zhang  
Wei Hong Yang  
Chungen Shen  
Ren-Cang Li**

**Technical Report 2014-19**

# A Krylov Subspace Method for Large Scale Second Order Cone Linear Complementarity Problem

Lei-Hong Zhang\*    Wei Hong Yang<sup>†</sup>    Chungen Shen<sup>‡</sup>    Ren-Cang Li<sup>§</sup>

November 8, 2014

## Abstract

Recently, a bisection-Newton (BN) iteration [Zhang and Yang, *Math. Comp.*, 83:1701–1726 (2013)] was proposed for solving the second order cone linear complementarity problem (SOCLCP) and was shown to be very efficient for small-to-medium size problems. However, for large scale problems, difficulty arises in the BN iteration in two aspects: 1) a specific eigenpair needs to be computed accurately, and 2) there are many large scale linear systems each with different coefficient matrices to solve. In this paper, we present an efficient method based on Krylov subspace approximation. Although there are still many large scale linear systems arise in the new method, most of them share the same coefficient matrices, a fact we will fully exploit to significantly cut down the number of, often the most expensive part in the computation, LU decompositions needed in solving the linear systems. Here, we first show that SOCLCP can be solved by finding a positive zero  $s_* \in \mathbb{R}$  of a particular rational function  $h(s)$ , and then propose a Krylov subspace method to reduce  $h(s)$  to  $h_\ell(s)$  as in the model reduction. The zero  $s_*$  of  $h(s)$  can be accurately approximated by that of  $h_\ell(s) = 0$  which itself can be casted as a small eigenvalue problem. The new method is made possible by our complete description of the curve  $h(s)$ , and it does not need the eigen-information of a large matrix previously required in the BN iteration and is suitable for large scale problems. The new method is tested and compared against the BN

---

\*School of Mathematics, Shanghai University of Finance and Economics, 777 Guoding Road, Shanghai 200433, People's Republic of China. Email: [longzlh@163.com](mailto:longzlh@163.com). The work of this author was supported in part by the National Natural Science Foundation of China NSFC-11101257, NSFC-11371102, and the Basic Academic Discipline Program, the 11th five year plan of 211 Project for Shanghai University of Finance and Economics.

<sup>†</sup>School of Mathematical Sciences, Fudan University, Shanghai, 200433, P. R. China. Email: [whyang@fudan.edu.cn](mailto:whyang@fudan.edu.cn). The work of this author was supported by the National Natural Science Foundation of China NSFC-11371102.

<sup>‡</sup>Department of Applied Mathematics, Shanghai Finance University, Shanghai 201209, China. Email: [shenchungen@gmail.com](mailto:shenchungen@gmail.com). The work of this author was supported in part by the National Natural Science Foundation of China (NSFC-11101281 and NSFC-11271259) and Innovation Program of Shanghai Municipal Education Commission (No.12YZ172).

<sup>§</sup>Department of Mathematics, University of Texas at Arlington, P.O. Box 19408, Arlington, TX 76019-0408, USA. Email: [rccli@uta.edu](mailto:rccli@uta.edu). This work was supported in part by NSF grants DMS-1115834 and DMS-1317330, by a Research Gift Grant from Intel Corporation, and by NSFC grant 11428104.

iteration and two other state-of-the-art packages: SDPT3 and SeDuMi. Our numerical results show that the method is very efficient both for small-to-medium dense problems as well as for large scale problems.

**Key words.** Second-order cone, linear complementarity problem, SOCLCP, globally uniquely solvable property, GUS, Krylov subspace, Model reduction, Linear complementarity problem via Arnoldi process (LCPvA)

**AMS subject classifications.** 90C33, 65K05, 65F99, 65F15, 65F30, 65P99

## 1 Introduction

For a given matrix-vector pair  $(M, \mathbf{q}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^n$ , in this paper, we will consider computational methods for the *second-order cone linear complementarity problem* (SOCLCP):

$$\mathbf{x} \in \mathbb{K}^n, \quad \mathbf{q} + M\mathbf{x} \in \mathbb{K}^n, \quad \mathbf{x}^\top(\mathbf{q} + M\mathbf{x}) = 0, \quad (1.1)$$

where  $\mathbb{K}^n$  is the *second-order cone* (SOC), also known as the *Lorentz cone* or the *ice-cream cone*, defined by

$$\mathbb{K}^n := \left\{ [x_1, \mathbf{x}_2^\top]^\top \in \mathbb{R} \times \mathbb{R}^{n-1} : \|\mathbf{x}_2\|_2 \leq x_1 \right\}.$$

To simplify our presentation, we will denote (1.1) by  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  whose set of solutions will be denoted by  $\text{SOL}(M, \mathbb{K}^n, \mathbf{q})$ .

The SOLCLP (1.1) is a special case of the linear complementarity problem over a Cartesian product of *multiple second-order cones*  $\mathbb{K}_{\times m} := \mathbb{K}^{n_1} \times \mathbb{K}^{n_2} \times \cdots \times \mathbb{K}^{n_m}$ :

$$\mathbf{x} \in \mathbb{K}_{\times m}, \quad \mathbf{q} + M\mathbf{x} \in \mathbb{K}_{\times m}, \quad \mathbf{x}^\top(\mathbf{q} + M\mathbf{x}) = 0. \quad (1.2)$$

We refer to, e.g., [1, 6, 8, 13, 17, 18, 20, 25, 33] for various applications as well as for the theoretical and numerical studies. As it is a more general problem, many theoretical properties that we will discuss, develop, and exploit for SOLCLP (1.1) are no longer admitted by (1.2). These properties are crucial for us to design our efficient algorithms for (1.1), however. It is hoped that our development here for (1.1) will help us to efficiently deal with this more general (1.2) in the future.

Both SOCLCPs (1.1) and (1.2) are similar to the classical linear complementarity problem over the cone  $\mathbb{R}_+^n = \{\mathbf{x} = [x_i] \in \mathbb{R}^n : x_i \geq 0, i = 1, 2, \dots, n\}$ :

$$\mathbf{x} \in \mathbb{R}_+^n, \quad \mathbf{q} + M\mathbf{x} \in \mathbb{R}_+^n, \quad \mathbf{x}^\top(\mathbf{q} + M\mathbf{x}) = 0,$$

which arises from many areas and has a wealth of development, in both theory and implementation (see e.g., [4, 7, 10, 21, 23, 29]).

An efficient approach for SOCLCP (1.1) is not only of value in its own right, but also may shed lights on solving the more general linear complementarity problem (1.2). A prime example is the matrix splitting method [44] for (1.2), whose efficiency largely comes from the efficiency of the bisection-Newton (BN) iteration [43] for (1.1). The BN

iteration is a recent development based on the basic algebraic-geometric properties for SOCLCP (1.1) with  $M$  having the so-called *globally uniquely solvable property*<sup>1</sup> (GUS) [16, 17, 18, 42] (see Definition 2.1 below); BN was employed and imbedded successfully into the matrix-splitting framework in [44] to solve the linear complementary problem (1.2).

The BN iteration solves SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) from an entirely different perspective from other methods in the literature (e.g., [5, 6, 8, 13, 20, 22, 24, 32, 37, 39, 40, 41]), and it is quite efficient for small-to-medium scale problems [43]. However, for large scale problems, there are two bottlenecks: it needs to solve the unique positive eigenvalue of  $M^\top J_n$  and the associated eigenvector accurately, and repeatedly it has to solve many linear systems with coefficient matrices  $M - sJ_n$ , where  $s$  is some varying shift and

$$J_n := \text{diag}(1, -1, \dots, -1) \in \mathbb{R}^{n \times n}. \quad (1.3)$$

In this paper, we will propose an efficient algorithm for SOCLCP (1.1). Our numerical tests show that it is competitive to the BN iteration for small-to-medium scale problems, but more importantly it is suitable for large scale problems.

First, we will show that SOCLCP (1.1) with  $M$  having the GUS property can be solved by finding a particular positive zero  $s_*$  of a rational function  $h(s)$ . We will show that this function  $h(s)$  enjoys many nice geometry properties on  $s \in (0, +\infty)$  and indeed it can only admit one or two zeros on  $(0, +\infty)$ . Based on the rich techniques in the model reduction, we then propose a Krylov subspace method to approximate  $h(s)$  by a Padé-type approximation  $h_\ell(s)$ . With the aid of the geometry properties of  $h(s)$ , the positive zero  $s_*$  of  $h(s)$  can be efficiently computed by solving  $h_\ell(s) = 0$  through a small eigenvalue problem. Our new method does not rely upon the eigen-information of  $M^\top J_n$  and is suitable for large scale problems. We tested our method and compared with the BN iteration and two available MATLAB software packages: SDPT3 [39, 40, 41] and SeDuMi [37], for small-to-medium dense problems as well as for large scale sparse problems. Preliminary numerical experiments show that our method is rather efficient in either case in terms of accuracy and speed.

The rest of this article is organized as follows. In section 2, we present some basic results for SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ), primarily with the GUS property, and then give a brief description of the BN iteration [43] in section 3. We take a new look at SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) from the view point of finding the positive zero  $s_*$  of  $h(s)$  in section 4, which is followed by a detailed theoretical analysis of its geometry property in section 5. Section 6 is dedicated to the discussion for finding the zero  $s_* > 0$  of  $h(s)$ : the Padé-type approximation  $h_\ell(s)$  of  $h(s)$  via Krylov subspace techniques is first introduced, and then the strategy of obtaining an approximation of  $s_*$  is presented. Our numerical tests and comparison with the BN iteration, SDPT3 and SeDuMi are reported in section 7. Final concluding remarks are drawn in section 8.

**Notation:** Throughout this paper, all vectors are column vectors and are typeset in bold lower case letters. For  $A \in \mathbb{R}^{n \times m}$  (the set of all  $m \times n$  real matrices),  $A^\top$  denotes its transpose, and  $\mathcal{R}(A)$  and  $\mathcal{N}(A)$  represent the range and kernel of  $A$ , respectively. Thus

---

<sup>1</sup>A complete algebraic-geometric characterization of this property will be detailed in Theorem 2.2.

$\mathcal{R}(A)^\perp = \mathcal{N}(A^\top)$ , where  $\mathcal{R}(A)^\perp$  denotes the orthogonal complement of  $\mathcal{R}(A)$ . As usual, the identity matrix in  $\mathbb{R}^{n \times n}$  will be denoted by  $I_n \equiv [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n]$ , where  $\mathbf{e}_i$  is its  $i$  column. We shall also adopt MATLAB-like convention to access the entries of vectors and matrices. For example,  $\mathbf{x}_{(i)}$  is the  $i$ th element of  $\mathbf{x}$  and  $A_{(i,j)}$  is the  $(i, j)$ th entry of  $A$ , where  $(i : j)$  stands for the set of integers from  $i$  to  $j$  inclusive, and  $A_{(k:\ell, i:j)}$  is the submatrix of  $A$  that consists of intersections from row  $k$  to row  $\ell$  and column  $i$  to column  $j$ . For a matrix already with a subscript, e.g.  $Y_r$ , we write  $Y_{r;(k:\ell, i:j)}$  for  $(Y_r)_{(k:\ell, i:j)}$ . Notation  $\|\cdot\|_2$  is either the matrix spectral norm or the Euclidean vector norm, depending on its arguments:

$$\|\mathbf{x}\|_2 = \sqrt{\sum_i |\mathbf{x}_{(i)}|^2}, \quad \|A\|_2 := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|_2}{\|\mathbf{x}\|_2}.$$

The  $\ell$ th Krylov subspace generated by  $A \in \mathbb{R}^{n \times n}$  on  $\mathbf{x} \in \mathbb{R}^n$  is defined as

$$\mathcal{K}_\ell(A, \mathbf{x}) = \text{span}(\mathbf{x}, A\mathbf{x}, \dots, A^{\ell-1}\mathbf{x}).$$

## 2 Basic Properties of SOCLCP

In this section, we review several basic results for  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$ . These results form the theoretical basis of our development later. For the ease of presentation, we define

$$M_s := M - sJ_n, \quad \mathbb{K}_s := M_s \mathbb{K}^n = \{M_s \mathbf{x} : \mathbf{x} \in \mathbb{K}^n\}. \quad (2.1)$$

Evidently,  $\mathbb{K}_0 = M\mathbb{K}^n$ . Since  $\mathbb{K}_s$  approaches  $-\mathbb{K}^n$  as  $s \rightarrow +\infty$  ([42, Lemma 9]), we define

$$\mathbb{K}_\infty := -\mathbb{K}^n.$$

Finally,  $\partial(\mathbb{K}^n)$  and  $\text{int}(\mathbb{K}^n)$  stand for the boundary and the interior of  $\mathbb{K}^n$ , respectively.

**Lemma 2.1** ([42]). *The following statements hold:*

- (i) *For nonzero vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n$ ,  $\mathbf{x}^\top \mathbf{y} = 0$  if and only if  $\mathbf{x} \in \partial(\mathbb{K}^n)$ ,  $\mathbf{y} \in \partial(\mathbb{K}^n)$  and  $\mathbf{y} = sJ_n \mathbf{x}$  for some  $s > 0$ .*
- (ii) *Let  $\mathbf{x} \in \mathbb{K}^n$ . Then  $\mathbf{x}^\top \mathbf{y} > 0$  for all nonzero vector  $\mathbf{y} \in \mathbb{K}^n$  if and only if  $\mathbf{x} \in \text{int}(\mathbb{K}^n)$ .*

The next theorem characterizes the three mutually exclusive cases for  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$ . We point out that if  $M$  is nonsingular, (C2) can be equivalently restated as  $-M^{-1}\mathbf{q} \in \mathbb{K}^n$  (which implies that  $\mathbf{x} = -M^{-1}\mathbf{q}$  is the solution). In (C3),  $M - s_*J_n$  may or may not be singular.

**Theorem 2.1** ([43, 44]). *There are three mutually exclusive cases for the solution  $\mathbf{x} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$ , namely:*

- (C1)  $\mathbf{q} \in \mathbb{K}^n$  (which implies that  $\mathbf{x} = \mathbf{0}$  is the solution);
- (C2)  $\text{SOL}(M, \mathbb{K}^n, \mathbf{q}) \supseteq \{\mathbf{x} \in \mathbb{K}^n : M\mathbf{x} + \mathbf{q} = \mathbf{0}\} \neq \emptyset$ ;

(C3) *there exists  $s_* > 0$  such that  $M\mathbf{x} + \mathbf{q} = s_* J_n \mathbf{x} \in \partial(\mathbb{K}^n)$ .*

Next we introduce the so-called GUS property on  $M$ , previously discussed in [16, 17, 18, 42].

**Definition 2.1.**  $M \in \mathbb{R}^{n \times n}$  is said to have the *globally uniquely solvable* (GUS) property if the solution  $\mathbf{x}$  of  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  is unique for any given  $\mathbf{q} \in \mathbb{R}^n$ .

The next theorem provides an algebraic-geometric characterization of such a property. The reader is referred to [42, 43, 44] for proofs and more.

**Theorem 2.2.**  $M \in \mathbb{R}^{n \times n}$  has the GUS property if and only if the following three statements hold<sup>2</sup>.

(i)  $MJ_n$  has exactly one positive eigenvalue  $\tau$  in  $[0, +\infty)$ . Moreover,

$$\text{rank}(MJ_n - \tau I_n) = \text{rank}(M_\tau J_n) = n - 1 \quad (2.2)$$

and  $MJ_n$  has an eigenvector  $\mathbf{w} \in \text{int}(\mathbb{K}^n)$  associated with  $\tau$ ;

(ii)  $M^\top J_n$  has an eigenvector  $\mathbf{v} \in \text{int}(\mathbb{K}^n)$  associated with  $\tau$ , and  $\mathcal{R}(M_\tau) = (J_n \mathbf{v})^\perp$ ;

(iii) for all  $\mathbf{a} \in \partial(\mathbb{K}^n)$ ,  $\mathbf{a}^\top M \mathbf{a} \geq 0$  and  $\mathbf{a}^\top M^{-1} \mathbf{a} \geq 0$ .

When  $M$  has the GUS property, the following statements hold

(a) For  $0 < s < \tau < t$ ,  $\mathbb{K}_0 \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_s)$ , and  $\mathbb{K}_\infty \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_t)$ . Moreover,  $\mathbb{K}_s$  and  $\mathbb{K}_t$  lie in the two opposite sides of  $\mathcal{R}(M_\tau)$  which is a hyperplane in  $\mathbb{R}^n$ .

(b)  $\mathbb{K}_t \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_s)$  for  $0 < t < s < \tau$  and for  $\tau < s < t$ .

(c) If  $\mathbf{x} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$  is not in the cases (C1) and (C2), i.e.,  $\mathbf{q} \notin (-M\mathbb{K}^n) \cup \mathbb{K}^n$ , then there exists a unique  $s_* > 0$  such that  $M\mathbf{x} + \mathbf{q} = s_* J_n \mathbf{x}$  and

$$s_* \begin{cases} = \tau, & \text{if } (-\mathbf{q})^\top J_n \mathbf{v} = 0 \text{ or equivalently, } \mathbf{q} \in \mathcal{R}(M_\tau), \\ < \tau, & \text{if } (-\mathbf{q})^\top J_n \mathbf{v} > 0, \\ > \tau, & \text{if } (-\mathbf{q})^\top J_n \mathbf{v} < 0. \end{cases}$$

We remark that Theorem 2.2(i) doesn't say anything about the algebraic multiplicity of  $\tau$ , i.e, it doesn't exclude the case that  $\tau$  is a multiple eigenvalue algebraically. However, if  $\tau$  indeed is an algebraically multiple eigenvalue, the corresponding eigenvector of  $MJ_n$  is unique, modulo a nonzero scalar factor, because of (2.2).

Theorem 2.2(i) implies that 0 cannot be an eigenvalue of  $MJ_n$  and thus  $M$  must be nonsingular if  $M$  has the GUS property. the GUS property or not. But an immediate consequence of this theorem is that *a symmetric positive definite  $M$  has the GUS property*. In fact, more can be said: 1)  $M$  such that  $M + M^\top$  is positive definite has the GUS property

---

<sup>2</sup>We note that  $MJ_n$  has the same eigenvalues as  $(MJ_n)^\top = J_n M^\top$  which has the same eigenvalues as  $J_n^{-1}(J_n M^\top)J_n = M^\top J_n$ .

[18, Theorem 17], and 2) a symmetric  $M$  has the GUS property if and only if it is positive definite (implied by [18, Theorem 21]).

In the case when  $M$  is symmetric positive definite,  $M - \lambda J_n$  is the so-called positive definite pencil [28, Definition 1.1] since  $M - 0 \cdot J_n = M$  is positive definite. It follows from a more general result there that  $M - \lambda J_n$  is diagonalizable and has one simple positive eigenvalue and  $n - 1$  negative eigenvalues [28, Lemma 3.8]. So in this case,  $\tau$  in Theorem 2.2(i) is a simple eigenvalue.

By Theorem 2.1, we know that if  $\mathbf{x} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$  is not in the cases (C1) and (C2), then it must be in the case (C3), i.e., there exists  $s_* > 0$  such that  $M\mathbf{x} + \mathbf{q} = s_* J_n \mathbf{x}$ . If also  $M$  has the GUS property, then  $\mathbf{x}$  is unique and thus,  $s_*$  is unique, too. By Theorem 2.2(c), we can locate  $s_*$  with the information of the eigenpair  $(\tau, \mathbf{v})$  of  $M^\top J_n$ . This is exactly what the BN iteration was designed to do [43].

Throughout the rest of this paper, our target is the second-order cone linear complementarity problem  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  with  $M$  having the GUS property, and the assignments to  $\tau, s_*, \mathbf{q}, \mathbf{v}, \mathbf{w}$  in Theorem 2.2 will all be kept.

### 3 The BN iteration

The BN iterative method [43] is based on Theorem 2.1, where the cases (C1) and (C2) are rather trivial. It is the case (C3) that is conceivably more common and needs most work. The BN iteration combines bisection and the Newton method to generate a sequence  $\{s_k\}$  that converges to  $s_*$  by investigating whether

$$\mathbf{x}(s_k) := -(M - s_k J_n)^{-1} \mathbf{q}$$

is in  $\mathbb{K}^n$  or not. We refer to [43, Algorithm 3] and [44, Algorithm 2] for more detail.

The iteration needs to know the eigenpair  $(\tau, \mathbf{v})$  of  $M^\top J_n$  associated with the unique positive eigenvalue  $\tau$  of  $M^\top J_n$ . Therefore, the main computational burden is to find the eigenpair  $(\tau, \mathbf{v})$  and to compute  $\mathbf{x}(s_k)$  in each step. The latter has a cost of  $O(n^3)$  flops in general, but can be much less for certain special structured  $M$ . For example, when  $M$  is in the upper Hessenberg form, it costs only  $O(n^2)$  to compute  $\mathbf{x}(s_k)$ . In general we can always reduce  $M$  to an upper Hessenberg matrix through an orthogonal transformation (more detail in subsection 6.5). This reduction of  $M$  to an upper Hessenberg matrix itself still costs  $O(n^3)$  flops, but it needs to be done only once.

The BN iteration with a pre-processing: transforming  $M$  to  $Q^\top M Q$  is rather efficient for small-to-medium scale problems. However, for large scale problems, such a transformation is too costly because it costs  $O(n^3)$  flops. Thus there are two computational bottlenecks in this approach, namely,

1. the need of the eigenpair  $(\tau, \mathbf{v})$  of  $M^\top J_n$ , and
2. the transformation of  $M$  to  $Q^\top M Q$  which in general cannot exploit, e.g., much of the sparsity structure of  $M$  for a large scale  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$ .

In this paper, we will attempt to circumvent both computational bottlenecks by developing numerical methods that are efficient for large scale  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  for which

either  $M$  is sparse or the matrix-vector product with  $M$  is fast (at a cost much less than  $O(n^2)$  flops that is needed usually for a dense  $M$ ). Our new method still focuses on the case (C3).

## 4 Transform SOCLCP into a zero-finding problem

Our new method relies on the following simple observation.

**Theorem 4.1.** *Suppose  $s \in \mathbb{R}$  is not an eigenvalue of  $MJ_n$ . Let*

$$\mathbf{x}(s) \equiv \begin{bmatrix} x_1(s) \\ \mathbf{x}_2(s) \end{bmatrix} := -(M - sJ_n)^{-1}\mathbf{q}, \quad (4.1)$$

where  $J_n$  is given by (1.3). Then  $\mathbf{x}(s) \in \partial(\mathbb{K}^n)$  if and only if  $x_1(s) > 0$  and  $h(s) = 0$ , where

$$h(s) := \mathbf{x}(s)^\top J_n \mathbf{x}(s) = \mathbf{q}^\top (M - sJ_n)^{-\top} J_n (M - sJ_n)^{-1} \mathbf{q}. \quad (4.2)$$

*Proof.* It is evident by the definition of  $\mathbb{K}^n$ .  $\square$

According to Theorem 4.1, we know that finding  $s_*$  for the case (C3) in Theorem 2.1 is equivalent to finding the zero  $s_*$  of  $h(s)$ . Notice that

$$h(s) = \mathbf{q}^\top [s^2 J_n - (M + M^\top)s + MJ_n M^\top]^{-1} \mathbf{q} \quad (4.3)$$

is a rational function that involves the inverse of an  $n \times n$  quadratic matrix-valued function. Such a function has the same form as the so-called *second-order transfer functions* that are sought to be reduced in, e.g., [2, 3, 9, 12, 26, 36, 38], by the name of the *second-order model reduction*. The basic idea, in the context of  $h(s)$  here, is to approximate  $h(s)$  by a Padé-type approximation  $h_\ell(s)$  having the same form but involving the inverse of a much smaller, say  $\ell \times \ell$ , quadratic matrix-valued function. Usually  $\ell \ll n$ , say  $\ell = 10$  up to 50 for example. Some of the *poles* (i.e., values of  $s$  at which  $h_\ell(s)$  becomes  $\infty$ ) and zeros of  $h_\ell(s)$  are often good approximations to some of those of  $h(s)$ , respectively. Note from (4.1) and (4.2) that generically, the unique positive eigenvalue  $\tau$  of  $MJ_n$  (which is also the one of  $M^\top J_n$ ) is a pole of  $h(s)$ , and thus conceivably one of the poles of  $h_\ell(s)$  should provide a good approximation of  $\tau$  (see the details in sections 5 and 6 below).

Thanks to powerful techniques developed in the model reduction, the (rational) Krylov subspace techniques can be employed to construct  $h_\ell(s)$  at price of only solving certain linear systems involving  $M$ . This makes it possible for us to numerically solve large scale SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) because the involved linear systems can be solved iteratively through exploiting the sparsity of  $M$  or fast matrix-vector multiplications by  $M$ . Furthermore, this approach does not require the eigenpair  $(\tau, \mathbf{v})$  of  $M^\top J_n$ , and it even can produce accurate approximations to  $\tau$  by the poles of  $h_\ell(s)$  (see Theorem 5.1 below). We shall discuss these issues in detail in section 6.

In Theorem 2.2(c), the case  $s_* \neq \tau$  is generic and should occur more frequently than otherwise and thus the case should be considered more carefully. On the other hand, Zhang and Yang [43] explained how to compute the solution  $\mathbf{x} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$  by a direct method [43, Algorithm 2] for the special case  $s_* = \tau$ . Therefore we will not invest much effort in the special case any further from now on.



## 5 Analyze the curve of $h(s)$

In this section, we shall present a theoretical analysis of the curve of  $h(s)$  defined in (4.2). By relying upon the special structure of  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$ , our analysis reveals some basic geometry properties of  $h(s)$  both for the generic case  $s_* \neq \tau$  and for the special case  $s_* = \tau$ . These findings are helpful in guiding us to compute the wanted zero and pole of  $h(s)$ .

### 5.1 The general case $s_* \neq \tau$

We first analyze the curve  $h(s)$  of (4.2) for the generic case:  $s_* \neq \tau$ . According to Theorem 2.2(c),  $s_* \neq \tau$  if and only if  $\mathbf{v}^\top J_n \mathbf{q} \neq 0$  or, equivalently,  $\mathbf{q} \notin \mathcal{R}(M_\tau)$ .

As a convention, in what follows, we will use “ $h(s) = -$ ” to mean  $h(s) < 0$  and likewise “ $h(s) = +$ ” to mean  $h(s) > 0$ .

**Lemma 5.1.** *Suppose that  $M$  has the GUS property and  $\mathbf{q} \notin (-M\mathbb{K}^n) \cup \mathbb{K}^n$ . We have the following statements.*

(a) If  $s_* \in (0, \tau)$ , then

$$h(s) = \begin{cases} - & \text{for } s \in (0, s_*), \\ 0 & \text{for } s = s_*, \\ + & \text{for } s \in (s_*, \tau); \end{cases}$$

(b) If  $s_* \in (\tau, \infty)$ , then

$$h(s) = \begin{cases} + & \text{for } s \in (\tau, s_*), \\ 0 & \text{for } s = s_*, \\ - & \text{for } s \in (s_*, \infty). \end{cases}$$

*Proof.* First, we know that under the conditions, the solution  $\mathbf{x}(s_*) \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$  is on the boundary  $\partial(\mathbb{K}^n)$  of  $\mathbb{K}^n$ , where  $\mathbf{x}(s)$  is defined by (4.1).

Consider (a), where  $s \in (0, s_*)$ . We claim that  $\mathbf{x}(s) = -M_s^{-1}\mathbf{q} \notin \mathbb{K}^n \cup (-\mathbb{K}^n)$ , which immediately implies  $h(s) \equiv [x_1(s)]^2 - \|\mathbf{x}_2(s)\|_2^2 < 0$ . Otherwise, if  $\mathbf{x}(s) \in \mathbb{K}^n$ , by Theorem 2.2(b), we have

$$-\mathbf{q} = M_s \mathbf{x}(s) \in \mathbb{K}_s \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_{s_*}),$$

implying  $-\mathbf{q} = M_{s_*} \mathbf{y}$  for some  $\mathbf{y} \in \text{int}(\mathbb{K}^n)$ . Thus  $\mathbf{x}(s_*) = -M_{s_*}^{-1}\mathbf{q} = \mathbf{y} \in \text{int}(\mathbb{K}^n)$ , contradicting  $\mathbf{x}(s_*) \in \partial(\mathbb{K}^n)$ . Similarly,

$$\mathbf{x}(s) \in -\mathbb{K}^n \Rightarrow \mathbf{q} = M_s(-\mathbf{x}(s)) \in \mathbb{K}_s \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_{s_*}) \Rightarrow -\mathbf{x}(s_*) \in \text{int}(\mathbb{K}^n),$$

contradicting again  $\mathbf{x}(s_*) \in \partial(\mathbb{K}^n)$ . Now, for  $s_* < s < \tau$ , from Theorem 2.2(b), we have

$$-\mathbf{q} = M_{s_*} \mathbf{x}(s_*) \in \mathbb{K}_{s_*} \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_s) \Rightarrow \mathbf{x}(s) \in \text{int}(\mathbb{K}_s) \Rightarrow h(s) > 0.$$

The case (b) can be proved similarly. For  $s > s_*$ , we claim that  $\mathbf{x}(s) = -M_s^{-1}\mathbf{q} \notin \mathbb{K}^n \cup (-\mathbb{K}^n)$ , which directly implies  $h(s) < 0$ . Otherwise, by Theorem 2.2(b), we have

$$\mathbf{x}(s) \in \mathbb{K}^n \Rightarrow -\mathbf{q} = M_s \mathbf{x}(s) \in \mathbb{K}_s \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_{s_*}) \Rightarrow \mathbf{x}(s_*) \in \text{int}(\mathbb{K}^n),$$

contradicting  $\mathbf{x}(s_*) \in \partial(\mathbb{K}^n)$ . Similarly,

$$\mathbf{x}(s) \in -\mathbb{K}^n \Rightarrow \mathbf{q} = M_s(-\mathbf{x}(s)) \in \mathbb{K}_s \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_{s_*}) \Rightarrow -\mathbf{x}(s_*) \in \text{int}(\mathbb{K}^n),$$

contradicting again  $\mathbf{x}(s_*) \in \partial(\mathbb{K}^n)$ . Now, for  $\tau < s < s_*$ , from Theorem 2.2(b), we have

$$-\mathbf{q} = M_{s_*}\mathbf{x}(s_*) \in \mathbb{K}_{s_*} \setminus \{\mathbf{0}\} \subset \text{int}(\mathbb{K}_s) \Rightarrow \mathbf{x}(s) \in \text{int}(\mathbb{K}_s) \Rightarrow h(s) > 0.$$

The proof is completed.  $\square$

We now present our main theorem in this subsection. It gives a completed picture of  $h(s)$  on  $(0, \infty)$  as detailed in Table 5.1 for the current case. In particular, it shows that  $h(s)$  have at least one but at most two positive zeros in  $(0, \infty)$ . For cases 3 and 4 in the table,  $\mathbf{q} \notin \mathcal{R}(M_\tau)$  is redundantly added for clarity since it is in fact implied by  $\mathbf{q} \in (-\mathbb{K}^n) \setminus (-M\mathbb{K}^n)$  for case 3 and by  $\mathbf{q} \in M\mathbb{K}^n \setminus \mathbb{K}^n$  for case 4.

Table 5.1: The sign of  $h(s)$  on  $(0, \infty)$  in terms of location of  $\mathbf{q}$

cases	where is $\mathbf{q}$ ?	solution $\mathbf{x}_*$ or $h(s)$
1	$\mathbf{q} \in \mathbb{K}^n$	$\mathbf{0} = \mathbf{x}_* \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$
2	$\mathbf{q} \in -M\mathbb{K}^n$	$-M^{-1}\mathbf{q} = \mathbf{x}_* \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$
3	$\mathbf{q} \in (-\mathbb{K}^n) \setminus (-M\mathbb{K}^n),$ $\mathbf{q} \notin \mathcal{R}(M_\tau)$	$h(s) = \begin{cases} - & \text{for } s \in (0, s_*), \\ 0 & \text{for } s = s_*, \\ + & \text{for } s \in (s_*, \tau), \\ + & \text{for } s \in (\tau, \infty). \end{cases}$
4	$\mathbf{q} \in M\mathbb{K}^n \setminus \mathbb{K}^n,$ $\mathbf{q} \notin \mathcal{R}(M_\tau)$	$h(s) = \begin{cases} + & \text{for } s \in (0, \tau), \\ + & \text{for } s \in (\tau, s_*), \\ 0 & \text{for } s = s_*, \\ - & \text{for } s \in (s_*, \infty). \end{cases}$
5	$\mathbf{q} \notin (-M\mathbb{K}^n) \cup M\mathbb{K}^n \cup \mathbb{K}^n \cup (-\mathbb{K}^n),$ $\mathbf{q} \notin \mathcal{R}(M_\tau)$	$h(s) = \begin{cases} - & \text{for } s \in (0, s_{*;1}), \\ 0 & \text{for } s = s_{*;1}, \\ + & \text{for } s \in (s_{*;1}, \tau), \\ + & \text{for } s \in (\tau, s_{*;2}), \\ 0 & \text{for } s = s_{*;2}, \\ - & \text{for } s \in (s_{*;2}, \infty). \end{cases}$

**Theorem 5.1.** *Suppose  $M$  has the GUS property. If  $\mathbf{q} \notin \mathcal{R}(M_\tau)$  (and thus  $s_* \neq \tau$ ), then*

- (1) *the sign of  $h(s)$  for  $s \in (0, +\infty)$  can be characterized by cases 3, 4, and 5 in Table 5.1. In cases 3 and 4,  $h(s)$  has one positive root  $s_*$ , and in case 5, it has two positive roots  $s_{*;1}$  and  $s_{*;2}$  one of which is  $s_*$ , and*
- (2)  *$\lim_{s \rightarrow \tau} h(s) = +\infty$ , which combining with Table 5.1, implies that  $\tau$  is the unique pole of  $h(s)$  in  $[0, +\infty)$ .*

*Proof.* (1) We first remark that the five cases of  $\mathbf{q}$  presented in Table 5.1 are mutually exclusive and also complete for  $\mathbf{q} \in \mathbb{R}^n \setminus \mathcal{R}(M_\tau)$ . The first two cases are (C1) and (C2) in Theorem 2.1.

We consider case 3. Note that  $\mathbf{q} \in (-\mathbb{K}^n) \setminus (-M\mathbb{K}^n)$  implies  $\mathbf{q}^\top J_n \mathbf{v} < 0$  and by Theorem 2.2(c), we know that  $0 < s_* < \tau$ . Thus, using Lemma 5.1, we only need to show  $h(s) > 0$  for  $s > \tau$ , which is true because  $-\mathbb{K}^n \setminus \{\mathbf{0}\} \subseteq \text{int}(\mathbb{K}_s)$  by Theorem 2.2(a), and thus

$$\mathbf{q} \in (-\mathbb{K}^n) \setminus (-M\mathbb{K}^n) \subseteq \text{int}(\mathbb{K}_s) \Rightarrow -\mathbf{x}(s) \in \text{int}(\mathbb{K}^n) \Rightarrow h(s) > 0.$$

For case 4, we note that for  $s \in (0, \tau)$ , it follows from Theorem 2.2(a,b) that

$$\mathbf{q} \in M\mathbb{K}^n \setminus \mathbb{K}^n = \mathbb{K}_0 \setminus \mathbb{K}^n \subset \text{int}(\mathbb{K}_s) \Rightarrow -\mathbf{x}(s) \in \text{int}(\mathbb{K}^n) \Rightarrow h(s) > 0,$$

which also implies that  $s_* > \tau$ . Therefore, combining it with Lemma 5.1, yields the conclusion.

For the last case, note that  $\mathbf{q} \notin (-M\mathbb{K}^n) \cup M\mathbb{K}^n \cup \mathbb{K}^n \cup (-\mathbb{K}^n) \cup \mathcal{R}(M_\tau)$  implies both  $\mathbf{q} \notin (-M\mathbb{K}^n) \cup \mathbb{K}^n$  and  $-\mathbf{q} \notin (-M\mathbb{K}^n) \cup \mathbb{K}^n$ . We now prove the result for  $\mathbf{v}^\top J_n \mathbf{q} < 0$  and  $\mathbf{v}^\top J_n \mathbf{q} > 0$ .

- a)  $\mathbf{v}^\top J_n \mathbf{q} < 0$ . By Theorem 2.2(c), we know  $0 < s_* < \tau$ . So  $s_{*;1} = s_*$ . On the other hand, consider the problem SOCLCP( $M, \mathbb{K}^n, -\mathbf{q}$ ). Since  $-\mathbf{q} \notin (-M\mathbb{K}^n) \cup \mathbb{K}^n$ , we know also from Theorem 2.2(c) that there is another  $s_{*;2} > 0$  such that

$$-(M - s_{*;2}J_n)^{-1}(-\mathbf{q}) = -\mathbf{x}(s_{*;2}) \in \text{SOL}(M, \mathbb{K}^n, -\mathbf{q}).$$

Moreover, because  $\mathbf{v}^\top J_n(-\mathbf{q}) > 0$ , by applying Theorem 2.2(c) to the problem SOCLCP( $M, \mathbb{K}^n, -\mathbf{q}$ ), we conclude that  $s_{*;2} > \tau$ . Consequently, based on Lemma 5.1, the assertion follows.

- b)  $\mathbf{v}^\top J_n \mathbf{q} > 0$ . The proof for this situation is similar to the first situation a), except that  $s_*$  now is  $s_{*;2}$ . The details are omitted.

(2) Pre-multiplying  $(M - sJ_n)\mathbf{x}(s) = -\mathbf{q}$  by  $\mathbf{v}^\top J_n$  on both sides and using  $(J_n \mathbf{v})^\top M = \tau \mathbf{v}^\top$ , we know that for any  $0 < s \neq \tau$ ,

$$(\tau - s)\mathbf{v}^\top \mathbf{x}(s) = -\mathbf{v}^\top J_n \mathbf{q} \neq 0 \tag{5.1}$$

and hence,

$$\lim_{s \rightarrow \tau} |\mathbf{v}^\top \mathbf{x}(s)| = +\infty \quad \text{and} \quad \lim_{s \rightarrow \tau} \|\mathbf{x}(s)\|_2 = +\infty.$$

Let  $\mathbf{z}$  be any accumulation point of  $\frac{\mathbf{x}(s)}{\|\mathbf{x}(s)\|_2}$  as  $s \rightarrow \tau$ , i.e., there exists a sequence  $s_1, s_2, \dots$  such that  $s_i \rightarrow \tau$  and  $\frac{\mathbf{x}(s_i)}{\|\mathbf{x}(s_i)\|_2} \rightarrow \mathbf{z}$  as  $i \rightarrow +\infty$ . Dividing  $(M - sJ_n)\mathbf{x}(s) = -\mathbf{q}$  by  $\|\mathbf{x}(s)\|_2$ , we have for  $0 < s \neq \tau$

$$(MJ_n - sI_n)J_n \frac{\mathbf{x}(s)}{\|\mathbf{x}(s)\|_2} = -\frac{\mathbf{q}}{\|\mathbf{x}(s)\|_2}. \tag{5.2}$$

Putting  $s = s_i$  in (5.2) and letting  $i \rightarrow +\infty$ , we have  $(MJ_n - \tau I_n)J_n \mathbf{z} = \mathbf{0}$ . Since  $\|J_n \mathbf{z}\|_2 = 1$ , by Theorem 2.2(i),

$$\text{either } J_n \mathbf{z} = \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \quad \text{or} \quad J_n \mathbf{z} = -\frac{\mathbf{w}}{\|\mathbf{w}\|_2}.$$

In other words, the only possible accumulation points of  $\frac{\mathbf{x}(s)}{\|\mathbf{x}(s)\|_2}$  are  $\mathbf{z} = \pm J_n \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ . This fact together with  $\lim_{s \rightarrow \tau} \|\mathbf{x}(s)\|_2 = +\infty$  are sufficient for us to conclude  $\lim_{s \rightarrow \tau} h(s) = +\infty$ . In fact, if the opposite were true, then there would be a sequence  $s_1, s_2, \dots$  converging to  $\tau$  such that  $\{h(s_i)\}$  is bounded. We can assume without loss of generality that  $\lim_{i \rightarrow +\infty} h(s_i) = \varphi < +\infty$ . Note that we can choose a subsequence  $s_{i_1}, s_{i_2}, \dots$  such that

$$\lim_{j \rightarrow +\infty} \frac{\mathbf{x}(s_{i_j})}{\|\mathbf{x}(s_{i_j})\|_2} = J_n \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \quad \text{or} \quad \lim_{j \rightarrow +\infty} \frac{\mathbf{x}(s_{i_j})}{\|\mathbf{x}(s_{i_j})\|_2} = -J_n \frac{\mathbf{w}}{\|\mathbf{w}\|_2}.$$

But in either case,

$$\varphi = \lim_{j \rightarrow +\infty} h(s_{i_j}) = \lim_{j \rightarrow +\infty} \frac{\mathbf{x}(s_{i_j})^\top J_n \mathbf{x}(s_{i_j})}{\|\mathbf{x}(s_{i_j})\|_2^2} \|\mathbf{x}(s_{i_j})\|_2^2 = \frac{\mathbf{w}^\top J_n \mathbf{w}}{\|\mathbf{w}\|_2^2} \lim_{j \rightarrow +\infty} \|\mathbf{x}(s_{i_j})\|_2^2 = +\infty,$$

a contradiction, where we have used the fact that  $\mathbf{w} \in \text{int}(\mathbb{K}^n)$  and thus  $\mathbf{w}^\top J_n \mathbf{w} > 0$ . This completes the proof.  $\square$

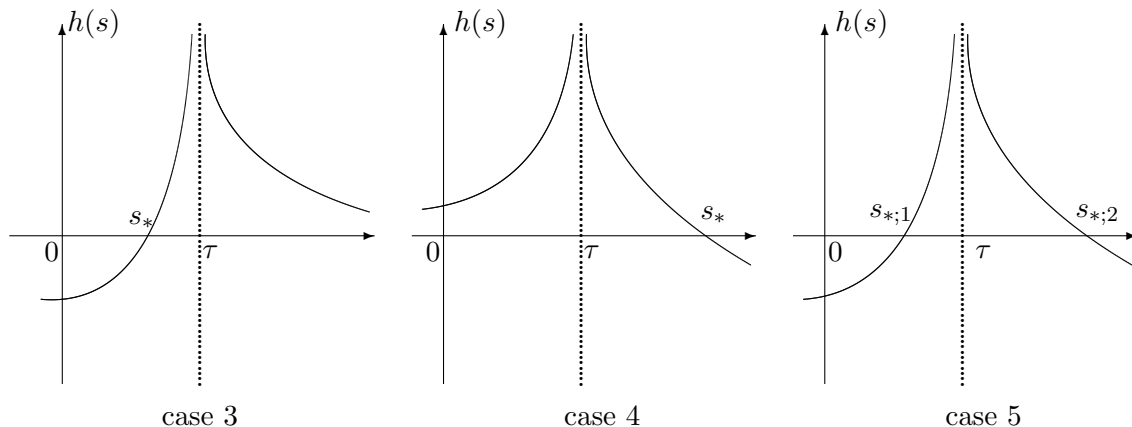


Figure 5.1:  $h(s)$  corresponding to cases 3, 4, and 5 in Table 5.1.

**Remark 5.1.** Except for the trivial cases: cases 1 and 2, Theorem 5.1 reveals a complete picture of  $h(s)$ , which is illustrated in Figure 5.1. We point out that the three patterns of  $h(s)$  in Figure 5.1 are one-to-one corresponding to the three locations of  $\mathbf{q}$ . These one-to-one correspondences are important for correctly finding the right approximation of  $s_*$  via the reduced  $h_\ell(s)$  (see section 6.1). Before reducing  $h(s)$  to  $h_\ell(s)$ , we will first check if it is in case 1 or case 2. If neither case occurs, then it must fall into one of the cases 3, 4, and 5. Now the location of  $\mathbf{q}$  and, thereby, the sign pattern of  $h(s)$  or, equivalently, the relationship between  $\tau$  and  $s_*$ , are then clear.

## 5.2 The special case $s_* = \tau$

We now analyze the curve  $h(s)$  for the special case  $s_* = \tau$ . The next proposition implies that  $h(s) < 0$  for all  $0 < s \neq \tau$ .

**Theorem 5.2.** *Suppose  $M$  has the GUS property and  $\mathbf{q} \in \mathcal{R}(M_\tau)$ . Then  $h(s) < 0$  for all  $0 < s \neq \tau$ .*

*Proof.* We show that  $\mathbf{x}(s) \notin \mathbb{K}^n \cup (-\mathbb{K}^n)$  and therefore  $h(s) < 0$  for  $0 < s \neq \tau$ . Suppose on the contrary that  $\mathbf{x}(s) \in \mathbb{K}^n \cup (-\mathbb{K}^n)$ . Then either  $\mathbf{x}(s) \in \mathbb{K}^n \Rightarrow \mathbf{q} \in \mathbb{K}_s$  which contradicts  $\mathbf{q} \in \mathcal{R}(M_\tau)$  because  $\mathbb{K}_s \cap \mathcal{R}(M_\tau) = \emptyset$ , or  $\mathbf{x}(s) \in -\mathbb{K}^n \Rightarrow \mathbf{q} \in -\mathbb{K}_s$  which contradicts  $\mathbf{q} \in \mathcal{R}(M_\tau)$  because  $(-\mathbb{K}_s) \cap \mathcal{R}(M_\tau) = \emptyset$ . This establishes our assertion.  $\square$

Recall our discussion immediately after Theorem 2.2 that we cannot exclude the possibility that algebraically  $\tau$  is a multiple eigenvalue of  $M^\top J_n$ . But if  $\tau$  is indeed a *simple* eigenvalue of  $M^\top J_n$ , the following theorem says that  $\tau$  is no longer a pole of  $h(s)$ , i.e., it is a removable singularity. In such a case, by Theorem 5.2, we know  $\lim_{s \rightarrow \tau} h(s) \leq 0$  with a possibility of being an equality.

**Theorem 5.3.** *Suppose  $M$  has the GUS property. If  $\mathbf{q} \in \mathcal{R}(M_\tau)$  (and thus  $s_* = \tau$ ) and if  $\tau$  is a simple eigenvalue of  $M^\top J_n$ , then  $h(s)$  is analytic in  $(0, +\infty)$ . In particular, if  $M$  is symmetric positive definite, then  $h(s)$  is analytic in  $(0, +\infty)$ .*

*Proof.* It suffices to prove  $h(s)$  is analytic at  $s = \tau$ . To this end, we first note from Theorem 2.2 that when  $M$  has the GUS property, we have

$$M^\top J_n \mathbf{v} = \tau \mathbf{v} \Rightarrow \mathbf{v}^\top J_n M = \tau \mathbf{v}^\top.$$

Therefore, if  $\tau$  is a *simple* eigenvalue of  $M^\top J_n$ , by the Jordan canonical decomposition of a matrix, there exists a nonsingular matrix  $X \in \mathbb{R}^{n \times n}$  such that

$$X J_n M = \begin{bmatrix} \tau & \mathbf{0} \\ \mathbf{0} & B \end{bmatrix} X \quad \text{with} \quad X_{(1,:)} = \mathbf{v}^\top,$$

i.e., the first row of  $X$  is  $\mathbf{v}^\top$ . Thus,

$$\begin{aligned} \mathbf{x}(s) &= -(M - sJ_n)^{-1} \mathbf{q} \\ &= -(J_n M - sI_n)^{-1} J_n \mathbf{q} \\ &= -X^{-1} \begin{bmatrix} \frac{1}{\tau-s} & \mathbf{0} \\ \mathbf{0} & (B - sI_{n-1})^{-1} \end{bmatrix} X J_n \mathbf{q} \\ &= -X^{-1} \begin{bmatrix} \frac{1}{\tau-s} & \mathbf{0} \\ \mathbf{0} & (B - sI_{n-1})^{-1} \end{bmatrix} \begin{bmatrix} 0 \\ \mathbf{z}_2 \end{bmatrix} \\ &= -X^{-1} \begin{bmatrix} 0 \\ (B - sI_{n-1})^{-1} \mathbf{z}_2 \end{bmatrix}, \end{aligned}$$

where we have used  $XJ_n\mathbf{q} = \begin{bmatrix} \mathbf{v}^\top J_n\mathbf{q} \\ X_{(2:n,:)}J_n\mathbf{q} \end{bmatrix}$  and  $\mathbf{v}^\top J_n\mathbf{q} = 0$  by Theorem 2.2(c), and set  $\mathbf{z}_2 := X_{(2:n,:)}J_n\mathbf{q}$ . Therefore,

$$\begin{aligned} h(s) &= \mathbf{x}(s)^\top J_n\mathbf{x}(s) \\ &= \begin{bmatrix} 0, \mathbf{z}_2^\top (B - sI_{n-1})^{-\top} \end{bmatrix} X^{-\top} J_n X^{-1} \begin{bmatrix} 0 \\ (B - sI_{n-1})^{-1} \mathbf{z}_2 \end{bmatrix} \\ &= \mathbf{z}_2^\top (B - sI_{n-1})^{-\top} C (B - sI_{n-1})^{-1} \mathbf{z}_2, \end{aligned} \quad (5.3)$$

where  $C = (X^{-\top} J_n X^{-1})_{(2:n,2:n)} \in \mathbb{R}^{(n-1) \times (n-1)}$ . Since every eigenvalue of  $B$  is also an eigenvalue of  $M^\top J_n$  and  $\tau$  is assumed a simple eigenvalue of  $M^\top J_n$ , we conclude that  $\tau$  cannot be an eigenvalue of  $B$ , i.e.,  $B - \tau I_{n-1}$  is nonsingular, and thus  $\tau$  is not a pole of  $h(s)$ .  $\square$

**Remark 5.2.** Theorem 5.3 reveals the similarity between  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  and the trust-region subproblem [31, Chapter 4]:

$$\min_{\mathbf{x}^\top \mathbf{x} = \delta > 0} f(\mathbf{x}) \quad \text{with} \quad f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^\top M \mathbf{x} + \mathbf{q}^\top \mathbf{x}. \quad (5.4)$$

Indeed, when  $M$  is symmetric,  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  in (1.1) is just the KKT condition of the following quadratic second-order cone programming

$$\min_{\mathbf{x}^\top J_n \mathbf{x} = 0, x_1 > 0} f(\mathbf{x}) \quad \text{with} \quad f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^\top M \mathbf{x} + \mathbf{q}^\top \mathbf{x}, \quad (5.5)$$

and  $s$  in  $h(s)$  is the Lagrangian multiplier. By comparing (5.4) to (5.5), we note both have the same objective function but different constraints. Moreover, for (5.4), a real-valued rational function  $\|p(\lambda)\|_2$  of the Lagrangian multiplier  $\lambda$ , as the counterpart of  $h(s)$  here, can also be defined [31, (4.37)]. It is well-know that in the general case,  $-\lambda_{\min}(M)$  (the negative of the smallest eigenvalue of  $M$ ) is a pole of  $\|p(\lambda)\|_2$ , but, like the special case  $s = s_*$  in  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$ , the Lagrangian multiplier associated with the optimal solution is  $\lambda_* = -\lambda_{\min}(M)$  and  $\lambda_*$  is not a pole in the ‘‘hard case’’.

**Remark 5.3.** Theorems 5.1, 5.2 and 5.3 are important because they together provide us with ways to detect the special case  $\mathbf{q} \in \mathcal{R}(M_\tau)$ . In particular, Theorems 5.1 and 5.2 simply suggest that  $\mathbf{q} \in \mathcal{R}(M_\tau)$  if  $h(s) < 0$  for  $0 < s \neq \tau$ ; so one practical strategy is to check if  $h(s) > 0$  for  $s$  near  $\tau$  (see Figure 5.1).

## 6 Compute the positive zero $s_*$ of $h(s)$

### 6.1 Reduce $h(s)$ via the Arnoldi process

There are several approaches (see e.g., [2, 36, 38]) to reduce the ‘‘transfer function’’  $h(s)$  in (4.2) which can be related to a time-invariant dynamical system. In particular in the second-order form as in (4.3), there are two existing treatments:

1. linearizing  $h(s)$  and reducing it *via* the bi-orthogonal Lanczos process by Gallivan, Grimme, and Van Dooren [14] and by Feldman and Freund [12], and
2. reducing  $h(s)$  *via* SOAR (Second-Order Arnoldi) by Bai and Su [2, 3].

Note our “transfer function”  $h(s)$  bears the factor form (4.2), but the reduced  $h_\ell(s)$  by either treatment above does not preserve the factor form. In addition, *serious breakdowns* may occur.

We opt to design a special reduction procedure, similar to [27], which respects the factor form in (4.2) and uses only the classical Arnoldi process. We will also show that our reduced  $h_\ell(s)$  corresponds to another SOCLCP in the form of (1.1) of much smaller scale. Moreover, the Taylor series of  $h_\ell(s)$  agrees with that of  $h(s)$  in their first  $\ell$  terms at the expansion point  $s_0$ . More detailed discussions on the advantages of our factor-form-preserving reduction can be found in Remarks 6.1 and 6.2 below.

Before we present our reduction procedure, we first discuss a commonly used shifting technique. Suppose<sup>3</sup>  $0 < s_0 \in \mathbb{R}$  is an arbitrary but otherwise fixed expansion point. Let

$$A_{s_0} = M_{s_0}^{-1} J_n, \quad \mathbf{b}_{s_0} = -M_{s_0}^{-1} \mathbf{q}, \quad (6.1)$$

where  $M_{s_0}$  is defined by (2.1). Write  $s = s_0 + \sigma$ . We have<sup>4</sup>

$$\begin{aligned} \mathbf{x}(s) &= -(M - sJ_n)^{-1} \mathbf{q} \\ &= -(M_{s_0} - \sigma J_n)^{-1} \mathbf{q} \\ &= -(I_n - \sigma M_{s_0}^{-1} J_n)^{-1} (M_{s_0}^{-1} \mathbf{q}) \\ &= (I_n - \sigma A_{s_0})^{-1} \mathbf{b}_{s_0}, \\ h(s) &= \mathbf{x}(s)^\top J_n \mathbf{x}(s) \\ &= \mathbf{b}_{s_0}^\top (I_n - \sigma A_{s_0})^{-\top} J_n (I_n - \sigma A_{s_0})^{-1} \mathbf{b}_{s_0} \\ &=: g(\sigma). \end{aligned}$$

Suppose  $s_0$  is neither a pole nor a zero of  $h(s)$  (in other words,  $h(s_0) \neq 0$  and is finite). This implies  $h(s_0) = \mathbf{b}_{s_0}^\top J_n \mathbf{b}_{s_0} \neq 0$ . Our reduction process begins with the Arnoldi process as in Algorithm 6.1 to generate the Krylov subspace  $\mathcal{K}_\ell(A_{s_0}, \mathbf{b}_{s_0})$ .

It can be seen that if the process does not break down at line 8 till  $j = \ell$ , then we have in exact arithmetic

$$A_{s_0} Y_\ell = Y_\ell H_\ell + H_{(\ell+1, \ell)} \mathbf{y}_{\ell+1} \mathbf{e}_\ell^\top \quad \text{and} \quad Y_\ell^\top A_{s_0} Y_\ell = H_\ell, \quad (6.2)$$

where  $Y_\ell = [\mathbf{y}_1, \dots, \mathbf{y}_\ell]$ ,  $H_\ell := H_{(1:\ell, 1:\ell)} \in \mathbb{R}^{\ell \times \ell}$  is an upper Hessenberg matrix. Now, our reduced  $h_\ell(s)$  is given by

$$h_\ell(s) = \|\mathbf{b}_{s_0}\|_2^2 \mathbf{e}_1^\top (I_\ell - \sigma H_\ell)^{-\top} Y_\ell^\top J_n Y_\ell (I_\ell - \sigma H_\ell)^{-1} \mathbf{e}_1 =: g_\ell(\sigma). \quad (6.3)$$

<sup>3</sup>As far as the Arnoldi process in Algorithm 6.1 and the approximation property in Theorem 6.1 are concerned, it is not necessary to have  $s_0 \in \mathbb{R}$ , not to mention  $s_0 > 0$ . We are making this assumption here mainly because only such an  $s_0$  interests us because it is intended to approximate  $s_*$  in the end.

<sup>4</sup>Because of  $s = s_0 + \sigma$ , a function in  $s$  is a function in  $\sigma$ , too, and vice versa. For this reason, we conveniently write  $h(s)$ , when regarded as a function in  $\sigma$ , as  $g(\sigma)$ . Similarly for their reduced ones, we have notations  $h_\ell(s)$  and  $g_\ell(\sigma)$ .

---

**Algorithm 6.1** The Arnoldi Process

---

Given  $s_0$  which is neither a pole nor a zero of  $h(s)$ , this procedure computes an orthonormal basis matrix of  $\mathcal{K}_\ell(A_{s_0}, \mathbf{b}_{s_0})$ , where  $A_{s_0}$  and  $\mathbf{b}_{s_0}$  are defined by (6.1).

---

- 1: solve  $M_{s_0} \mathbf{b}_{s_0} = -\mathbf{q}$  for  $\mathbf{b}_{s_0}$ ;
  - 2:  $\mathbf{y}_1 = \mathbf{b}_{s_0} / \|\mathbf{b}_{s_0}\|_2$ ;  $Y_1 = [\mathbf{y}_1]$ ;
  - 3: **for**  $j = 1, 2, \dots, \ell$  **do**
  - 4:   solve  $M_{s_0} \mathbf{w} = J_n \mathbf{y}_j$  for  $\mathbf{w}$ ;
  - 5:    $H_{(1:j,j)} = Y_j^\top \mathbf{w}$ ,  $\mathbf{w} = \mathbf{w} - Y_j H_{(1:j,j)}$ ;
  - 6:    $\beta = \|\mathbf{w}\|_2$ ;
  - 7:   **if**  $\beta = 0$  **then**
  - 8:     **BREAK**;
  - 9:   **else**
  - 10:      $H_{(j+1,j)} = \|\mathbf{w}\|_2$ ,  $\mathbf{y}_{j+1} = \mathbf{w}/\beta$ ,  $Y_{j+1} = [Y_j, \mathbf{y}_{j+1}]$ ;
  - 11:   **end if**
  - 12: **end for**
  - 13: **return**  $Y_\ell = [\mathbf{y}_1, \dots, \mathbf{y}_\ell]$ , an orthonormal basis matrix of  $\mathcal{K}_\ell(A_{s_0}, \mathbf{b}_{s_0})$ , and  $H$ .
- 

The following theorem implies that the Taylor series of  $h_\ell(s)$  at  $s_0$  matches that of  $h(s)$  also at  $s_0$  in their first  $\ell$  terms. The proof follows similarly to that in [27], but for completeness, we include it here.

**Theorem 6.1.** *Suppose Algorithm 6.1 runs to its completion without breakdown to produce  $H_\ell$ . Let  $h(s)$  be defined by (4.2) and  $h_\ell(s)$  by (6.3). Then*

$$h(s) = h_\ell(s) + O(|s - s_0|^\ell). \quad (6.4)$$

*Proof.* The key for proving (6.4) is

$$A_{s_0}^j \mathbf{y}_1 = Y_\ell H_\ell^j \mathbf{e}_1 \quad \text{for } j \leq \ell - 1 \quad (6.5)$$

which can be proved by induction as follows. First, for  $j = 1$ , by (6.2), we have  $A_{s_0} \mathbf{y}_1 = A_{s_0} Y_\ell \mathbf{e}_1 = Y_\ell H_\ell \mathbf{e}_1$ . Suppose (6.5) is true for  $j - 1 \leq \ell - 2$ . Then

$$\begin{aligned} A_{s_0}^j \mathbf{y}_1 &= A A_{s_0}^{j-1} \mathbf{y}_1 = A_{s_0} Y_\ell H_\ell^{j-1} \mathbf{e}_1 \\ &= Y_\ell H_\ell H_\ell^{j-1} \mathbf{e}_1 + H_{(\ell+1,\ell)} \mathbf{y}_{\ell+1} \mathbf{e}_\ell^\top H_\ell^{j-1} \mathbf{e}_1 \\ &= Y_\ell H_\ell^j \mathbf{e}_1, \end{aligned}$$

where we have used  $\mathbf{e}_\ell^\top H_\ell^{j-1} \mathbf{e}_1 = 0$  because  $H_\ell^{j-1}$  is a banded matrix with a lower bandwidth  $j$ . Thus, for sufficiently tiny  $|\sigma|$ ,

$$\begin{aligned} \mathbf{x}(s) &= (I - \sigma A_{s_0})^{-1} \mathbf{b}_{s_0} = \sum_{i=1}^{\infty} \sigma^i A_{s_0}^i \mathbf{b}_{s_0} \\ &= \|\mathbf{b}_{s_0}\|_2 \sum_{i=1}^{\infty} \sigma^i A_{s_0}^i \mathbf{y}_1 \end{aligned}$$



$$\begin{aligned}
&= \|\mathbf{b}_{s_0}\|_2 \left[ \sum_{i=1}^{\ell-1} \sigma^i Y_\ell H_\ell^i \mathbf{e}_1 + O(\sigma^\ell) \right] \\
&= \|\mathbf{b}_{s_0}\|_2 Y_\ell \left[ \sum_{i=1}^{\ell-1} \sigma^i H_\ell^i \mathbf{e}_1 + O(\sigma^\ell) \right] \\
&= \|\mathbf{b}_{s_0}\|_2 Y_\ell \left[ (I_\ell - \sigma H_\ell)^{-1} \mathbf{e}_1 + O(\sigma^\ell) \right].
\end{aligned}$$

Substituting the above relation into  $h(s) = \mathbf{x}(s)^\top J_n \mathbf{x}(s)$  yields (6.4).  $\square$

**Remark 6.1.** By comparing our factor-form-preserving reduction procedure with the two existing treatments, i.e., linearization+reduction *via* the bi-orthogonalization Lanczos process or reduction *via* SOAR, the computational costs of ours is lower in order to obtain the same order reduction. In particular, it can be seen that only one linear system of size  $n$  is needed for each order of accuracy achieved in our factor-form-preserving reduction procedure, while two linear systems of size  $n$  are needed in the other two treatments. For the storage, since both our reduction procedure and SOAR are based on the Arnoldi process, the storage requirements are almost the same. The linearization+reduction *via* the bi-orthogonalization Lanczos process is based on the three-term recurrence, and the storage is smaller. Serious breakdowns may occur in the two existing methods, and that makes them less stable or more involved if the look-ahead strategy [35] is used as a cure. On the contrary, a breakdown (i.e.,  $\beta = 0$  occurs at line 7) in the Arnoldi process in Algorithm 6.1 is always welcome, because it then finds an invariant subspace and thus makes  $h(s) \equiv h_\ell(s)$ .

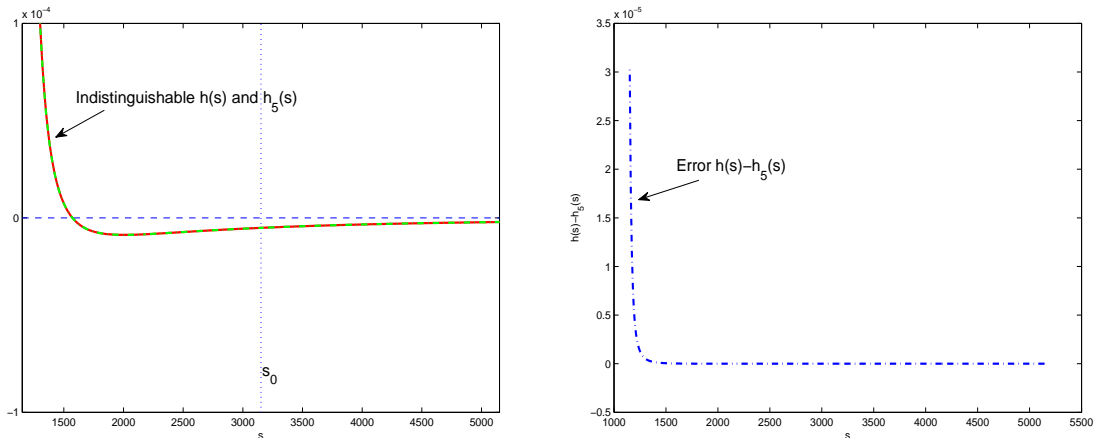


Figure 6.1: Example 6.1:  $h(s)$  vs.  $h_\ell(s)$ . *Left:* indistinguishable  $h(s)$  and  $h_5(s)$ ; *Right:* error  $h(s) - h_5(s)$ .

**Example 6.1.** As an illustration, we test this reduction on  $M \in \mathbb{R}^{66 \times 66}$  of BCSSTK02 from the matrix market<sup>5</sup>. This  $M$  is symmetric positive definite. Although it is small

<sup>5</sup><http://math.nist.gov/MatrixMarket/>.



The poles of  $h(s) = (J_n \mathbf{q})^\top (MJ_n - sI_n)^{-\top} (MJ_n - sI_n)^{-1} \mathbf{q}$  correspond to some of the eigenvalues of  $MJ_n$ . Depending on  $\mathbf{q}$ , some eigenvalues of  $MJ_n$  may give removable poles of  $h(s)$ , but generically, the poles and the eigenvalues are the same. By Theorem 2.2, we know that if  $M$  has the GUS property,  $\tau$  is the unique positive pole of  $h(s)$  (see Theorem 5.1 for more detailed discussions on the pole  $\tau$ ). This observation suggests that we can obtain an approximation of  $\tau$  by finding the positive pole(s) of  $h_\ell(s)$ , if any. From (6.3), the poles  $\tau_1$  of  $h_\ell(s)$  relate to the eigenvalues  $\nu$  of  $H_\ell$  by  $\tau_1 = s_0 + 1/\nu$ . Note that  $H_\ell$  might have complex eigenvalues, and as a practical strategy, we choose

$$\eta = \max \{ \operatorname{Re}(\mu) : 1/\mu \in \operatorname{eig}(H_\ell), |\operatorname{Im}(\mu)| \leq 10^{-6} \} \quad (6.9)$$

and use  $\tau_1 = \eta + s_0$  as an approximation of  $\tau$ .

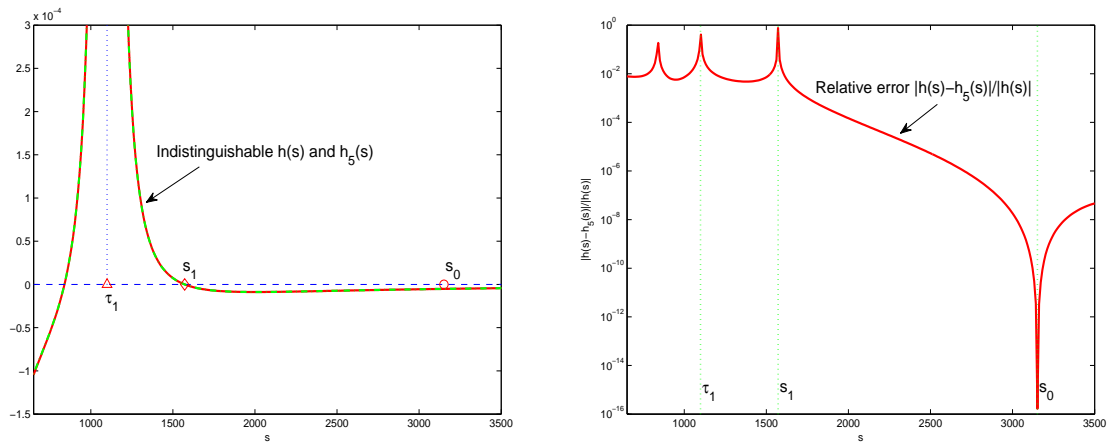


Figure 6.2: Example 6.2. *Left*: indistinguishable  $h(s)$  and  $h_5(s)$  (with limited  $y$ -axis range so that the useful details can be recognized), the expansion point  $s_0$ , and the approximate zero  $s_1$  and pole  $\tau_1$  by  $h_5(s)$ ; *Right*: small relative error  $|h(s) - h_5(s)|/|h(s)|$  away from zeros and poles over a wide range of  $s$ .

To demonstrate the effectiveness of the above-mentioned ways for the approximating the poles and zeros of  $h(s)$ , we now revisit Example 6.1 (more numerical results will be reported in section 7).

**Example 6.2.** Continuing with Example 6.1, Figure 6.2 compares  $h_5(s)$  with  $h(s)$  over an even larger range of  $s$  to the left of  $s_0$  to include a pole of  $h_5(s)$  with  $h(s)$ . Now, by computing the eigenvalues of  $Q(\sigma)$  and  $\widehat{Q}(\sigma)$ , we find the approximation  $\tau_1 \approx 1.0992 \times 10^3$  of  $\tau$  and the approximation  $s_1 \approx 1.5721 \times 10^3$  of  $s_*$  with  $h(s_1) \approx 2.82 \times 10^{-8}$ . Interestingly, Figure 6.2 shows another zero around  $\hat{s}_1 = 8.3912 \times 10^2$  of  $h(s)$ . This zero is not the solution for SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) since the first element of  $\mathbf{x}(\hat{s}_1) = -(M - \hat{s}_1 J_n)^{-1} \mathbf{q}$  is negative, i.e.,  $\mathbf{x}(\hat{s}_1) \notin \partial(\mathbb{K}^n)$ . Indeed, this is an example for the case 5 in Theorem 5.1:  $8.3912 \times 10^2$  and  $1.5721 \times 10^3$  are the approximations to  $s_{*;1}$  and  $s_{*;2}$ , respectively, but  $s_* = s_{*;2}$  is the desired one. So this example falls into the case 5 there.  $\diamond$

Next we look at an example for the special case  $s_* = \tau$ . Theorem 5.2 claims that  $h(s) < 0$  for  $0 < s \neq \tau$  and Theorem 5.3 says  $\tau$  is no longer a pole if  $\tau$  is a simple eigenvalue of  $MJ_n$ .

**Example 6.3.** We again use the matrix  $M \in \mathbb{R}^{66 \times 66}$  of BCSSTK02 as an illustration. To construct the special case, we first compute the largest eigenvalue  $\tau \approx 1.0996 \times 10^3$  of  $MJ_n$ , and then generate the solution

$$\mathbf{x} = [\sqrt{n-1}, 1, \dots, 1]^\top \in \mathbb{R}^{66},$$

and finally set  $\mathbf{q} = -(M - \tau J_n)\mathbf{x}$ . The function  $h(s)$  and the reduced  $h_5(s)$  associated with the example are plotted in Figure 6.3 which shows that indeed  $h(s) < 0$  in the plotted range of  $s > 0$  and  $\tau$  is not a pole any more, and also  $h_5(s)$  has no poles on  $(0, \infty)$ . This would be the case for which [43, Algorithm 2] is triggered to solve this SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ).  $\diamond$

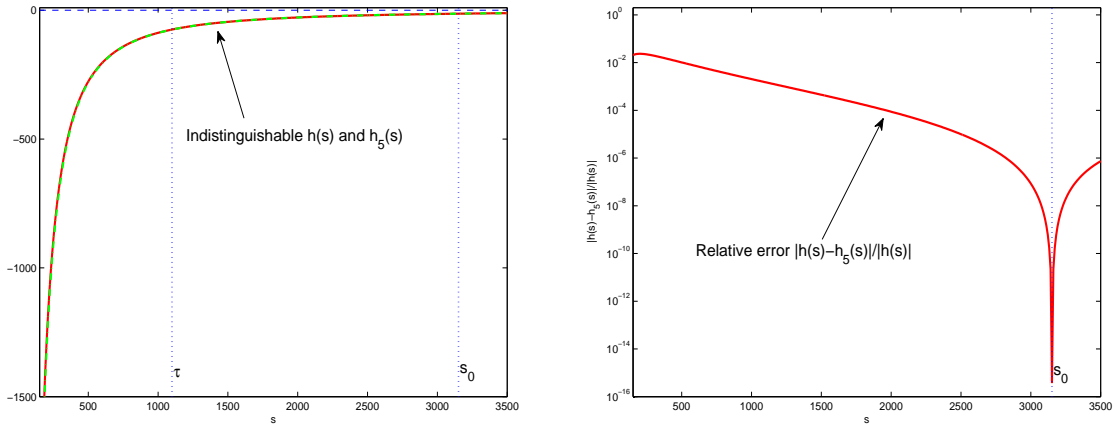


Figure 6.3: Example 6.3. The special  $s_* = \tau$  which is no longer a pole any more. *Left:* indistinguishable  $h(s)$  and  $h_5(s)$ ; *Right:* relative error  $|h(s) - h_5(s)|/|h_5(s)|$ .

### 6.3 The Reduced SOCLCP

In this subsection, we will establish an interesting connection of the factor-form-preserving reduced  $h_\ell(s)$  to an SOCLCP in the form of (1.1) of a much smaller size. To this end, we first establish the following proposition.

**Lemma 6.1.** *Let  $Y_\ell = [\mathbf{y}_1, \dots, \mathbf{y}_\ell] \in \mathbb{R}^{n \times \ell}$  ( $1 \leq \ell < n$ ) be the orthonormal basis matrix by Algorithm 6.1, i.e.,  $Y_\ell^\top Y_\ell = I_\ell$ . Then  $Y_\ell^\top J_n Y_\ell$  has an eigenvalue  $\varrho := 2\|Y_{\ell;(1,:)}\|_2^2 - 1 \in [-1, 1]$ , and all other eigenvalues are  $-1$ .*

*Proof.* Note  $J_n = 2\mathbf{e}_1\mathbf{e}_1^\top - I_n$  and  $Y_\ell^\top Y_\ell = I_\ell$ . We have

$$Y_\ell^\top J_n Y_\ell = 2Y_\ell^\top \mathbf{e}_1\mathbf{e}_1^\top J_n Y_\ell - I_\ell = 2Y_{\ell;(1,:)}^\top Y_{\ell;(1,:)} - I_\ell,$$

where  $Y_{\ell;(1,:)}$  is the first row of  $Y_\ell$ . Hence,  $Y_\ell^\top J_n Y_\ell$  has an eigenvalue  $\varrho = 2\|Y_{\ell;(1,:)}\|_2^2 - 1$ , and all other eigenvalues are  $-1$ . Due to  $1 \leq \ell < n$ , we have that  $0 \leq \|Y_{\ell;(1,:)}\|_2^2 \leq 1$  and thus  $-1 \leq \varrho \leq 1$ .  $\square$

Since  $Y_\ell$  has orthogonal columns,  $\varrho > 0$  for sufficiently large  $\ell$ . Whenever  $\varrho > 0$ , the matrix  $Y_\ell^\top J_n Y_\ell$  admits the following decomposition:

$$\begin{aligned} Y_\ell^\top J_n Y_\ell &= P_\ell \operatorname{diag}(\varrho, -I_{\ell-1}) P_\ell^\top \quad (\text{eigen-decomposition}) \\ &= P_\ell \operatorname{diag}(\sqrt{\varrho}, I_{\ell-1}) J_\ell \operatorname{diag}(\sqrt{\varrho}, I_{\ell-1}) P_\ell^\top, \end{aligned} \quad (6.10)$$

where  $P_\ell \in \mathbb{R}^{\ell \times \ell}$  is orthogonal, i.e.,  $P_\ell^\top P_\ell = I_\ell$ . Substituting (6.10) into  $h_\ell(s)$  yields

$$\begin{aligned} h_\ell(s) &= \|\mathbf{b}_{s_0}\|_2^2 \mathbf{e}_1^\top (\widehat{P}_\ell - \sigma H_\ell \widehat{P}_\ell)^{-\top} J_\ell (\widehat{P}_\ell - \sigma H_\ell \widehat{P}_\ell)^{-1} \mathbf{e}_1 \\ &= \|\mathbf{b}_{s_0}\|_2^2 \mathbf{q}_\ell^\top (M_\ell - \sigma J_\ell)^{-\top} J_\ell (M_\ell - \sigma J_\ell)^{-1} \mathbf{q}_\ell \end{aligned} \quad (6.11)$$

which corresponds to SOCLCP( $M_\ell, \mathbb{K}^\ell, \mathbf{q}_\ell$ ), where  $\widehat{P}_\ell = P_\ell \operatorname{diag}(1/\sqrt{\varrho}, I_{\ell-1})$ , and

$$M_\ell = J_\ell \widehat{P}_\ell^{-1} H_\ell^{-1} \widehat{P}_\ell \quad \text{and} \quad \mathbf{q}_\ell = J_\ell \widehat{P}_\ell^{-1} H_\ell^{-1} \mathbf{e}_1. \quad (6.12)$$

**Theorem 6.2.** *In Lemma 6.1, if  $\varrho > 0$ , then  $h_\ell$  corresponds to SOCLCP( $M_\ell, \mathbb{K}^\ell, \mathbf{q}_\ell$ ) with  $M_\ell$  and  $\mathbf{q}_\ell$  given by (6.12).*

**Remark 6.2.** As was explained, by connecting the function  $h_\ell(s)$  with an SOCLCP in the form of (1.1), we know from (6.11) that solving  $h_\ell(s) = g_\ell(\sigma) = 0$  is equivalent to solving a much smaller SOCLCP( $M_\ell, \mathbb{K}^\ell, \mathbf{q}_\ell$ ). Therefore, if we have an efficient algorithm for solving small to medium size SOCLCPs in the form of (1.1),  $h_\ell(s) = 0$  can also be tackled efficiently. From this point of view, *our factor-form-preserving reduction procedure can be viewed as a projection technique that projects the original large scale SOCLCP (1.1) to a smaller SOCLCP (1.1) with the property that the Taylor series of the corresponding functions  $h(s)$  and  $h_\ell(s)$  at the expansion point match till the  $\ell$ th order term.*

## 6.4 Algorithmic framework: LCPvA

Our foregoing discussions naturally lead to our main algorithmic framework to solve SOCLCP (1.1) by Krylov subspace reduction *via* the Arnoldi process. It is outlined in Algorithm 6.2. However, in order to make the algorithm more efficient and stable, there are still some issues needed to be addressed. In subsection 6.5, we shall first suggest a strategy to improve the efficiency for the small-to medium-size dense problems, and in subsection 6.6, we will discuss stability issue for large scale SOCLCP.

## 6.5 Transform dense SOCLCP to upper Hessenberg SOCLCP

The matrix  $M$  in this subsection is assumed to be of small-to-medium sizes and is dense. We note that the main computationally heavy part of Algorithm 6.1 lie at its line 4:

$$\text{solve } M_{s_0} \mathbf{w} = J_n \mathbf{y}_j \text{ for } \mathbf{w}. \quad (6.13)$$

---

**Algorithm 6.2** LCPvA: Linear Complementarity Problem via the Arnold process for SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ )

---

**input:**  $M, \mathbf{q}$ , and  $k_{\max}$  (maximum number of outer-iterations allowed);

**output:** an approximation to  $\mathbf{x} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$ .

---

```

1: if  $\mathbf{q} \in \mathbb{K}^n$  then  $\mathbf{x} = \mathbf{0}$  and return;
2: if  $-M^{-1}\mathbf{q} \in \mathbb{K}^n$  then  $\mathbf{x} = -M^{-1}\mathbf{q}$  and return;
3: select  $s_0$  and set  $k = -1$ ;
4: repeat
5:    $k = k + 1$ ;
6:   apply Algorithm 6.1 (with full re-orthogonalization) to generate  $H_\ell$  and  $Y_\ell$  associated
   with  $\mathcal{K}_\ell(A_{s_k}, \mathbf{b}_{s_k})$ ;
7:   solve QEP (6.7) for the approximation  $s_{k+1}$  of  $s_*$  according to (6.8);
8:   solve  $(M - s_{k+1}J_n)\mathbf{x}(s_{k+1}) = -\mathbf{q}$  for  $\mathbf{x}(s_{k+1})$ ;
9:   if  $h(s_{k+1}) \leq \epsilon \|\mathbf{x}(s_{k+1})\|_2^2$  and  $x_1(s_{k+1}) > 0$  then
10:    return  $\mathbf{x}(s_{k+1})$  as an approximate solution of SOCLCP( $M, \mathbb{K}^n, \mathbf{w}$ ), and LCPvA
    stops;
11:  end if
12: until  $k \geq k_{\max}$ ;
13: if  $k \geq k_{\max}$  and there is no positive zero of  $h_\ell(s)$  then
14:   call [43, Algorithm 2] for the special case  $s_* = \tau$ ;
15: else
16:   LCPvA fails to find an approximate solution.
17: end if

```

---

These are linear systems like  $M_{s_0}\mathbf{t} = \mathbf{r}$  and usually cost  $O(n^3)$  flops for every first use of a new expansion point  $s_0$  and  $O(n^2)$  flops afterwards. Since  $s_0$  will have to be updated as it has to move towards  $s_*$ , linear systems  $M_{s_0}\mathbf{t} = \mathbf{r}$  will have to be solved for a few different  $s_0$ , and thus they could be computationally expensive.

Previously, we mentioned that if  $M$  is in the upper/lower Hessenberg form, the upper/lower triangular form included, then (6.13) takes only  $O(n^2)$  flops to compute, independent of  $s_0$ . This is good.

But what happens when  $M$  is not in the upper/lower Hessenberg form? When the size of  $M$  is modest, say up to a few thousands, and  $M$  is also dense, we can preprocess SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) to transform it into one for which  $M$  is in the upper Hessenberg form. Here is the detail. Find  $Q = \text{diag}(1, Q_0)$  with orthogonal  $Q_0 \in \mathbb{R}^{(n-1) \times (n-1)}$  such that  $Q^\top M Q$  is upper Hessenberg. Such a transformation from  $M$  to  $Q^\top M Q$  is not new and in fact it is the first step by the name of the *Hessenberg reduction* in solving a dense eigenvalue problem by the QR algorithm [11]. Important for us here is that it does not change the second-order cone  $\mathbb{K}^n$ :  $Q\mathbb{K}^n = \mathbb{K}^n$  and at the same time  $Q^\top J_n Q = J_n$ . Therefore, we can solve SOCLCP( $Q^\top M Q, \mathbb{K}^n, Q^\top \mathbf{q}$ ) instead: any  $\mathbf{z} \in \text{SOL}(Q^\top M Q, \mathbb{K}^n, Q^\top \mathbf{q})$  gives  $\mathbf{x} = Q\mathbf{z} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$  and vice versa.

It should be pointed out that the Hessenberg reduction itself, although only needed to be done once, still costs  $O(n^3)$  flops which can be too expensive to be practical for a

very large scale problem. Nevertheless, for large scale and sparse  $M$ , the involved linear systems can be solved iteratively through exploiting the sparsity of  $M$  or fast matrix-vector multiplications by  $M$  (at a cost usually much less than  $O(n^2)$  flops that is needed for a dense  $M$ ).

Comparing with BN in [43], the new approach has the following noticeable advantage: it does not need to compute the eigenpair  $(\tau, \mathbf{v})$  of  $M^\top J_n$ , while the computational complexity is almost the same as that of BN.

## 6.6 Reorthogonalization to prevent loss of orthogonality

A potential problem in the Arnoldi process is loss of orthogonality. There are extensive discussions on this issue, for example in [11, Chapter 7] and [34] and the references therein. It is argued that the loss of orthogonality does not cause the algorithm to behave completely unpredictable. Indeed, loss of orthogonality could result in multiple copies of the same eigenvalues. There are also several cures for such loss of orthogonality, and the Arnoldi algorithm with full reorthogonalization is one. Another alternative is selective reorthogonalization [34] by orthogonalizing any new Arnoldi vector against already converged Ritz vectors. Our choice is simply to do full reorthogonalization. Our goal is to find the zero  $s_*$  of  $h(s)$  by iteratively (the outer-loop) choosing new and better shift  $s_0$  obtained from the eigenvalues of the QEP for  $\widehat{Q}(\sigma)$  associated with the previous Arnoldi process; from this point of view, we do not need Krylov subspaces of very large dimension in each outer-loop, and therefore cost for doing full reorthogonalization is usually manageable.

## 7 Numerical experiments

We coded Algorithm 6.2 (LCPvA) in MATLAB 7.13.0 (R2011b) and tested it against other efficient algorithms on an iMac ME086CH/A with Intel Core i5@2.7GHz and 8GB memory. As we shall see, preliminary results demonstrate the high efficiency of LCPvA for SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ).

We divide our experiments into two parts: (1) numerical tests for (medium sized) dense  $M$ , and (2) for large scale sparse  $M$ . In our current implementation, for dense  $M$ , we preprocess the associated SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) as we discussed in subsection 6.5 to make sure all subsequent linear systems are solved in  $O(n^2)$  flops. But for sparse  $M$ , we still use MATLAB's sparse LU decomposition to solve the linear system at line 8 of Algorithm 6.2 as well as the one at line 4 of Algorithm 6.1 in which one LU should be computed before the **for**-loop and subsequently we just use triangular system solving. This is because an LU decomposition costs potentially much more than two triangular system solvings. The use of the LU decomposition in the sparse case limits the size of  $M$  that we can test on. In the future, we will explore solving these linear systems iteratively to allow even larger  $M$ .

All tested  $M$  are symmetric and positive definite (implying  $M$  has the GUS property) and their condition numbers are varied from 10 to  $10^5$ . As a comparison, we also report the numerical results from three other methods: the BN iteration [43], and

two popular MATLAB software packages: SDPT3 [39, 40, 41] (version 4)<sup>6</sup> and SeDuMi [37] (version 1.3). We point out that both SDPT3 and SeDuMi are designed for general semidefinite-quadratic-linear programming, not particularly targeting at SOCLCP in the form of (1.1), but when  $M$  is symmetric and positive definite, solving  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  is equivalent to solving the following convex quadratic second-order cone programming (because  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  is indeed its KKT conditions):

$$\min_{\mathbf{x} \in \mathbb{K}^n} \frac{1}{2} \mathbf{x}^\top M \mathbf{x} + \mathbf{q}^\top \mathbf{x}. \quad (7.1)$$

In fact, if  $M$  is decomposed into or already takes the form

$$M = \widetilde{M}^\top \widetilde{M}, \quad (7.2)$$

then according to [40, section 4.6] and [37], (7.1) can be converted equivalently to the standard programming problems that can be solved by SDPT3 and SeDuMi. Following [40, section 4.6] and [37], we outline the conversion procedure in Appendix A.

We run SDPT3 with its default setups and run the BN iteration with the parameters used in the numerical testing in [43]; for SeDuMi, we set `pars.eps` =  $10^{-10}$  as the desired accuracy. As for our method LCPvA, the involved parameters and their recommended ranges are summarized in Table 7.1. The use of the particular  $s_0$  as in the table is drawn from our own experience on numerous tests.

Table 7.1: Parameters in LCPvA

parameter	value	description
$\varepsilon$	$10^{-7} \sim 10^{-4}$	Stopping criterion at line 9 in Algorithm 6.2
$s_0$	$\frac{\ M\ _1}{5}$	initial shift in Algorithm 6.2
$k_{\max}$	$\geq 3$	maximal number of iterations for the outer loop of Algorithm 6.2
$\text{iter}_{\text{Arnoldi}}$	$\mathbb{Z}^{k_{\max}}$	$\text{iter}_{\text{Arnoldi}}(i)$ is the number of Arnoldi steps in the $i$ th outer loop of Algorithm 6.2

Finally, in order to measure the accuracy of a computed solution  $\mathbf{x} \in \text{SOL}(M, \mathbb{K}^n, \mathbf{q})$ , we define the following three relative errors

$$\chi_{\text{rel}_1} = \frac{\max\{\|\mathbf{x}_{(2:n)}\|_2 - x_1, 0\}}{\|\mathbf{x}\|_2}, \quad (7.3a)$$

$$\chi_{\text{rel}_2} = \frac{\max\{\|\mathbf{g}_{(2:n)}\|_2 - g_1, 0\}}{\|M\|_1 \|\mathbf{x}\|_2 + \|\mathbf{q}\|_2}, \quad (7.3b)$$

$$\chi_{\text{rel}_3} = \frac{|\mathbf{x}^\top \mathbf{g}|}{\|\mathbf{x}\|_2 (\|M\|_1 \|\mathbf{x}\|_2 + \|\mathbf{q}\|_2)}, \quad (7.3c)$$

and the total relative error

$$\chi_{\text{rel}} = \chi_{\text{rel}_1} + \chi_{\text{rel}_2} + \chi_{\text{rel}_3}, \quad (7.4)$$

<sup>6</sup>The MATLAB package of SDPT3 is available at [www.math.nus.edu.sg/~matttohc/sdpt3.html](http://www.math.nus.edu.sg/~matttohc/sdpt3.html) and that of SeDuMi is available at [sedumi.ie.lehigh.edu/downloads](http://sedumi.ie.lehigh.edu/downloads).



where  $\mathbf{g} = M\mathbf{x} + \mathbf{q}$ . Our rationale in deciding the normalizing factors in (7.3) is the following. Let  $\text{fl}(\cdot)$  denote the computed result of an expression and  $\mathbf{u}$  be the machine's unit roundoff. Then for exact  $\mathbf{x}$  and  $\mathbf{q}$ , we have  $\text{fl}(\mathbf{g}) = M\mathbf{x} + \mathbf{q} + O(\mathbf{u}[\|M\|_1\|\mathbf{x}\|_2 + \|\mathbf{q}\|_2])$  and thus

$$\begin{aligned}\text{fl}(\|\mathbf{x}_{(2:n)}\|_2 - x_1) &= \|\mathbf{x}_{(2:n)}\|_2 - x_1 + O(\mathbf{u}\|\mathbf{x}\|_2), \\ \text{fl}(\|\mathbf{g}_{(2:n)}\|_2 - g_1) &= \|\mathbf{g}_{(2:n)}\|_2 - g_1 + O(\mathbf{u}[\|M\|_1\|\mathbf{x}\|_2 + \|\mathbf{q}\|_2]), \\ \text{fl}(\mathbf{x}^\top \mathbf{g}) &= \mathbf{x}^\top \mathbf{g} + O(\mathbf{u}\|\mathbf{x}\|_2[\|M\|_1\|\mathbf{x}\|_2 + \|\mathbf{q}\|_2]).\end{aligned}$$

## 7.1 Numerical tests for dense $M$

This subsection is dedicated to numerical tests for medium scale dense problems. Following [44], in this set of tests, we generate pairs  $(M = \widetilde{M}^\top \widetilde{M}, \mathbf{q})$ , where entries of  $\mathbf{q}$  are uniformly distributed in  $[-1, 1]$  and

$$\widetilde{M} = \text{diag}\left(1, \sqrt{1 + \delta}, \sqrt{1 + 2\delta}, \dots, \sqrt{1 + (n-1)\delta}\right) Q,$$

where  $Q = \text{orth}(\text{randn}(n, n))$  and  $\delta = \frac{\text{cond}}{n}$ . The way we construct  $M$  makes it convenient to

1. control the condition number of  $M$ , which is  $(1 - \frac{1}{n})\text{cond} + 1 \approx \text{cond}$ , and
2. easily formulate equivalent problems that SDPT3 and SeDuMi apply, as both of them work on  $\widetilde{M}$ , not  $M = \widetilde{M}^\top \widetilde{M}$ .

To evaluate the performance of the methods, for each  $n$ , we vary  $\text{cond}$  from 10 to  $10^5$ , and set  $\epsilon = 10^{-7}$  as the stopping criteria at line 9 of Algorithm 6.2, and choose  $k_{\max} = 3$  and  $\text{iter}_{\text{Arnoldi}} = [30, 20, 10]$  meaning that the calls to Algorithm 6.1 at line 6 of Algorithm 6.2 are with  $\ell = 30, 20, 10$  as  $k = 0, 1, 2$ , respectively.

For every given pair  $(n, \delta)$ , we generate 5 random SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) together with randomly generated initial points  $\mathbf{x}_0$  for BN, SDPT3 and SeDuMi. Feeding these problems into the four methods<sup>7</sup>: BN, SDPT3, SeDuMi, and LCPvA, in Table 7.2 we report average numbers of iterations ( $\# \text{ iter}$ ), CPU times (measured by MATLAB function `cputime`) and relative errors  $\chi_{\text{rel}}$  defined by (7.3) for each method, over the 5 random problems for each pair  $(n, \delta)$ . In particular, we mention that  $\# \text{ iter}$  for LCPvA is the number of total Arnoldi steps in all calls to Algorithm 6.1 for each run of LCPvA, and the number of the iterations for BN is listed as  $\text{iter}_{\text{bisection}}/\text{iter}_{\text{Newton}}$ , where  $\text{iter}_{\text{bisection}}$  is the number of the bisection steps and  $\text{iter}_{\text{Newton}}$  is that of the Newton steps.

In applying BN and LCPvA, we first perform upper Hessenberg reduction on  $M$  as explained in subsection 6.5.

Table 7.2 clearly shows the effectiveness of LCPvA for dense  $M$  in considering both speed and relative accuracy. But BN comes out the best for  $n = 3000$  and bigger, although it is comparable to LCPvA for  $n = 1000$ . But we note that being an interpretative programming language, MATLAB will consume much more time in executing the **for**-loop in

<sup>7</sup>For SDPT3 and SeDuMi, the inputs are  $\widetilde{M}$ ,  $\mathbf{q}$ , and  $\mathbf{x}_0$ .

Table 7.2: Average numbers of iterations, CPU times and relative errors for dense  $M$

	$n$	1000			3000			5000		
	cond	10	$10^3$	$10^5$	10	$10^3$	$10^5$	10	$10^3$	$10^5$
# iter	BN	10.0/2.8	16.4/1.8	23.2/1.2	16.4/2.2	5.8/5.4	8.6/6.2	9.8/3.6	17.2/2.0	23.5/1.2
	SDPT3	18.0	16.0	15.2	16.0	16.0	18.0	16.0	18.0	20.2
	SeDuMi	13.8	12.8	15.8	13.8	13.8	17.2	15.0	13.6	16.4
	LCPvA	30.0	30.0	30.0	30.0	30.0	30.0	30.0	30.0	30.0
CPU(s)	BN	2.1	2.0	2.2	33.2	33.5	33.1	127.9	128.5	131.1
	SDPT3	43.2	38.8	41.5	979.6	1014.7	1135.3	5144.9	5751.8	6039.5
	SeDuMi	11.9	10.8	15.3	312.1	287.5	360.4	1521.5	1349.0	1646.8
	LCPvA	1.8	1.8	1.7	35.5	35.4	35.3	161.0	160.4	161.5
$\chi_{\text{rel}1}$	BN	$6.0e-17$	$4.2e-16$	0	$6.5e-17$	$1.7e-16$	0	$5.9e-16$	$3.0e-16$	0
	SDPT3	0	0	0	0	0	0	0	0	0
	SeDuMi	0	0	0	0	0	0	0	0	0
	LCPvA	$1.8e-14$	$2.7e-12$	$1.2e-12$	$9.2e-14$	$2.4e-10$	$6.3e-12$	$2.4e-11$	$6.8e-12$	0
$\chi_{\text{rel}2}$	BN	$1.5e-18$	$1.1e-17$	0	$2.0e-17$	$4.1e-17$	0	$3.2e-17$	$1.8e-17$	0
	SDPT3	$4.1e-07$	$3.4e-07$	0	0	$5.7e-07$	$2.7e-06$	0	$4.4e-07$	$3.3e-06$
	SeDuMi	$2.4e-06$	0	0	0	0	0	$3.7e-06$	0	0
	LCPvA	$1.0e-15$	$1.2e-13$	$6.0e-14$	$3.2e-15$	$6.9e-12$	$1.6e-13$	$7.0e-13$	$1.5e-13$	0
$\chi_{\text{rel}3}$	BN	$5.0e-15$	$1.4e-13$	$1.1e-13$	$3.0e-16$	$2.2e-15$	$7.3e-13$	$7.6e-16$	$1.8e-15$	$7.4e-12$
	SDPT3	$2.9e-07$	$2.4e-07$	$6.1e-06$	$9.8e-06$	$4.0e-07$	$1.9e-06$	$1.8e-05$	$3.1e-07$	$2.4e-06$
	SeDuMi	$2.4e-06$	$1.9e-06$	$1.3e-06$	$2.7e-06$	$1.2e-06$	$4.5e-06$	$2.6e-06$	$1.7e-06$	$1.0e-05$
	LCPvA	$2.3e-15$	$2.1e-12$	$7.5e-08$	$2.2e-14$	$9.8e-12$	$1.1e-08$	$1.0e-12$	$3.3e-13$	$2.6e-10$
$\chi_{\text{rel}}$	BN	$5.1e-15$	$1.4e-13$	$1.1e-13$	$3.9e-16$	$2.4e-15$	$7.3e-13$	$1.3e-15$	$2.1e-15$	$7.4e-12$
	SDPT3	$7.1e-07$	$5.9e-07$	$6.1e-06$	$9.9e-06$	$9.7e-07$	$4.6e-06$	$1.8e-05$	$7.5e-07$	$5.7e-06$
	SeDuMi	$4.9e-06$	$1.9e-06$	$1.3e-06$	$2.8e-06$	$1.2e-06$	$4.5e-06$	$6.4e-06$	$1.7e-06$	$1.0e-05$
	LCPvA	$2.2e-14$	$4.9e-12$	$7.5e-08$	$1.2e-13$	$2.6e-10$	$1.1e-08$	$2.6e-11$	$7.3e-12$	$2.6e-10$

Algorithm 6.1 than it does when coded in the language C and FORTRAN. This may be the reason that LCPvA runs slower than BN here for  $n = 3000$  and bigger. In any case, LCPvA is far more efficient than SDPT3 and SeDuMi.

## 7.2 Numerical tests for sparse $M$

For sparse problems, we can afford to treat  $M$  that are larger in size than those in the previous subsection because they are sparse.

Our investigation is carried out on two testing sets. The test matrices  $M$  in the first set are randomly generated while the matrices in the second are from the Matrix Market<sup>8</sup>. The same stopping criteria as in the dense problems are used for SDPT3 and SeDuMi, while for LCPvA, we relax the stopping criterion  $\epsilon$  to  $\epsilon = 10^{-5}$  because with it, LCPvA already yields solutions with comparable relative errors to or much smaller relative errors than what SDPT3 and SeDuMi give. For the eigenpair  $(\tau, \mathbf{v})$  that BN needs, we use the MATLAB function `eigs` with the default options. This alone makes BN much less

<sup>8</sup><http://math.nist.gov/MatrixMarket/index.html>

competitive on the first random testing set. Due to the destruction to the sparsity of  $M$ , we do not transform  $M$  to an upper Hessenberg matrix as we have done in the dense case.

In our first testing set, we choose to first generate the sparse  $\widetilde{M}$  with a prescribed sparsity and condition number, and then form  $M = \widetilde{M}^\top \widetilde{M}$ . It is clear that the number of nonzero entries in  $M$  generated from this procedure is usually much larger than that of  $\widetilde{M}$ . Limiting by our current hardware environment, we choose  $n = 10000$  and generate  $(M = \widetilde{M}^\top \widetilde{M}, \mathbf{q})$ , where  $\mathbf{q}$  is randomly generated with elements uniformly distributed in the interval  $[-1, 1]$  and  $\widetilde{M}$  is generated by the MATLAB function

$$\widetilde{M} = \text{sprandsym}(n, \text{density}, \text{rc}, \text{kind}).$$

It returns an  $n$ -by- $n$  symmetric and positive definite random matrix with approximately the prescribed `density` and an approximate condition number  $\frac{1}{\text{rc}}$ , implying that the condition number `cond` of  $M$  is approximately  $(\frac{1}{\text{rc}})^2$ . The value of the input `kind` can be either 1 or 2, where the former means that  $\widetilde{M}$  is generated by random Jacobi rotations on a positive definite diagonal matrix, while the latter indicates that  $\widetilde{M}$  is a shifted sum of outer products. Both cases are tested with `density` = 0.0005 and different `rc` to make `cond` vary from  $10^2$  to  $10^5$ .

For each prescribed pair `(rc, kind)`, we generate 5 triples  $(M = \widetilde{M}^\top \widetilde{M}, \mathbf{q}, \mathbf{x}_0)$ , where  $\mathbf{x}_0$  is the initial point for the three algorithms other than LCPvA, and record the numerical results by the four methods in Table 7.3. Also reported are the average sparse densities  $d_M$  and  $d_{\widetilde{M}}$  of  $M$  and  $\widetilde{M}$  for each pair `(rc, kind)`, respectively.

We have several observations from Table 7.3:

- (1) The numerical results show that each algorithm performs quite differently for the case `kind=1` and `kind=2`;
- (2) For `kind=1`, it turns out that LCPvA has the best numerical performance in terms of speed and accuracy. Note that there are almost twice as many number of nonzeros in  $M$  as in  $\widetilde{M}$ . It is also observed that BN performs the worst due to the requirement in finding an accurate eigenpair  $(\tau, \mathbf{v})$ ;
- (3) For `kind=2`, SDPT3 is the best in time; but its solutions are often less accurate than the ones by LCPvA.

Our second set of testing matrices  $M$  are from the Matrix Market through searching all real sparse symmetric, positive definite and non-diagonal matrices with  $5000 \leq n \leq 90000$ . We found nine matrices, namely S1RMQ4M1, S1RMT3M1, S2RMQ4M1, S2RMT3M1, S3RMQ4M1, S3RMT3M1, S3RMT3M3, BCSSTK17 and BCSSTK18, whose basic properties are listed at the top of Table 7.4. The corresponding  $\widetilde{M}$  for SDPT3 and SeDuMi is computed by  $\widetilde{M} = \text{chol}(M)$  (the CPU time for this is not counted in), and we also report  $d_M$  and  $d_{\widetilde{M}}$ , the sparse densities of  $M$  and  $\widetilde{M}$ , respectively. To form a specific SOCLCP (1.1), we take  $\mathbf{q} = [1, 1, \dots, 1]^\top \in \mathbb{R}^n$  and generate a random  $\mathbf{x}_0$  as the starting vector for BN, SDPT3 and SeDuMi.

Table 7.4 summarizes the numerical results for the second test set, where ‘-’ for BN means that it fails in `eigs` for solving the eigenpair  $(\tau, \mathbf{v})$ . We observe that LCPvA has a

Table 7.3: Numerical results for sparse problems in the first testing set

	$(\frac{1}{rc})^2 (\approx \text{cond})$	kind=1			kind=2		
		10 <sup>2</sup>	10 <sup>4</sup>	10 <sup>5</sup>	10 <sup>2</sup>	10 <sup>4</sup>	10 <sup>5</sup>
		$\widetilde{d}_M$	$d_M$	$\widetilde{d}_M$	$d_M$	$\widetilde{d}_M$	$d_M$
# iter	BN	7.4/5.4	6.2/6.6	9.0/10.0	13.4/2.6	5.8/5.4	8.6/6.2
	SDPT3	21.2	21.0	20.0	16.0	20.0	20.6
	SeDuMi	17.4	18.6	18.0	14.4	17.8	18.4
	LCPvA	30.0	42.0	30.0	30.0	52.0	46.0
CPU(s)	BN	225.9	278.9	420.6	269.7	301.3	371.3
	SDPT3	1.5	1.1	1.8	52.6	64.8	67.2
	SeDuMi	2.9	3.3	3.7	824.7	1002.9	1067.1
	LCPvA	1.3	2.5	1.5	374.2	528.3	716.2
$\chi_{\text{rel}1}$	BN	1.4e-15	5.6e-01	1.4e-00	5.8e-16	9.1e-16	2.8e-01
	SDPT3	0	0	0	0	0	0
	SeDuMi	0	0	0	0	0	0
	LCPvA	4.2e-13	1.4e-14	0	0	1.6e-11	2.1e-12
$\chi_{\text{rel}2}$	BN	4.9e-16	1.1e-02	5.9e-01	1.6e-17	1.1e-18	3.9e-06
	SDPT3	2.7e-05	5.8e-04	6.6e-05	2.7e-06	1.0e-05	7.3e-06
	SeDuMi	0	3.4e-07	0	3.8e-07	2.1e-06	1.2e-06
	LCPvA	1.7e-13	8.6e-06	0	0	8.9e-15	2.4e-13
$\chi_{\text{rel}3}$	BN	6.0e-16	1.5e-03	6.7e-16	3.1e-16	1.3e-18	1.1e-15
	SDPT3	2.0e-05	4.1e-04	4.7e-05	1.9e-06	7.4e-06	5.1e-06
	SeDuMi	1.1e-06	3.2e-06	2.0e-06	1.0e-06	2.8e-06	1.1e-06
	LCPvA	1.7e-13	9.9e-06	9.7e-10	2.4e-10	2.5e-10	4.7e-06
$\chi_{\text{rel}}$	BN	2.5e-15	5.8e-01	2.0e-00	9.0e-16	9.2e-16	2.8e-01
	SDPT3	4.7e-05	9.9e-04	1.1e-04	4.6e-06	1.7e-05	1.2e-05
	SeDuMi	1.1e-06	3.6e-06	2.0e-06	1.4e-06	4.9e-06	2.3e-06
	LCPvA	5.4e-12	1.8e-05	9.7e-10	2.4e-10	2.7e-10	4.7e-06

very good performance in all the cases except for BCSSTK18. The failure is due to that the choice  $s_0 \approx 1.0252 \times 10^{10}$  by Table 7.1 is too far from the true  $s_* \approx 3.8768 \times 10^3$ ; as a result,  $h_\ell(s)$  provides a very poor approximation of  $h(s)$  near  $s_*$  and the next shift satisfying (6.8) cannot be found. We experimented with using different initial shifts and were able to make LCPvA run successfully with  $s_0 = \|\mathbf{q}\|_2/5$  to give  $\chi_{\text{rel}1} = \chi_{\text{rel}2} = 0$  and  $\chi_{\text{rel}3} = 1.9 \times 10^{-11}$  in 6.15 seconds. Although the  $s_0$  given in Table 7.1 has been working so well until BCSSTK18, this suggests that there are cases for which more effective  $s_0$  is needed for LCPvA to succeed. We shall investigate this issue in our future study.

Finally, it is interesting to note that for BCSSTK17, LCPvA uses only 20 Arnoldi steps which implies that a friendly breakdown occurred and, thereby, result in  $h_{20}(s) \equiv h(s)$ ; consequently,  $s_*$  is an exact zero of  $h_{20}(s)$ . This is partially revealed by the tiny relative error  $\chi_{\text{rel}} = 8.9e-14$ .

## 8 Concluding remarks

In this paper, we have proposed a Krylov subspace method for the linear complementarity problem over a single second-order cone. The method is in the spirit of projection and

Table 7.4: Numerical results for sparse problems in the second testing set

	$M$	s1RMQ4M1	s1RMT3M1	s2RMQ4M1	s2RMT3M1	s3RMQ4M1	s3RMT3M1	s3RMT3M3	BCSSTK17	BCSSTK18
	$n$	5489	5489	5489	5489	5489	5489	5357	10974	11948
	cond	$3.2e+06$	$5.4e+06$	$1.2e+08$	$5.8e+08$	$1.3e+10$	$1.3e+10$	$2.6e+10$	$6.5e+01$	$6.5e+01$
	2-norm	$6.9e+05$	$9.7e+05$	$6.9e+04$	$9.7e+04$	$6.9e+03$	$9.7e+03$	$9.6e+03$	$1.3e+10$	$4.3e+10$
	$d_M$	$8.7e-03$	$7.2e-03$	$8.7e-03$	$7.2e-03$	$8.7e-03$	$7.2e-03$	$7.2e-03$	$3.6e-03$	$1.0e-03$
	$d_{\tilde{M}}$	$3.3e-02$	$3.2e-02$	$3.3e-02$	$3.3e-02$	$3.3e-02$	$3.3e-02$	$7.7e-02$	$1.3e-02$	$2.0e-02$
# iter	BN	25/1	25/1	21/1	22/2	18/2	18/2	18/2	1/10	—
	SDPT3	15	16	14	13	12	18	14	18	32
	SeDuMi	20	21	17	18	16	17	17	14	20
	LCPvA	30	30	30	30	30	30	30	20	30
CPU(s)	BN	16.6	13.3	14.8	20.6	21.6	20.4	18.4	588.3	—
	SDPT3	11.4	12.1	10.7	10.1	9.2	13.7	247.4	25.9	199.3
	SeDuMi	36.7	34.9	22.7	28.6	16.1	18.8	102.4	25.5	151.1
	LCPvA	9.7	5.8	9.6	6.0	9.1	6.2	5.2	6.3	6.7
$\chi_{rel1}$	BN	0	0	0	0	0	0	0	$1.4e-00$	—
	SDPT3	0	0	0	0	0	0	0	0	0
	SeDuMi	0	0	0	0	0	0	0	0	0
	LCPvA	0	0	$2.4e-09$	0	$5.9e-10$	0	$1.2e-09$	$8.9e-14$	$1.0e-00$
$\chi_{rel2}$	BN	0	0	0	0	0	0	0	$6.6e-11$	—
	SDPT3	$1.7e-05$	$8.8e-06$	$7.7e-07$	$1.7e-05$	$6.5e-06$	0	0	0	0
	SeDuMi	0	0	$2.1e-08$	$1.4e-07$	$2.0e-06$	$9.1e-07$	$9.3e-07$	0	0
	LCPvA	0	0	$1.3e-10$	0	$3.5e-11$	0	$5.9e-11$	$4.2e-24$	$7.6e-23$
$\chi_{rel3}$	BN	$1.1e-13$	$2.4e-13$	$1.8e-12$	$7.5e-17$	$1.4e-16$	$2.6e-17$	$1.9e-17$	$2.3e-25$	—
	SDPT3	$1.2e-05$	$6.2e-06$	$5.5e-07$	$1.2e-05$	$4.6e-06$	$1.2e-08$	$4.0e-06$	$4.6e-07$	$8.1e-08$
	SeDuMi	$3.9e-07$	$4.3e-07$	$1.5e-08$	$9.8e-08$	$1.4e-06$	$6.5e-07$	$6.6e-07$	$1.4e-07$	$2.6e-08$
	LCPvA	$3.0e-09$	$2.7e-05$	$1.9e-10$	$1.1e-10$	$4.9e-11$	$1.8e-09$	$8.3e-11$	$6.6e-24$	$1.0e-27$
$\chi_{rel}$	BN	$1.1e-13$	$2.4e-13$	$1.8e-12$	$7.5e-17$	$1.4e-16$	$2.6e-17$	$1.9e-17$	$1.4e-00$	—
	SDPT3	$2.8e-05$	$1.5e-05$	$1.3e-06$	$2.9e-05$	$1.1e-05$	$1.2e-08$	$4.0e-06$	$4.6e-07$	$8.1e-08$
	SeDuMi	$3.9e-07$	$4.3e-07$	$3.6e-08$	$2.4e-07$	$3.5e-06$	$1.6e-06$	$1.6e-06$	$1.4e-07$	$2.6e-08$
	LCPvA	$3.0e-09$	$2.7e-05$	$2.8e-09$	$1.1e-10$	$6.8e-10$	$1.8e-09$	$1.3e-09$	$8.9e-14$	$1.0e-00$

mimics the well-known Rayleigh-Ritz procedure [11, 15] for the large scale eigenvalue problems. It is known that the Rayleigh-Ritz procedure first seeks a good approximate subspace to an invariant subspace and then projects the eigenvalue problem to a much smaller one which is then solved to give approximate eigenpairs. The Krylov subspace technique is the key in the first phase. When it comes to SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ), it is not that evident that what would constitute a good approximate subspace that could be used as a projection subspace to yield a much smaller size SOCLCP (1.1). In this paper, through the special rational function  $h(s)$  and techniques in the model reduction, we have successfully defined good approximate subspaces upon which the original large scale SOCLCP( $M, \mathbb{K}^n, \mathbf{q}$ ) is projected to yield much smaller scale problems. Their accuracy in approximating the original problem can be controlled by the number  $\ell$  of Arnoldi steps and measured by  $|h(s) - h_\ell(s)|$  which is of order  $O(s^\ell)$ . Basing on this idea, we presented an algorithmic framework: LCPvA (Linear Complementarity Problem via Arnold) in Algorithm 6.2 which

was put into test against existing sophisticated solvers. Our preliminary numerical results demonstrate its efficiency.

Detailed theoretical analysis of the curve  $h(s)$  is also presented. It is of interest both in theory and in computations. The success of our proposed method here in large part is due to our complete knowledge on how  $h(s)$  behaves. That knowledge makes it possible for us to select the right one among all zeros of the reduced  $h_\ell(s)$  to approximate the desired one of  $h(s)$ .

We have identified two problems that warrant our further investigation. One is to iteratively solve the linear systems arising in the Arnoldi process in Algorithm 6.1 for extremely large scale SOCLCP. Our current implementation use MATLAB's sparse LU for simplicity. The other is to find better initial  $s_0$ . So far the one given in Table 7.1 works pretty well for us, except for one example in BCSSTK18 for which we found a better initial  $s_0$  is needed for convergence.

Our major assumption for  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  is that  $M$  has the GUS property upon which all our developments are built. What happens when  $M$  does not have the GUS property? This will be yet another problem that warrants our attention in the future.

We believe that the idea and technique we presented in this paper can be adapted to efficiently solve other critical optimization problems such as the trust-region subproblem [30, 31] and for more general linear complementarity problems such as (1.2). We will investigate these problems, too.

## A Convert SOCLCP (1.1) to the format for SDPT3 and SeDuMi

In this appendix, we outlines how to convert  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  to the format used in SDPT3 and SeDuMi. For this purpose, it suffices to rewrite  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  into the (primal) mathematical programming in the form of

$$\text{(Primal)} \quad \min_{A\mathbf{y}=\mathbf{b}, \mathbf{y} \in \mathcal{K}} \mathbf{c}^\top \mathbf{y}, \quad (\text{A.1})$$

and its dual form

$$\text{(Dual)} \quad \max_{A^\top \mathbf{z} + \mathbf{s} = \mathbf{c}, \mathbf{s} \in \check{\mathcal{K}}} \mathbf{b}^\top \mathbf{z}, \quad (\text{A.2})$$

where  $\mathcal{K}$  is certain cone and  $\check{\mathcal{K}}$  stands for the dual cone of  $\mathcal{K}$ .

To this end, we first note that when  $M$  is symmetric and positive definite, solving  $\text{SOCLCP}(M, \mathbb{K}^n, \mathbf{q})$  is equivalent to the following quadratic second-order cone programming:

$$\min_{\mathbf{x} \in \mathbb{K}^n} \frac{1}{2} \mathbf{x}^\top M \mathbf{x} + \mathbf{q}^\top \mathbf{x}. \quad (\text{A.3})$$

Now if  $M = \widetilde{M}^\top \widetilde{M}$ , according to [39, 41], we can rewrite (A.3) equivalently as

$$\min_{\xi, \mathbf{x}} \frac{1}{2} \xi + \mathbf{q}^\top \mathbf{x} \quad \text{subject to} \quad \|\widetilde{M} \mathbf{x}\|_2^2 \leq \xi, \quad \mathbf{x} \in \mathbb{K}^n, \quad \xi \in \mathbb{R}_+. \quad (\text{A.4})$$

By introducing a new variable  $\mathbf{t} \in \mathbb{K}^{n+2}$ , (A.4) is equivalent to

$$\begin{aligned} & \min_{\xi, \mathbf{x}, \mathbf{t}} \frac{1}{2} \xi + \mathbf{q}^\top \mathbf{x} \\ \text{subject to} & \begin{bmatrix} -\frac{1}{2} & 0 \\ -\frac{1}{2} & 0 \\ 0 & \widetilde{M} \end{bmatrix} \begin{bmatrix} \xi \\ \mathbf{x} \end{bmatrix} + \mathbf{t} = \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \\ 0 \end{bmatrix}, \quad \mathbf{t} \in \mathbb{K}^{n+2}, \mathbf{x} \in \mathbb{K}^n, \xi \in \mathbb{R}_+, \end{aligned}$$

which is in the form of (A.1) with

$$A = \left[ \begin{array}{cc|c} -\frac{1}{2} & 0 & \\ -\frac{1}{2} & 0 & I_{n+2} \\ 0 & \widetilde{M} & \end{array} \right], \quad \mathbf{y} = \begin{bmatrix} \xi \\ \mathbf{x} \\ \mathbf{t} \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} \frac{1}{2} \\ \mathbf{q} \\ 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \\ \mathbf{0} \end{bmatrix},$$

and  $\mathcal{K} = \mathbb{R}_+ \times \mathbb{K}^n \times \mathbb{K}^{n+2} = \check{\mathcal{K}}$ . In our numerical experiments, we employ **SeDuMi** to solve the primal problem (A.1) while use **SDPT3** to solve the dual problem (A.2).

## References

- [1] A. Auslender. Variational inequalities over the cone of semidefinite positive symmetric matrices and over the Lorentz cone. *Opt. Methods Soft.*, 18:359–376, 2003.
- [2] Z. Bai and Y. Su. Dimension reduction of large-scale second-order dynamical systems via a second-order Arnoldi method. *SIAM J. Sci. Comput.*, 25(5):1692–1709, 2005.
- [3] Z. Bai and Y. Su. SOAR: A second-order Arnoldi method for the solution of the quadratic eigenvalue problem. *SIAM J. Matrix Anal. Appl.*, 26(3):640–659, 2005.
- [4] J. Burke and S. Xu. The global linear convergence of a noninterior path-following algorithm for linear complementarity problems. *Math. Oper. Res.*, 23:719–734, 1998.
- [5] J.-S. Chen and S. H. Pan. A descent method for a reformulation of the second-order cone complementarity problem. *J. Comput. Appl. Math.*, 213:547–558, 2008.
- [6] J.-S. Chen and P. Tseng. An unconstrained smooth minimization reformulation of the second-order cone complementarity problem. *Math. Programming*, 104:293–327, 2005.
- [7] X. Chen, L. Qi, and D. F. Sun. Global and superlinear convergence of the smoothing Newton method and its application to general box constrained variational inequalities. *Math. Comp.*, 67:519–540, 1998.
- [8] X. D. Chen, D. F. Sun, and J. Sun. Complementarity functions and numerical experiments on some smoothing Newton methods for second-order-cone complementarity problems. *Comput. Opt. Appl.*, 25:39–56, 2003.
- [9] E. Chiprout and M.S. Nakhla. *Asymptotic Waveform Evaluation and Moment Matching for Interconnect Analysis*. Kluwer Academic Publishers, 1994.
- [10] R. W. Cottle, J.-S. Pang, and R. E. Stone. *The Linear Complementarity Problem*. Computer Science and Scientific Computing, Academic Press, Inc., Boston, MA, 1992.
- [11] J. Demmel. *Applied Numerical Linear Algebra*. SIAM, Philadelphia, PA, 1997.

- [12] P. Feldman and R. W. Freund. Efficient linear circuit analysis by Padé approximation via the Lanczos process. *IEEE Trans. Computer-Aided Design*, 14:639–649, 1995.
- [13] M. Fukushima, Z.-Q. Luo, and P. Tseng. Smoothing functions for second-order-cone complementarity problems. *SIAM J. Optim.*, 12:436–460, 2002.
- [14] K. Gallivan, E. Grimme, and P. Van Dooren. Asymptotic waveform evaluation via a Lanczos method. *Appl. Math. Lett.*, 7(5):75–80, 1994.
- [15] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, Maryland, 3rd edition, 1996.
- [16] M. S. Gowda and R. Sznajder. Automorphism invariance of P- and GUS-properties of linear transformations on Euclidean Jordan algebras. *Math. Oper. Res.*, 31:109123, 2006.
- [17] M. S. Gowda and R. Sznajder. Some global uniqueness and solvability results for linear complementarity problems over symmetric cones. *SIAM J. Optim.*, 18:461–481, 2007.
- [18] M. S. Gowda, R. Sznajder, and J. Tao. Some P-properties for linear transformations on Euclidean Jordan algebras. *Linear Algebra Appl.*, 393:203–232, 2004.
- [19] E. J. Grimme. *Krylov Projection Methods For Model Reduction*. PhD thesis, University of Illinois at Urbana-Champaign, Urbana, Illinois, 1997.
- [20] S. Hayashi, N. Yamashita, and M. Fukushima. A combined smoothing and regularization method for monotone second-order cone complementarity problems. *SIAM J. Optim.*, 15:593–615, 2005.
- [21] T. Heyn, M. Anitescu, A. Tasora, and D. Negrut. Using Krylov subspace and spectral methods for solving complementarity problems in many-body contact dynamics simulation. *Internat. J. Numer. Methods Eng.*, 95:541–561, 2013.
- [22] Z. H. Huang and T. Ni. Smoothing algorithms for complementarity problems over symmetric cones. *Comput. Opt. Appl.*, 45:557–579, 2010.
- [23] C. Kanzow. Some noninterior continuation methods for linear complementarity problems. *SIAM J. Matrix Anal. Appl.*, 17:851–868, 1996.
- [24] M. Kojima, S. Shindoh, and S. Hara. Interior-point methods for the monotone semidefinite linear complementarity problem in symmetric matrices. *SIAM J. Optim.*, 7:86–125, 1997.
- [25] L. C. Kong, J. Sun, and N. H. Xiu. A regularized smoothing Newton method for symmetric cone complementarity problems. *SIAM J. Optim.*, 19:1028–1047, 2008.
- [26] R.-C. Li and Z. Bai. Structure-preserving model reduction using a Krylov subspace projection formulation. *Comm. Math. Sci.*, 3(2):179–199, 2005.
- [27] R.-C. Li and Q. Ye. Simultaneous similarity reductions for a pair of matrices to condensed forms. *Comm. Math. Stat.*, 2014. to appear.
- [28] X. Liang, R.-C. Li, and Z. Bai. Trace minimization principles for positive semi-definite pencils. *Linear Algebra Appl.*, 438:3085–3106, 2013.
- [29] J. L. Morales, J. Nocedal, and M. Smelyanskiy. An algorithm for the fast solution of symmetric linear complementarity problems. *Numer. Math.*, 111:251–266, 2008.
- [30] J. Moré and D. Sorensen. Computing a trust region step. *SIAM J. Sci. Statist. Comput.*, 4(3):553–572, 1983.



- [31] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.
- [32] S. Pan and J.-S. Chen. A damped Gauss-Newton method for the second-order cone complementarity problem. *Appl. Math. Opt.*, 59:293–318, 2009.
- [33] J.-S. Pang, D. F. Sun, and J. Sun. Semismooth homeomorphisms and strong stability of semidefinite and Lorentz complementarity problems. *Math. Oper. Res.*, 28:39–63, 2003.
- [34] B. N. Parlett. *The Symmetric Eigenvalue Problem*. SIAM, Philadelphia, 1998.
- [35] B. N. Parlett, D. R. Taylor, and Z. A. Liu. A look-ahead Lanczos algorithm for unsymmetric matrices. *Math. Comp.*, 44(169):105–124, 1985.
- [36] W. H. A. Schilders, Henk A. van der Vorst, and J. Rommes (editors). *Model Order Reduction: Theory, Research Aspects and Applications*. Springer, Boston, 2008.
- [37] J. F. Sturm. Using SeDuMi 1.02, a Matlab toolbox for optimization over symmetric cones. *Opt. Methods Soft.*, 11:625–653, 1999.
- [38] T.-J. Su and Jr. R. R. Craig. Model reduction and control of flexible structures using Krylov vectors. *J. Guidance, Control, and Dynamics*, 14:260–267, 1991.
- [39] K. C. Toh, M. J. Todd, and R. H. Tütüncü. SDPT3—a MATLAB software package for semidefinite programming. *Opt. Methods Soft.*, 11:545–581, 1999.
- [40] K. C. Toh, M. J. Todd, and R. H. Tütüncü. On the implementation and usage of SDPT3—a MATLAB software package for semidefinite-quadratic-linear programming, version 4.0. Technical report, Department of Mathematics, National University of Singapore, 2010.
- [41] R. H. Tütüncü, K. C. Toh, and M. J. Todd. Solving semidefinite-quadratic-linear programs using SDPT3. *Math. Programming*, 95:189–217, 2003.
- [42] W. H. Yang and X. M. Yuan. The GUS-property of second-order cone linear complementarity problems. *Math. Programming*, 141:295–317, 2013.
- [43] L.-H. Zhang and W. H. Yang. An efficient algorithm for second-order cone linear complementarity problems. *Math. Comp.*, 83:1701–1726, 2013.
- [44] L.-H. Zhang and W. H. Yang. An efficient matrix splitting method for the second-order cone complementarity problem. *SIAM J. Optim.*, 24(3):1178–1205, 2014.