

# Error Analysis of the symplectic Lanczos Method for the symplectic Eigenvalue Problem

Heike Faßbender \*

December 29, 1999

## Abstract

An error analysis of the symplectic Lanczos algorithm for the symplectic eigenvalue problem in finite-precision arithmetic is given, if no breakdown occurs. The analysis shows that the restriction of preserving the symplectic structure does not destroy the characteristic feature of nonsymmetric Lanczos processes. An analog of Paige's theory on the relationship between the loss of orthogonality among the Lanczos vectors and the convergence of Ritz values in the symmetric Lanczos algorithm is discussed. As to be expected, it follows that (under certain assumptions) the computed  $J$ -orthogonal Lanczos vectors lose  $J$ -orthogonality when some Ritz values begin to converge.

**Key words :** symplectic Lanczos method, symplectic matrix, eigenvalues, error analysis.

**AMS(MOS) subject classifications :** 65G05, 65F15, 65F50

## 1 Introduction

The Lanczos algorithm proposed by Lanczos in 1950 [11] is a procedure for the successive reduction of a given general matrix  $A \in \mathbb{R}^{n \times n}$  to a tridiagonal matrix  $T$ . In the  $j$ th step the Lanczos algorithm generates two  $n \times j$  matrices  $Q_j$  and  $P_j$

$$Q_j = [q_1, q_2, \dots, q_j], \quad P_j = [p_1, p_2, \dots, p_j]$$

which satisfy

$$P_j^T Q_j = I$$

and

$$\begin{aligned} (1) \quad A Q_j &= Q_j T_j + \beta_{j+1} q_{j+1} e_j^T, \\ (2) \quad A^T P_j &= P_j T_j^T + \gamma_{j+1} p_{j+1} e_j^T, \end{aligned}$$

where  $e_j = [0, \dots, 0, 1] \in \mathbb{R}^j$  and  $T_j$  is the tridiagonal matrix

$$T_j = \begin{bmatrix} \alpha_1 & \gamma_2 & & & \\ \beta_2 & \alpha_2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & \beta_j & \alpha_j \end{bmatrix}.$$

---

\*Universität Bremen, Fachbereich 3 - Mathematik und Informatik, Zentrum für Technomathematik, 28334 Bremen, FRG. E-mail: heike@math.uni-bremen.de

The eigenvalues of the intermediate tridiagonal matrices  $T_j$  of smaller dimension typically approximate some of the eigenvalues of  $A$  (often the ones largest in magnitude). During the iteration the matrix  $A$  is referenced only through matrix-vector products  $Ax$  and  $A^T x$ ; hence the algorithm is useful for finding a few eigenvalues of very large and sparse matrices. A wide range of Lanczos papers appeared since the 1960s, see, e.g., the references in [7].

Recently, there has been considerable interest in structure-preserving Lanczos algorithms for the symplectic eigenproblem. These eigenproblems arise in applications like the problem of solving algebraic Riccati equations or  $H_\infty$ -norm computations. See, e.g. [10, 12, 17]. In some of these applications the symplectic matrix is very large and sparse, and only a few eigenvalues and the corresponding invariant subspace are desired.

A structure-preserving Lanczos-like method for the symplectic eigenproblem was first proposed by Banse [2]. The symplectic matrix is reduced to a symplectic butterfly matrix. Banse presents a look-ahead version of the method which overcomes breakdown by giving up the strict butterfly form. Benner and Faßbender [3, 4] suggest to combine the idea of the symplectic Lanczos method with the idea of implicitly restarted Lanczos methods in order to deal with the numerical difficulties inherent to any nonsymmetric Lanczos-like method.

Here we give an error analysis of the symplectic Lanczos method for the symplectic eigenproblem. Numerical experiments show that, just like in the conventional Lanczos algorithm, information about the extreme eigenvalues tends to emerge long before the symplectic Lanczos process is completed. The effect of finite-precision arithmetic is discussed. Using Bai's work [1] on the nonsymmetric Lanczos algorithm, an analog of Paige's theory [13] on the relationship between the loss of orthogonality among the computed Lanczos vectors and the convergence of a Ritz value is discussed. The symplectic Lanczos algorithm is reviewed in Section 2. Stopping criteria are discussed. In Section 3 a rounding error analysis of the symplectic Lanczos algorithm in finite-precision arithmetic is presented. Section 4 discusses convergence of the symplectic Lanczos algorithm versus the loss of  $J$ -orthogonality of the computed Lanczos vectors. Numerical experiments are presented in Section 5. All proofs are deferred to the Appendix, due to their highly technical nature.

## 2 The Symplectic Lanczos Algorithm

A matrix  $M \in \mathbb{R}^{2n \times 2n}$  is called *symplectic* if

$$(3) \quad M J^{2n,2n} M^T = J^{2n,2n}$$

(or equivalently,  $M^T J^{2n,2n} M = J^{2n,2n}$ ), where

$$(4) \quad J^{2n,2n} = \begin{bmatrix} 0 & I^{n,n} \\ -I^{n,n} & 0 \end{bmatrix},$$

and  $I^{n,n}$  is the  $n \times n$  identity matrix. If the dimension of  $I^{n,n}$ , or  $J^{2n,2n}$ , is clear from the context, we leave off the superscript. We denote by  $Z^{2k}$  the first  $2k$  columns of a  $2n \times 2n$  matrix  $Z$ .

The symplectic matrices form a group under multiplication. The eigenvalues of symplectic matrices occur in reciprocal pairs: if  $\lambda$  is an eigenvalue of  $M$  with right eigenvector  $x$ , then  $\lambda^{-1}$  is an eigenvalue of  $M$  with left eigenvector  $(Jx)^T$ .

In exact arithmetic and without breakdown, the symplectic Lanczos methods proposed by Banse [2] and Benner and Faßbender [4] reduce  $M$  to a symplectic butterfly matrix. A symplectic matrix

$$B = \begin{bmatrix} B_1 & B_2 \\ B_3 & B_4 \end{bmatrix} = \begin{bmatrix} \diagdown & \equiv \\ \diagup & \equiv \end{bmatrix}$$

is called a *butterfly matrix* if  $B_1, B_3 \in \mathbb{R}^{n \times n}$  are diagonal matrices and  $B_2, B_4 \in \mathbb{R}^{n \times n}$  are tridiagonal matrices. An *unreduced butterfly matrix* is one for which the tridiagonal matrix  $B_4$  is

unreduced, see [4, 5]. Using the definition of a symplectic matrix, one easily verifies that if  $B$  is unreduced, then the diagonal submatrix  $B_3$  is nonsingular. This allows the parameterization of  $B$  in the following form (see [4, 5]):

$$B = (K^{2n,2n})^{-1} N^{2n,2n}$$

$$= \left[ \begin{array}{c|c} a_1^{-1} & b_1 \\ & \ddots \\ & a_n^{-1} & b_n \\ \hline & a_1 & & \\ & & \ddots & \\ & & & a_n \end{array} \right] \left[ \begin{array}{c|cccc} & & & & -1 \\ & & & & \ddots \\ & & & & \ddots \\ & & & & \ddots & -1 \\ \hline 1 & & & c_1 & d_2 & \\ & \ddots & & d_2 & \ddots & \ddots \\ & & \ddots & & \ddots & \ddots & d_n \\ & & & & & & d_n & c_n \\ & & & & & & & 1 \end{array} \right].$$

Given  $s_1 \in \mathbb{R}^{2n}$  and a symplectic matrix  $M \in \mathbb{R}^{2n \times 2n}$  the symplectic Lanczos algorithm generates a sequence of symplectic butterfly matrices  $B^{2k,2k} \in \mathbb{R}^{2k \times 2k}$  such that (if no breakdown occurs)

$$(5) \quad MS^{2k} = S^{2k} B^{2k,2k} + r_{k+1} e_{2k}^T, \quad k = 1, 2, \dots, n,$$

where  $S^{2k} \in \mathbb{R}^{2n \times 2k}$ ,  $S^{2k} e_1 = s_1$ , and the columns of  $S^{2k}$  are orthogonal with respect to the indefinite inner product defined by  $J$  as in (4). That is, the columns of  $S^{2k}$  are  $J$ -orthogonal. The eigenvalues of the intermediate matrices  $B^{2k,2k}$  are progressively better estimates of  $M$ 's eigenvalues. For  $k = n$  the algorithm computes a symplectic matrix  $S$  such that  $S$  transforms  $M$  into butterfly form:  $S^{-1}MS = B$ .

In order to simplify the notation we use in the following permuted versions of  $M$ ,  $B$ , and  $S$ . Let

$$Z_P := PZP^T$$

with the permutation matrix

$$P := [e_1, e_3, \dots, e_{2n-1}, e_2, e_4, \dots, e_{2n}] \in \mathbb{R}^{2n \times 2n}.$$

Using the permutation matrix  $P$ , the matrix  $J$  can be permuted to the  $2n \times 2n$  block diagonal matrix

$$J_P := PJP^T = \text{diag}\left(\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}\right).$$

$M_P$ ,  $B_P$ , and  $S_P$  are permuted symplectic matrices, in other words, they are  $J_P$ -orthogonal.

Using the permuted versions of  $M_P$ ,  $B_P$ , and  $S_P$ , the structure preserving Lanczos method generates a sequence of permuted symplectic matrices

$$S_P^{2k} := [v_1, w_1, v_2, w_2, \dots, v_k, w_k] \in \mathbb{R}^{2n \times 2k}$$

satisfying

$$(6) \quad M_P S_P^{2k} = S_P^{2k} B_P^{2k,2k} + d_{k+1} (b_{k+1} v_{k+1} + a_{k+1} w_{k+1}) e_{2k}^T.$$

The symplectic Lanczos algorithm for symplectic matrices is summarized in Table 1. For a derivation of the algorithm and a detailed discussion of various aspects see [3, 4, 5]. There is some freedom in the choice of the parameters that occur in the algorithm. Essentially, the parameters  $b_k$  can be chosen freely. Here we set  $b_k = 1$ . A different choice of the parameters  $a_k$  and

Algorithm : Symplectic Lanczos method	
Choose an initial vector $\tilde{v}_1 \in \mathbb{R}^{2n}, \tilde{v}_1 \neq 0$ .	
Set $v_0 = 0 \in \mathbb{R}^{2n}$ .	
Set $d_1 = \ \tilde{v}_1\ _2$ and $v_1 = \frac{1}{d_1}\tilde{v}_1$ .	
for $m = 1, 2, \dots$ do	
(update of $w_m$ )	
set	
$\tilde{w}_m = M_P v_m - b_m v_m$	
$a_m = v_m^T J_P M_P v_m$	
$w_m = \frac{1}{a_m} \tilde{w}_m$	
(computation of $c_m$ )	
$c_m = -a_m^{-1} w_m^T J_P M_P^{-1} v_m$	
(update of $v_{m+1}$ )	
$\tilde{v}_{m+1} = -d_m v_{m-1} - c_m v_m + w_m + a_m^{-1} M_P^{-1} v_m$	
$d_{m+1} = \ \tilde{v}_{m+1}\ _2$	
$v_{m+1} = \frac{1}{d_{m+1}} \tilde{v}_{m+1}$	

Table 1: Symplectic Lanczos Method for the Symplectic Eigenproblem

$d_k$  is possible. Note that  $M_P^{-1} = -J_P M_P^T J_P$ , since  $M$  is symplectic. Thus  $M_P^{-1} v_m$  is just a matrix-vector-product with the transpose of  $M_P$ .

Equivalent to (6), as  $B_P^{2k,2k} = (K_P^{2k,2k})^{-1} N_P^{2k,2k}$  and  $e_{2k}^T (N_P^{2k,2k})^{-1} = -e_{2k-1}^T$ , we have

$$(7) \quad M_P S_P^{2k} (N_P^{2k,2k})^{-1} = S_P^{2k} (K_P^{2k,2k})^{-1} - d_{k+1} (b_{k+1} v_{k+1} + a_{k+1} w_{k+1}) e_{2k-1}^T.$$

The vector  $r_{k+1} := d_{k+1} (b_{k+1} v_{k+1} + a_{k+1} w_{k+1})$  is the *residual vector* and is  $J_P$ -orthogonal to the columns of  $S_P^{2k}$ , the *Lanczos vectors*. The matrix  $B_P^{2k,2k}$  is the  $J_P$ -orthogonal projection of  $M_P$  onto the range of  $S_P^{2k}$

$$B_P^{2k,2k} = J_P^{2k,2k} (S_P^{2k})^T J_P M_P S_P^{2k}.$$

**Remark 2.1** *The usual nonsymmetric Lanczos algorithm generates two sequences of vectors  $\{q_j\}$  and  $\{p_j\}$  (see (1) and (2)). Adapting the usual nonsymmetric Lanczos algorithm to the situation considered here, the symplectic Lanczos process could have been stated as follows: Given  $v_1, t_1 \in \mathbb{R}^{2n}$  and a symplectic matrix  $M \in \mathbb{R}^{2n \times 2n}$ , the symplectic Lanczos algorithm produces matrices  $S_P^{2k} = [v_1, w_1, \dots, v_k, w_k] \in \mathbb{R}^{2n \times 2k}$  and  $W_P^{2k} = [t_1, \dots, t_{2k}] \in \mathbb{R}^{2n \times 2k}$  with  $J_P$ -orthogonal columns which satisfy*

$$(W_P^{2k})^T S_P^{2k} = I^{2k,2k},$$

and

$$\begin{aligned} M_P S_P^{2k} &= S_P^{2k} B_P^{2k,2k} + d_{k+1} r_{k+1} e_{2k}^T, \\ M_P^T W_P^{2k} &= W_P^{2k} (B_P^{2k,2k})^T + d_{k+1} \check{r}_{k+1} e_{2k}^T. \end{aligned}$$

As  $S_P$  is symplectic, we obtain from  $(W_P^{2k})^T S_P^{2k} = I^{2k,2k}$  that

$$W_P^{2k} = J_P^{2n,2n} S_P^{2k} J_P^{2k,2k} = [-J_P w_1, J_P v_1, \dots, -J_P w_k, J_P v_k].$$

Moreover,

$$r_{k+1} = M_P v_{k+1}, \quad \text{and} \quad \tilde{r}_{k+1} = J_P v_{k+1}.$$

Substituting the expressions for  $W_P^{2k}$  and  $\tilde{r}_{k+1}$  into the second recursion and pre- and post-multiplying with  $J_P$  yields that the two recursions are equivalent. Hence one of the two sequences can be eliminated here and thus work and storage can essentially be halved. (This property is valid for a broader class of matrices, see [6].)

Assume that we have performed  $k$  steps of the symplectic Lanczos method and thus obtained the identity (after permuting back)

$$MS^{2k} = S^{2k} B^{2k,2k} + d_{k+1}(b_{k+1}\hat{v}_{k+1} + a_{k+1}\hat{w}_{k+1})e_{2k}^T.$$

If the norm of the residual vector is small, the  $2k$  eigenvalues of  $B^{2k,2k}$  are approximations to the eigenvalues of  $M$ . Numerical experiments indicate that the norm of the residual rarely becomes small by itself. Nevertheless, some eigenvalues of  $B^{2k,2k}$  may be good approximations to eigenvalues of  $M$ . Let  $\lambda$  be an eigenvalue of  $B^{2k,2k}$  with the corresponding eigenvector  $y$ . Then the vector  $x = S^{2k}y$  satisfies

$$(8) \quad \begin{aligned} \|Mx - \lambda x\|_2 &= \|(MS^{2k} - S^{2k}B^{2k,2k})y\|_2 \\ &= |d_{k+1}| |e_{2k}^T y| \|b_{k+1}\hat{v}_{k+1} + a_{k+1}\hat{w}_{k+1}\|_2. \end{aligned}$$

The vector  $x$  is referred to as *Ritz vector* and  $\lambda$  as *Ritz value* of  $M$ . If the last component of the eigenvector  $y$  is sufficiently small, the right-hand side of (8) is small and the pair  $\{\lambda, x\}$  is a good approximation to an eigenvalue-eigenvector pair of  $M$ . Note that  $|e_{2k}^T y| > 0$  if  $B^{2k,2k}$  is unreduced (see, e.g., [5, Lemma 3.11]). The pair  $\{\lambda, x\}$  is exact for the nearby problem

$$(M + E)x = \lambda x \quad \text{where} \quad E = -d_{k+1}(b_{k+1}\hat{v}_{k+1} + a_{k+1}\hat{w}_{k+1})e_k^T (S^{2k})^T J^{2n,2n}.$$

In an actual implementation, typically the *Ritz estimate*

$$|d_{k+1}| |e_{2k}^T y| \|b_{k+1}\hat{v}_{k+1} + a_{k+1}\hat{w}_{k+1}\|_2$$

is used in order to decide about the numerical accuracy of an approximate eigenpair. This avoids the explicit formation of the residual  $(MS^{2k} - S^{2k}B^{2k,2k})y$ .

A small Ritz estimate is not sufficient for the Ritz pair  $\{\lambda, x\}$  to be a good approximation to an eigenvalue-eigenvector pair of  $M$ . It does not guarantee that  $\lambda$  is a good approximation to an eigenvalue of  $M$ . That is

$$\min_j |\lambda - \mu_j|, \quad \text{where} \quad \mu_j \in \sigma(M) = \{\mu \in \mathbb{C} \mid \exists x \in \mathbb{R}^{2n} \setminus \{0\} \ni Mx = \mu x\}$$

is not necessarily small when the Ritz estimate is small (see, e.g., [9, Section 3]). For nonnormal matrices the norm of the residual of an approximate eigenvector is not by itself sufficient information to bound the error in the approximate eigenvalue. It is sufficient however to give a bound on the distance to the nearest matrix to which the Ritz triplet  $\{\lambda, x, y\}$  is exact [9] (here  $y$  denotes the left Ritz vector of  $M$  corresponding to the Ritz value  $\lambda$ ). In the following, we will give a computable expression for the error. Assume that  $B^{2k,2k}$  is diagonalizable, i.e., there exists  $Y \in \mathbb{C}^{2k \times 2k}$  such that

$$Y^{-1} B^{2k,2k} Y = \left[ \begin{array}{c|c} \lambda_1 & \\ \hline & \lambda_k \\ \hline & \lambda_1^{-1} \\ & & \ddots \\ & & & \lambda_k^{-1} \end{array} \right] = \Lambda.$$

Let  $X = S^{2k}Y = [x_1, \dots, x_{2k}]$  and denote  $b_{k+1}\widehat{v}_{k+1} + a_{k+1}\widehat{w}_{k+1}$  by  $\widehat{r}_{k+1}$ . Since

$$MS^{2k} = S^{2k}B^{2k,2k} + d_{k+1}\widehat{r}_{k+1}e_{2k}^T,$$

it follows that

$$MS^{2k}Y = S^{2k}YY^{-1}B^{2k,2k}Y + d_{k+1}\widehat{r}_{k+1}e_{2k}^TY,$$

or

$$MX = X\Lambda + d_{k+1}\widehat{r}_{k+1}e_{2k}^TY.$$

Thus

$$Mx_i = \lambda_i x_i + y_{2k,i}d_{k+1}\widehat{r}_{k+1},$$

and

$$Mx_{k+i} = \lambda_i^{-1}x_{k+i} + y_{2k,k+i}d_{k+1}\widehat{r}_{k+1},$$

for  $i = 1, \dots, k$ . The last equation can be rewritten as

$$(Jx_{k+i})^T M = \lambda_i (Jx_{k+i})^T + y_{2k,k+i}\lambda_i d_{k+1}\widehat{r}_{k+1}^T JM.$$

Using Theorem 2' of [9] we obtain that  $\{\lambda_i, x_i, (Jx_{k+i})^T\}$  is an eigen-triplet of  $M - F_{\lambda_i}$  where

$$\|F_{\lambda_i}\|_2 = |d_{k+1}| \max_i \left\{ \frac{|y_{2k,i}| \|\widehat{r}_{k+1}\|_2}{\|x_i\|_2}, \frac{|y_{2k,k+i}\lambda_i| \|\widehat{r}_{k+1}^T JM\|_2}{\|Jx_{k+i}\|_2} \right\}.$$

Furthermore, if  $\|F_{\lambda_i}\|_2$  is small enough, then

$$|\theta_i - \lambda_j| \leq \text{cond}(\lambda_j) \|F_{\lambda_i}\|_2 + \mathcal{O}(\|F_{\lambda_i}\|_2^2),$$

where  $\theta_i$  is an eigenvalue of  $M$  and  $\text{cond}(\lambda_j)$  is the condition number of the Ritz value  $\lambda_j$

$$\text{cond}(\lambda_j) = \frac{\|x_i\|_2 \|Jx_{k+i}\|_2}{|x_{k+i}^T Jx_i|} = \|x_i\|_2 \|x_{k+i}\|_2.$$

Similarly, we obtain that  $\{\lambda_i^{-1}, x_{k+i}, (Jx_i)^T\}$  is an eigen-triplet of  $M - F_{\lambda_i^{-1}}$  where

$$\|F_{\lambda_i^{-1}}\|_2 = |d_{k+1}| \max_i \left\{ \frac{|y_{2k,k+i}| \|\widehat{r}_{k+1}\|_2}{\|x_{k+i}\|_2}, \frac{|y_{2k,i}\lambda_i^{-1}| \|\widehat{r}_{k+1}^T JM\|_2}{\|Jx_i\|_2} \right\}.$$

Consequently, as  $\lambda_i$  and  $\lambda_i^{-1}$  should be treated alike, the symplectic Lanczos algorithm should be continued until  $\|F_{\lambda_i}\|_2$  and  $\|F_{\lambda_i^{-1}}\|_2$  are small, and until  $\text{cond}(\lambda_j) \|F_{\lambda_i}\|_2$  and  $\text{cond}(\lambda_j) \|F_{\lambda_i^{-1}}\|_2$  are below a given threshold for accuracy. Note that as in the Ritz estimate, in the criteria derived here, the essential quantities are  $|d_{k+1}|$  and the last component of the desired eigenvectors  $|y_{2k,i}|$  and  $|y_{2k,k+i}|$ .

### 3 The Symplectic Lanczos Algorithm in Finite-Precision Arithmetic

In this section, we present a rounding error analysis of the symplectic Lanczos algorithm in finite-precision arithmetic. Our analysis will follow the lines of Bai's analysis of the nonsymmetric Lanczos algorithm [1]. It is in the spirit of Paige's analysis for the symmetric Lanczos algorithm

[13], except that we (as Bai) carry out the analysis component-wise rather than norm-wise. The component-wise analysis allows to measure each element of a perturbation relative to its individual tolerance, so that, unlike in the norm-wise analysis, the sparsity pattern of the problem under consideration can be exploited.

We use the usual model of floating-point arithmetic, as, e.g., in [7, 8]:

$$fl(x \circ y) = (x \circ y)(1 + \varepsilon)$$

where  $\circ$  denotes any of the four basic arithmetic operations  $+$ ,  $-$ ,  $*$ ,  $/$  and  $|\varepsilon| \leq \mathbf{u}$  with  $\mathbf{u}$  denoting the *unit roundoff*.

We summarize (as in [1]) all the results for basic linear algebra operations of sparse vectors and/or matrices that we need for our analysis:

*Saxpy operation:*

$$fl(\alpha x + y) = \alpha x + y + e, \quad |e| \leq \mathbf{u} (2|\alpha x| + |y|) + \mathcal{O}(\mathbf{u}^2),$$

*Inner product:*

$$fl(x^T y) = x^T y + e, \quad |e| \leq k\mathbf{u} |x|^T |y| + \mathcal{O}(\mathbf{u}^2),$$

*Matrix-vector multiplication:*

$$fl(Ax) = Ax + e, \quad |e| \leq m\mathbf{u} |A| |x| + \mathcal{O}(\mathbf{u}^2),$$

where  $k$  is the number of overlapping nonzero components in the vectors  $x$  and  $y$ , and  $m$  is the maximal number of nonzero elements of the matrix  $A$  in any row. For a vector  $x = [x_1, \dots, x_n]^T$ ,  $|x|$  denotes the vector  $[|x_1|, \dots, |x_n|]^T$ . Similar, for a matrix  $A = [a_{ij}]_{i,j=1}^n$ ,  $|A|$  denotes the  $n \times n$  matrix  $[|a_{ij}|]_{i,j=1}^n$ .

In this section, any computed quantity will be denoted by a hat, e.g.,  $\hat{\alpha}$  will denote a computed quantity that is affected by rounding errors. (Please note, that in the previous section, we used hatted quantities to denote the non-permuted symplectic Lanczos vectors.)

Analyzing one step of the symplectic Lanczos algorithm to see the effects of the finite-precision arithmetic we obtain the following theorem.

**Theorem 3.1** *Let  $M \in \mathbb{R}^{2n \times 2n}$  be a symplectic matrix with at most  $m$  nonzero entries in any row or column. If no breakdown occurs during the execution of  $k$  steps of the symplectic Lanczos algorithm as given in Table 1, the computed Lanczos vectors satisfy*

$$(9) \quad \hat{a}_j \hat{w}_j = M_P \hat{v}_j - \hat{v}_j + h_j,$$

$$(10) \quad \hat{d}_{j+1} \hat{v}_{j+1} = -\hat{d}_j \hat{v}_{j-1} - \hat{c}_j \hat{v}_j + \hat{w}_j + \hat{a}_j^{-1} M_P^{-1} \hat{v}_j + g_{j+1},$$

where

$$(11) \quad |h_j| \leq (m+2)\mathbf{u} |M_P| |\hat{v}_j| + 2\mathbf{u} |\hat{v}_j| + \mathcal{O}(\mathbf{u}^2),$$

$$(12) \quad |g_{j+1}| \leq (m+5)\mathbf{u} |\hat{a}_j^{-1}| |M_P^{-1}| |\hat{v}_j| + 4\mathbf{u} |\hat{w}_j| + 4\mathbf{u} |\hat{c}_j \hat{v}_j| + 3\mathbf{u} |\hat{d}_j \hat{v}_{j-1}| + \mathcal{O}(\mathbf{u}^2).$$

The computed matrices  $\hat{S}_P^{2k}$ ,  $\hat{N}_P^{2k,2k}$ , and  $\hat{K}_P^{2k,2k}$  satisfy

$$(13) \quad M_P \hat{S}_P^{2k} (\hat{N}_P^{2k,2k})^{-1} = \hat{S}_P^{2k} (\hat{K}_P^{2k,2k})^{-1} - \hat{d}_{k+1} M_P \hat{v}_{k+1} e_{2k-1}^T + E_k,$$

where

$$(14) \quad \begin{aligned} \|E_k\|_F &\leq \mathbf{u} \|\hat{S}^{2k}\|_F \left[ (m+5) \|\hat{K}^{2k,2k}\|_F \|M\|_F^2 + 4 \|\hat{N}^{2k,2k}\|_F \|M\|_F \right. \\ &\quad \left. + (m+6) \|M\|_F + 2 \|\hat{K}^{2k,2k}\|_F \right] + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

**Proof:** See Appendix. ✓

This indicates that the recursion equation  $M_P S_P^{2k} (N_P^{2k,2k})^{-1} = S_P^{2k} (K_P^{2k,2k})^{-1} - d_{k+1} M_P v_{k+1} e_{2k-1}^T$  is satisfied to working precision, if  $\|M\|_F^2, \|\widehat{S}^{2k}\|_F, \|\widehat{K}^{2k,2k}\|_F$ , and  $\|\widehat{N}^{2k,2k}\|_F \|M\|_F$  are of moderate size. But, unfortunately,  $\|\widehat{S}^{2k}\|_F$  may grow unboundedly in the case of near breakdown.

While the equation (9) is given by the  $(2j)$ th column of  $M_P S_P N_P^{-1} = S_P K_P^{-1}$ , the equation (10) corresponds to the  $(2j-1)$ th column of  $S_P N_P^{-1} = M_P^{-1} S_P K_P^{-1}$ . The upper bounds associated with (9) and (10) involve only  $\|M\|_F$  as to be expected, see (11) and (12). Recall that  $M_P^{-1} = -J_P M_P^T J_P$ , since  $M$  is symplectic. Thus  $|M_P^{-1}|$  does not introduce any problems usually involved by forming the inverse of a matrix. In order to summarize these results into one single equation, we define

$$(15) \quad E_k = [M_P g_2, -h_1, M_P g_3, -h_2, \dots, M_P g_{k+1}, -h_k].$$

Then (13) holds. Using the component-wise upper bounds for  $|h_j|$  and  $|g_{j+1}|$ , we obtain the upper bound for  $E_k$  as given in (14). As we summarize our results in terms of the equation  $M_P S_P N_P^{-1} = S_P K_P^{-1}$ , we have to pre-multiply the error bound associated with (10) by  $M_P$ , resulting in an artificial  $\|M\|_F^2$  term here. Hence combining all our findings into one single equation forces the  $\|M\|_F^2$  term.

For the nonsymmetric Lanczos algorithm, Bai obtains a similar result in [1]. The equations corresponding to our equations (9) and (10) are (see (1) and (2))

$$\begin{aligned} \widehat{\beta}_{j+1} \widehat{q}_{j+1} &= A \widehat{q}_j - \widehat{\alpha}_j \widehat{q}_j - \widehat{\gamma}_j \widehat{q}_j + h_j^{nonsymLan}, \\ \widehat{\gamma}_{j+1} \widehat{p}_{j+1} &= A^T \widehat{p}_j - \widehat{\alpha}_j \widehat{p}_j - \widehat{\beta}_j \widehat{p}_{j-1} + g_{j+1}^{nonsymLan}. \end{aligned}$$

The errors associated are given by

$$\begin{aligned} |h_j^{nonsymLan}| &\leq (3+m)\mathbf{u} |A| |\widehat{q}_j| + 4\mathbf{u} |\widehat{\alpha}_j| |\widehat{q}_j| + 3\mathbf{u} |\widehat{\gamma}_j| |\widehat{q}_{j-1}| + \mathcal{O}(\mathbf{u}^2), \\ |g_{j+1}^{nonsymLan}| &\leq (3+m)\mathbf{u} |A| |\widehat{p}_j| + 4\mathbf{u} |\widehat{\alpha}_j| |\widehat{p}_j| + 3\mathbf{u} |\widehat{\gamma}_j| |\widehat{p}_{j-1}| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

Hence, the symplectic Lanczos algorithms behaves essentially like the nonsymmetric Lanczos algorithm. The additional restriction of preserving the symplectic structure does not pose any additional problems concerning the rounding error analysis, the results of the analysis are essentially the same.

**Corollary 3.2** *In Remark 2.1 we have noted that the usual nonsymmetric Lanczos algorithm generates two sequences of vectors, but that due to the symplectic structure, the two recurrence relations of the standard nonsymmetric Lanczos algorithm are equivalent for the situation discussed here. It was noted that the equation which is not used is given by*

$$M_P^T W_P^{2k} (K_P^{2k,2k})^T = W_P^{2k} (N_P^{2k,2k})^T + d_{k+1} J_P v_{k+1} e_{2k}^T,$$

where

$$W_P^{2k} = J_P^{2n,2n} S_P^{2k} J_P^{2k,2k} = [-J_P w_1, J_P v_1, \dots, -J_P w_k, J_P v_k].$$

Instead of summarizing our findings into equation (13), we could have summarized

$$(16) \quad M_P^T \widehat{W}_P^{2k} (\widehat{K}_P^{2k,2k})^T = \widehat{W}_P^{2k} (\widehat{N}_P^{2k,2k})^T + \widehat{d}_{k+1} J_P \widehat{v}_{k+1} e_{2k}^T + F_k$$

where

$$(17) \quad \begin{aligned} \widehat{W}_P^{2k} &= J_P^{2n,2n} \widehat{S}_P^{2k} J_P^{2k,2k}, \\ F_k &= [M_P^T J_P h_1, J_P g_2, \dots, M_P^T J_P h_k, J_P g_{k+1}]. \end{aligned}$$

As an upper bound for  $\|F_k\|_F$  we obtain

$$(18) \quad \|F_k\|_F \leq \mathbf{u} \|\widehat{S}^{2k}\|_F \left[ (m+2) \|M\|_F^2 + (m+7) \|\widehat{K}^{2k,2k}\|_F \|M\|_F + 4 \|\widehat{N}^{2k,2k}\|_F + 4 \right] + \mathcal{O}(\mathbf{u}^2).$$

As before, the term  $\|M\|_F^2$  is introduced because we summarize all our findings into one single equation.

It is well-known, that in finite-precision arithmetic, orthogonality between the computed Lanczos vectors in the symmetric Lanczos process is lost. This loss of orthogonality is due to cancellation and is not the result of the gradual accumulation of roundoff error (see, e.g., [15, 16]). What can we say about the  $J$ -orthogonality of the computed symplectic Lanczos vectors? Obviously, rounding errors, once introduced into some computed Lanczos vectors, are propagated to future steps. Such error propagation for the nonsymmetric Lanczos process is analyzed by Bai [1].

Let us take a closer look at the  $J_P$ -orthogonality of the computed symplectic Lanczos vectors. Define

$$K = [\hat{v}_1, \hat{w}_1, \dots, \hat{v}_k, \hat{w}_k]^T J_P [\hat{v}_1, \hat{w}_1, \dots, \hat{v}_k, \hat{w}_k].$$

That is,

$$\begin{aligned} k_{2j-1,2m-1} &= \hat{v}_j^T J_P \hat{v}_m, & k_{2j-1,2m} &= \hat{v}_j^T J_P \hat{w}_m, \\ k_{2j,2m-1} &= \hat{w}_j^T J_P \hat{v}_m, & k_{2j,2m} &= \hat{w}_j^T J_P \hat{w}_m. \end{aligned}$$

In exact arithmetic we would have  $K = J_P$ , where  $J_P$  is block diagonal; each diagonal block is of the form  $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ . As  $x^T J_P x = 0$  for any vector  $x$ , we have

$$k_{2j,2j} = k_{2j-1,2j-1} = 0,$$

not depending on the loss of  $J_P$ -orthogonality between the computed symplectic Lanczos vectors. Examining the other elements of  $K$  we obtain the following lemma.

**Lemma 3.3** *The elements  $k_{jm}$  of  $K$  satisfy the following equations*

$$(19) \quad \begin{aligned} k_{jj} &= 0 & j &= 1, \dots, 2k \\ -k_{j,j+1} &= k_{j+1,j} = -1 + \kappa_j & j &= 1, \dots, 2k-1 \\ k_{jm} &= \mathcal{O}(\mathbf{u}) & j, m &= 1, \dots, 2k, \\ & & m &\neq j-1, j, j+1 \end{aligned}$$

where

$$(20) \quad |\kappa_j| \leq \mathbf{u} \frac{|\hat{v}_j|^T |J_P| \{2(m+n+2) |M_P| + 5\} |\hat{v}_j|}{|\hat{w}_j^T J_P \hat{v}_j|} + \mathcal{O}(\mathbf{u}^2).$$

**Proof:** See Appendix. ✓

Lemma 3.3 describes how  $J$ -orthogonality between the computed symplectic Lanczos vectors is lost. Our findings will be useful in the following section when discussing the question of loss of  $J$ -orthogonality versus convergence of a Ritz pair.

## 4 Convergence versus Loss of $J$ -Orthogonality

It is well-known that in the symmetric Lanczos procedure, loss of orthogonality between the computed Lanczos vectors implies convergence of a Ritz pair to an eigenpair, see, e.g., [14]. Here we will discuss the situation for the symplectic Lanczos algorithm, following the lines of Section 4 of Bai's analysis of the nonsymmetric Lanczos algorithm in [1]. We will see that a conclusion similar to the one for the symmetric Lanczos process holds here, subject to a certain condition.

From the previous section, we know that the computed symplectic Lanczos vectors obey the following equalities:

$$(21) \quad M_P \hat{S}_P^{2k} = \hat{S}_P^{2k} \hat{B}_P^{2k,2k} - \left[ \hat{d}_{k+1} \hat{r}_{k+1} e_{2k-1}^T - E_k \right] \hat{N}_P^{2k,2k},$$

$$(22) \quad M_P^T \hat{W}_P^{2k} = \hat{W}_P^{2k} (\hat{B}_P^{2k,2k})_P^T + \left[ \hat{d}_{k+1} J_P \hat{v}_{k+1} e_{2k}^T + F_k \right] (\hat{K}_P^{2k,2k})^{-T},$$

with

$$(23) \quad (\hat{S}_P^{2k})^T J_P^{2n,2n} \hat{S}_P^{2k} = K = J_P^{2k,2k} + C_k + \Delta_k - C_k^T,$$

where  $\widehat{B}_P^{2k,2k} = (\widehat{K}_P^{2k,2k})^{-1} \widehat{N}_P^{2k,2k}$ ,  $\widehat{W}_P^{2k} = J_P^{2n,2n} \widehat{S}_P^{2k} J_P^{2k,2k}$ , the rounding error matrices  $E_k$  and  $F_k$  are as in (15) and, (17), resp.,  $\Delta_k$  is a block diagonal matrix with  $2 \times 2$  block on the diagonal,

$$\Delta_k = \text{diag}\left(\begin{bmatrix} 0 & \kappa_1 \\ -\kappa_1 & 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 & \kappa_k \\ -\kappa_k & 0 \end{bmatrix}\right),$$

and  $C_k$  is a strictly lower block triangular matrix with block size 2. That is  $(C_k)_{\ell,j} = 0$  for  $\ell = 1, \dots, 2k, j = \ell, \dots, 2k$ , and  $(C_k)_{2\ell,2\ell-1} = 0$  for  $\ell = 1, \dots, k$ .

To simplify our discussion, we make two assumptions, which are also used in the analysis of the symmetric Lanczos process [15, p. 265] and in the analysis of the nonsymmetric Lanczos process [1]. The first assumption is *local  $J$ -orthogonality*, that is, the computed symplectic Lanczos vectors are  $J$ -orthogonal to their neighboring Lanczos vectors:

$$(24) \quad \begin{bmatrix} \widehat{v}_{j-1}^T \\ \widehat{w}_{j-1}^T \end{bmatrix} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} \widehat{v}_j \\ \widehat{w}_j \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

This implies that the  $2 \times 2$  block on the sub-diagonal of  $C_k$  are zero, yielding the following block-structure

$$C_k = \begin{bmatrix} 0 & 0 & 0 & \cdots & \cdots & 0 & 0 & 0 \\ 0 & 0 & 0 & \cdots & \cdots & 0 & 0 & 0 \\ X & 0 & 0 & \cdots & \cdots & 0 & 0 & 0 \\ X & X & 0 & \cdots & \cdots & 0 & 0 & 0 \\ X & X & X & \ddots & & 0 & 0 & 0 \\ & \vdots & & \ddots & \ddots & & \vdots & \\ X & X & X & \cdots & X & 0 & 0 & 0 \\ X & X & X & \cdots & X & X & 0 & 0 \end{bmatrix},$$

where the  $X$  denote  $2 \times 2$  blocks.

The second assumption is that the eigenvalue problem for the  $2k \times 2k$  butterfly matrix  $\widehat{B}_P^{2k,2k} = (\widehat{K}_P^{2k,2k})^{-1} \widehat{N}_P^{2k,2k}$  is solved exactly, that is,

$$(25) \quad Y_k^{-1} \widehat{B}_P^{2k,2k} Y_k = \text{diag}(\lambda_1, \lambda_1^{-1}, \dots, \lambda_k, \lambda_k^{-1}).$$

This implies that the computed Ritz vector for  $\lambda_j$  is given by

$$z_j = \widehat{S}_P^{2k} y_{2j-1},$$

while the computed Ritz vector for  $\lambda_j^{-1}$  is given by

$$x_j = \widehat{S}_P^{2k} y_{2j}.$$

**Theorem 4.1** *Assume that the symplectic Lanczos algorithm in finite-precision arithmetic satisfies (21) – (23). Let*

$$\begin{aligned} L_k^{(2)} + U_k^{(2)} &= J_P^{2k,2k} \Delta_k \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} \Delta_k, \\ L_k^{(4)} + U_k^{(4)} &= (\widehat{K}_P^{2k,2k})^{-1} F_k^T \widehat{S}_P^{2k} - (\widehat{W}_P^{2k})^T E_k \widehat{N}_P^{2k,2k}, \end{aligned}$$

where  $L_k^{(2)}$  and  $L_k^{(4)}$  are strictly lower block triangular matrices, and  $U_k^{(2)}$  and  $U_k^{(4)}$  are strictly upper block triangular matrices with block size 2. Then the computed Ritz vectors  $x_j = \widehat{S}_P^{2k} y_{2j}$  and  $z_j = \widehat{S}_P^{2k} y_{2j-1}$  satisfy

$$(26) \quad x_j^T J_P^{2n,2n} \widehat{r}_{k+1} = \frac{y_{2j}^T J_P^{2k,2k} [U_k^{(4)} - U_k^{(2)}] y_{2j-1}}{\widehat{d}_{k+1} (e_{2k}^T y_{2j-1})} =: \frac{\psi_1}{\widehat{d}_{k+1} (e_{2k}^T y_{2j-1})},$$

$$(27) \quad z_j^T J_P^{2n,2n} \widehat{r}_{k+1} = \frac{y_{2j-1}^T J_P^{2k,2k} [U_k^{(4)} - U_k^{(2)}] y_{2j}}{\widehat{d}_{k+1} (e_{2k}^T y_{2j})} =: \frac{\psi_2}{\widehat{d}_{k+1} (e_{2k}^T y_{2j})}.$$

**Proof:** See Appendix. ✓

The derived equations are similar to those obtained by Bai for the nonsymmetric Lanczos process. Hence we can interpret our findings analogously: equations (26) and (27) describe the way in which the  $J$ -orthogonality is lost. Recall that the scalar  $d_{k+1}$  and the last eigenvector components  $(e_{2k}^T y_{2j-1})$  and  $(e_{2k}^T y_{2j})$  are also essential quantities used as the backward error criteria for the computed Ritz triplets  $\{\lambda_i, z_i, (Jx_i)^T\}$  and  $\{\lambda_i^{-1}, x_i, (Jz_i)^T\}$  discussed in Section 2. (Also recall that  $|e_{2k}^T y_\ell| > 0$  if  $B^{2k,2k}$  is unreduced.) Hence, if the quantities  $|\psi_1|$  and  $|\psi_2|$  are bounded and bounded away from zero, then (26) and (27) reflect the reciprocal relation between the convergence of the symplectic Lanczos process (i.e., tiny  $\hat{d}_{k+1}(e_{2k}^T y_{2j-1})$  and  $\hat{d}_{k+1}(e_{2k}^T y_{2j})$ ) and the loss of  $J$ -orthogonality (i.e., large  $\hat{r}_{k+1}^T J_P x_j$  and  $\hat{r}_{k+1}^T J_P z_j$ ).

Let us conclude our analysis by estimating  $|\psi_1|$  and  $|\psi_2|$ . Let us assume (again analogous to Bai's analysis) that  $\Delta_k = 0$ , i.e.,  $\hat{w}_j^T J_P \hat{v}_j = -1$ , which simplifies the technical details of the analysis and appears to be the case in practice, up to the order of machine precision. Under this assumption, we have  $U_k^{(2)} = 0$ . Moreover, we have

$$|\psi_j| \leq \|U_k^{(4)}\|_F \|y_{2j}\|_2 \|y_{2j-1}\|_2,$$

for  $j = 1, 2$ . Let us derive an estimate for  $\|U_k^{(4)}\|_F$ .  $U_k^{(4)}$  is the strictly upper block triangular part of

$$(28) \quad (\hat{K}_P^{2k,2k})^{-1} F_k^T \hat{S}_P^{2k} - (\hat{W}_P^{2k})^T E_k \hat{N}_P^{2k,2k}.$$

A generous upper bound is therefore given by

$$\begin{aligned} \|U_k^{(4)}\|_F &\leq \|\hat{K}^{2k,2k}\|_F \|F_k^T\|_F \|\hat{S}^{2k}\|_F + \|\hat{W}^{2k}\|_F \|E_k\|_F \|\hat{N}^{2k,2k}\|_F \\ &\leq \|\hat{S}^{2k}\|_F \left[ \|\hat{K}^{2k,2k}\|_F \|F_k\|_F + \|E_k\|_F \|\hat{N}^{2k,2k}\|_F \right] \\ &\leq \mathbf{u} \|\hat{S}^{2k}\|_F^2 \left\{ (m+5) \|\hat{K}^{2k,2k}\|_F \|\hat{N}^{2k,2k}\|_F \|M\|_F^2 + 7 \|\hat{K}^{2k,2k}\|_F \|\hat{N}^{2k,2k}\|_F \right. \\ &\quad \left. + 4 \|\hat{K}^{2k,2k}\|_F + (m+2) \|\hat{K}^{2k,2k}\|_F \|M\|_F^2 + (m+8) \|\hat{K}^{2k,2k}\|_F^2 \|M\|_F \right. \\ &\quad \left. + 4 \|\hat{N}^{2k,2k}\|_F^2 \|M\|_F + (m+6) \|\hat{N}^{2k,2k}\|_F \|M\|_F \right\} + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

Summarizing, we obtain the following corollary, which gives an upper bound for  $|\psi_1|$  and  $|\psi_2|$ .

**Corollary 4.2** *Assume that  $\Delta_k = 0$  in Theorem 4.1. Then*

$$\begin{aligned} |\psi| &\leq \mathbf{u} \operatorname{cond}(\lambda_j) \left\{ (m+5) \|\hat{K}^{2k,2k}\|_F \|\hat{N}^{2k,2k}\|_F \|M\|_F^2 + 7 \|\hat{K}^{2k,2k}\|_F \|\hat{N}^{2k,2k}\|_F \right. \\ &\quad \left. + (m+2) \|\hat{K}^{2k,2k}\|_F \|M\|_F^2 + (m+8) \|\hat{K}^{2k,2k}\|_F^2 \|M\|_F + 4 \|\hat{K}^{2k,2k}\|_F \right. \\ &\quad \left. + 4 \|\hat{N}^{2k,2k}\|_F^2 \|M\|_F + (m+6) \|\hat{N}^{2k,2k}\|_F \|M\|_F \right\} + \mathcal{O}(\mathbf{u}^2), \end{aligned}$$

where  $\psi \in \{\psi_1, \psi_2\}$  and

$$\operatorname{cond}(\lambda_j) = \operatorname{cond}(\lambda_j^{-1}) = \|\hat{S}^{2k}\|_F^2 \|y_{2j}\|_2 \|y_{2j-1}\|_2$$

is the condition number of the Ritz values  $\lambda_j$  and  $\lambda_j^{-1}$ .

Note that this bound is too pessimistic. In order to derive an upper bound for  $\|U_k^{(4)}\|_F$ , an upper bound for the matrix (28) is used, as  $U_k^{(4)}$  is the strictly upper block triangular part of that matrix. This is a very generous upper bound for  $\|U_k^{(4)}\|_F$ . Moreover, the term

$$\|\hat{K}^{2k,2k}\|_F \|\hat{N}^{2k,2k}\|_F$$

is an upper bound for the norm of  $\widehat{B}_P^{2k,2k}$ . The squared terms  $\|\widehat{K}^{2k,2k}\|_F^2$  and  $\|\widehat{N}^{2k,2k}\|_F^2$  are introduced as the original equations derived in (13) and (16) are given in terms of  $\widehat{K}^{2k,2k}$  and  $\widehat{N}^{2k,2k}$ , but not in terms of  $\widehat{B}_P^{2k,2k}$ . Numerical examples show that the bound  $|\psi|$  is too pessimistic.

Unfortunately, for the symplectic Lanczos process (as for any nonsymmetric Lanczos-like process), near breakdown may cause the norms of the symplectic Lanczos vectors  $\|\widehat{v}_j\|_2$  and  $\|\widehat{w}_j\|_2$  to grow unboundedly. Accumulating the quantity  $\sum_{j=1}^k (\|\widehat{v}_j\|_2^2 + \|\widehat{w}_j\|_2^2)$ , which costs about  $4nk$  flops, we can obtain a computable bound for  $\text{cond}(\lambda_j)$  and  $\text{cond}(\lambda_j^{-1})$  in practise. Theorem 4.1 and Corollary 4.2 indicate that if the  $J$ -orthogonality between  $\widehat{r}_{k+1}$  and  $x_j$  (and  $z_j$ ) is lost, then the value  $\widehat{d}_{k+1}(e_{2k}^T y_{2j-1})$  is proportional to  $|\psi_1|$  (and the value  $\widehat{d}_{k+1}(e_{2k}^T y_{2j})$  is proportional to  $|\psi_2|$ ). Given the upper bound from Corollary 4.2, and supposing that  $\text{cond}(\lambda_j)$  is reasonably bounded, the loss of  $J$ -orthogonality implies that  $\widehat{d}_{k+1}(e_{2k}^T y_{2j-1})$  (and  $\widehat{d}_{k+1}(e_{2k}^T y_{2j})$ ) are small. Therefore, in the best case we can state that if the effects of finite-precision arithmetic,  $E_k$  and  $F_k$  in (21) and (22), are small, then small residuals tell us that the computed eigenvalues are eigenvalues of matrices close to the given matrix.

## 5 Numerical examples

As the results derived are not surprising we give just one example to demonstrate the practical behavior of the convergence of a Ritz value versus the loss of  $J$ -orthogonality among the symplectic Lanczos vectors. All computations were done using MATLAB<sup>1</sup> Version 5.3 on a Sun Ultra 1 with IEEE double-precision arithmetic and machine precision  $\epsilon = 2.2204 \times 10^{-16}$ .

Our code implements exactly the algorithm as given in Table 1. In order to detect convergence, the rather crude criterion

$$\|r_{k+1}\| \leq \|M\| * 10^{-6}$$

was used. Benign breakdown in the symplectic Lanczos process was detected by the criterion

$$\|\widetilde{v}_{m+1}\| \leq \epsilon * \|M\| \quad \text{or} \quad \|\widetilde{w}_{m+1}\| \leq \epsilon * \|M\|,$$

while a serious breakdown was detected by

$$v_{m+1} \neq 0, \quad w_{m+1} \neq 0, \quad |a_{m+1}| \leq \epsilon * \|M\|.$$

**Example 5.1** *This set of tests was done using a  $100 \times 100$  symplectic block-diagonal matrix*

$$(29) \quad M = \text{diag}(200, 100, 50, 47, \dots, 4, 3, \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}, \frac{1}{200}, \frac{1}{100}, \frac{1}{50}, \frac{1}{47}, \dots, \frac{1}{4}, \frac{1}{3}, \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}^{-1}).$$

*A random starting vector  $v_1$  is used. The symplectic Lanczos process generates a sequence of symplectic butterfly matrices  $B^{2k,2k}$  whose eigenvalues are increasingly better approximates to eigenvalues of  $M$ . The largest Ritz value approximates the largest eigenvalue  $\lambda_1 = 200$  of  $M$ . Table 2 illustrates the loss of  $J$ -orthogonality among the symplectic Lanczos vectors in terms of  $z_1^T J_P \widehat{r}_{k+1}$  and the convergence of a Ritz value in terms of the residual  $\widehat{d}_{k+1}(e_{2k}^T y_2)$ . As predicted by Theorem 4.1, the loss of  $J$ -orthogonality accompanies the convergence of a Ritz value to the largest eigenvalue  $\lambda_1$  (and the convergence of a Ritz value to the smallest eigenvalue  $\lambda_1^{-1}$ ) in terms of small residuals. The last two columns of Table 2 report the value for  $\psi_2$  and its upper bound  $|\psi|$  from Corollary 4.2. The upper bound  $|\psi|$  is too pessimistic.*

*When the symplectic Lanczos process is stopped at  $k = 16$ , the computed largest Ritz value  $\lambda_1$  has the relative accuracy*

$$\frac{|200 - \lambda_1|}{200} \approx 1.5632 * 10^{-15}.$$

*We note that in this example, the Ritz value corresponding to the largest eigenvalue of  $M$  is well conditioned, while the condition number for all eigenvalues of  $M$  is one, the condition number of the largest Ritz value is  $\approx 1.08$ . The results for  $w_1^T J_P \widehat{r}_{k+1}$  and  $\widehat{d}_{k+1}(e_{2k}^T y_1)$  are almost the same.*

<sup>1</sup>MATLAB is a registered trademark of The MathWorks, Inc.

Lanczos step	$z_1^T J_P \hat{r}_{k+1}$	$\hat{d}_{k+1} (e_{2k}^T y_2)$	$\psi_2$	$ \psi $
1	$-9.1290 * 10^{-17}$	$-6.4278 * 10^{-01}$	$5.8679 * 10^{-17}$	$4.0869 * 10^{-08}$
2	$1.6751 * 10^{-17}$	$3.5949 * 10^{-01}$	$6.0217 * 10^{-18}$	$3.6108 * 10^{-07}$
3	$-8.4297 * 10^{-18}$	$-8.1016 * 10^{-02}$	$6.8293 * 10^{-19}$	$8.4543 * 10^{-07}$
4	$2.6983 * 10^{-17}$	$-1.7984 * 10^{-02}$	$-4.8526 * 10^{-19}$	$1.5293 * 10^{-06}$
5	$2.8513 * 10^{-16}$	$-1.3822 * 10^{-03}$	$-3.9411 * 10^{-19}$	$2.3792 * 10^{-06}$
6	$4.7089 * 10^{-15}$	$-8.3119 * 10^{-05}$	$-3.9140 * 10^{-19}$	$3.5814 * 10^{-06}$
7	$6.8569 * 10^{-14}$	$-5.7074 * 10^{-06}$	$-3.9135 * 10^{-19}$	$4.7835 * 10^{-06}$
8	$-8.3995 * 10^{-13}$	$-4.6590 * 10^{-07}$	$3.9133 * 10^{-19}$	$6.0405 * 10^{-06}$
9	$-9.3850 * 10^{-12}$	$-4.1698 * 10^{-08}$	$3.9133 * 10^{-19}$	$8.5821 * 10^{-06}$
10	$9.0525 * 10^{-11}$	$-4.3229 * 10^{-09}$	$-3.9133 * 10^{-19}$	$1.4113 * 10^{-05}$
11	$-4.1822 * 10^{-10}$	$9.3571 * 10^{-10}$	$-3.9133 * 10^{-19}$	$1.0338 * 10^{-04}$
12	$6.8361 * 10^{-09}$	$-5.7230 * 10^{-11}$	$-3.9123 * 10^{-19}$	$4.9373 * 10^{-04}$
13	$-2.9881 * 10^{-07}$	$1.3010 * 10^{-12}$	$-3.8875 * 10^{-19}$	$8.9149 * 10^{-04}$
14	$5.7946 * 10^{-06}$	$6.5210 * 10^{-14}$	$3.7786 * 10^{-19}$	$9.3192 * 10^{-04}$
15	$-1.0299 * 10^{-04}$	$2.6478 * 10^{-15}$	$-2.7270 * 10^{-19}$	$9.5668 * 10^{-04}$
16	$1.5128 * 10^{-03}$	$-1.0915 * 10^{-15}$	$-1.6512 * 10^{-18}$	$9.7898 * 10^{-04}$

Table 2: loss of  $J$ -orthogonality versus convergence of Ritz value

## Appendix

In this appendix we present proofs for Theorem 3.1, Lemma 3.3, and Theorem 4.1.

**Proof** of Theorem 3.1:

We need to analyze one step of the symplectic Lanczos algorithm. After  $j - 1$  steps of the symplectic Lanczos algorithm, we have computed  $\hat{a}_{j-1}, \hat{w}_{j-1}, \hat{c}_{j-1}, \hat{d}_j, \hat{v}_j$ . During the  $j$ th step we will compute  $\hat{a}_j, \hat{w}_j, \hat{c}_j, \hat{d}_{j+1}$  and  $\hat{v}_{j+1}$ . Recall that we set  $b_k = 1$ , hence  $b_k$  is not a computed quantity.

At first we have to compute  $a_j = v_j^T J_P M_P v_j$ . Due to its special structure, multiplication by  $J_P$  does not cause any roundoff-error; hence it will not influence our analysis. Let  $M_P$  have at most  $m$  nonzero entries in any row or column. Then for the matrix-vector multiplication  $J_P M_P v_j$  we have

$$\hat{s}_1 = fl(J_P M_P \hat{v}_j) = J_P M_P \hat{v}_j + \hat{e}_1,$$

where

$$|\hat{e}_1| \leq m\mathbf{u} |J_P M_P| |\hat{v}_j| + \mathcal{O}(\mathbf{u}^2).$$

Then  $a_j$  is computed by an inner product

$$\hat{s}_2 = fl(\hat{v}_j^T \hat{s}_1) = \hat{v}_j^T \hat{s}_1 + \hat{e}_2,$$

where

$$|\hat{e}_2| \leq 2n\mathbf{u} |\hat{v}_j|^T |\hat{s}_1| + \mathcal{O}(\mathbf{u}^2),$$

assuming that  $\hat{v}_j$  and  $\hat{s}_1$  are full vectors. Overall, we have

$$(30) \quad \hat{a}_j = \hat{v}_j^T J_P M_P \hat{v}_j + \hat{f}_j^{[1]},$$

where the roundoff error  $\hat{f}_j^{[1]} = \hat{v}_j^T \hat{e}_1 + \hat{e}_2$  is bounded by

$$\begin{aligned} |\hat{f}_j^{[1]}| &\leq m\mathbf{u} |\hat{v}_j|^T |J_P M_P| |\hat{v}_j| + 2n\mathbf{u} |\hat{v}_j|^T |\hat{s}_1| + \mathcal{O}(\mathbf{u}^2) \\ &\leq (m + 2n)\mathbf{u} |\hat{v}_j|^T |J_P M_P| |\hat{v}_j| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

Next we have to compute  $w_j = (M_P v_j - v_j)/a_j$  (recall that we choose  $b_j = 1$ ). For the matrix-vector multiplication  $M_P v_j$  we obtain

$$\widehat{s}_3 = fl(M_P \widehat{v}_j) = M_P \widehat{v}_j + \widehat{e}_3,$$

where

$$|\widehat{e}_3| \leq m\mathbf{u} |M_P| |\widehat{v}_j| + \mathcal{O}(\mathbf{u}^2).$$

The saxpy operation  $\widehat{w}_j = M_P v_j - v_j$  yields

$$\widehat{s}_4 = fl(\widehat{s}_3 - \widehat{v}_j) = \widehat{s}_3 - \widehat{v}_j + \widehat{e}_4,$$

with

$$|\widehat{e}_4| \leq \mathbf{u} (|\widehat{v}_j| + |\widehat{s}_3|) + \mathcal{O}(\mathbf{u}^2).$$

Thus overall we have

$$(31) \quad \widehat{w}_j = M_P \widehat{v}_j - \widehat{v}_j + \widehat{f}_j^{[2]},$$

where the rounding error vector  $\widehat{f}_j^{[2]} = \widehat{e}_3 + \widehat{e}_4$  is bounded by

$$\begin{aligned} |\widehat{f}_j^{[2]}| &\leq m\mathbf{u} |M_P| |\widehat{v}_j| + \mathbf{u} (|\widehat{v}_j| + |\widehat{s}_3|) + \mathcal{O}(\mathbf{u}^2) \\ &\leq m\mathbf{u} |M_P| |\widehat{v}_j| + \mathbf{u} (|\widehat{v}_j| + |M_P| |\widehat{v}_j|) + \mathcal{O}(\mathbf{u}^2) \\ &\leq (m+1)\mathbf{u} |M_P| |\widehat{v}_j| + \mathbf{u} |\widehat{v}_j| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

The computation of  $w_j$  is completed by

$$(32) \quad \widehat{w}_j = fl(\widehat{w}_j/\widehat{a}_j) = \widehat{w}_j/\widehat{a}_j + \widehat{f}_j^{[3]}$$

where the rounding error vector  $\widehat{f}_j^{[3]}$  is bounded by

$$|\widehat{f}_j^{[3]}| \leq \mathbf{u} |\widehat{w}_j \widehat{a}_j^{-1}| + \mathcal{O}(\mathbf{u}^2).$$

The analysis of the computation of  $c_j = v_j^T J_P M_P w_j / a_j$  is entirely analogous to the analysis of the computation of  $a_j$ . We start with the matrix-vector multiplication  $J_P M_P w_j$

$$\widehat{s}_5 = fl(J_P M_P \widehat{w}_j) = J_P M_P \widehat{w}_j + \widehat{e}_5,$$

where

$$|\widehat{e}_5| \leq m\mathbf{u} |J_P M_P| |\widehat{w}_j| + \mathcal{O}(\mathbf{u}^2).$$

This is followed by an inner product  $v_j^T J_P M_P w_j$

$$\widehat{s}_6 = fl(\widehat{v}_j^T \widehat{s}_5) = \widehat{v}_j^T \widehat{s}_5 + \widehat{e}_6,$$

with

$$|\widehat{e}_6| \leq 2n\mathbf{u} |\widehat{v}_j|^T |\widehat{s}_5| + \mathcal{O}(\mathbf{u}^2).$$

Finally, the computation is completed by

$$\widehat{s}_7 = fl(\widehat{s}_6/\widehat{a}_j) = \widehat{s}_6/\widehat{a}_j + \widehat{e}_7,$$

where

$$|\hat{e}_7| \leq \mathbf{u} |\hat{s}_6 \hat{a}_j^{-1}| + \mathcal{O}(\mathbf{u}^2).$$

Overall, we have

$$\hat{c}_j = \hat{v}_j^T J_P M_P \hat{w}_j / \hat{a}_j + \hat{f}_j^{[4]},$$

where the roundoff error  $\hat{f}_j^{[4]} = \hat{a}_j^{-1} \hat{v}_j^T \hat{e}_5 + \hat{a}_j^{-1} \hat{e}_6 + \hat{e}_7$  is bounded by

$$\begin{aligned} |\hat{f}_j^{[4]}| &\leq m\mathbf{u} |\hat{v}_j|^T |J_P M_P| |\hat{w}_j| |\hat{a}_j^{-1}| + 2n\mathbf{u} |\hat{v}_j|^T |J_P M_P| |\hat{w}_j| |\hat{a}_j^{-1}| \\ &\quad + \mathbf{u} |\hat{v}_j|^T |J_P M_P| |\hat{w}_j| |\hat{a}_j^{-1}| + \mathcal{O}(\mathbf{u}^2) \\ &\leq (m + 2n + 1)\mathbf{u} |\hat{a}_j^{-1}| |\hat{v}_j|^T |J_P M_P| |\hat{w}_j| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

Finally, we have to compute  $\tilde{v}_{j+1} = -d_j v_{j-1} - c_j v_j + w_j + a_j^{-1} M_P^{-1} v_j$ ,  $d_{j+1} = \sqrt{\tilde{v}_{j+1}^T \tilde{v}_{j+1}}$  and  $v_{j+1} = \tilde{v}_{j+1} / d_{j+1}$ . Recall that, as  $M$  is symplectic, the inverse of  $M_P$  is given by  $M_P^{-1} = -J_P M_P^T J_P$ . Let us start us with the matrix-vector multiplication  $M_P^{-1} v_j$

$$\hat{s}_8 = fl(M_P^{-1} \hat{v}_j) = M_P^{-1} \hat{v}_j + \hat{e}_8$$

where

$$|\hat{e}_8| \leq m\mathbf{u} |M_P^{-1}| |\hat{v}_j| + \mathcal{O}(\mathbf{u}^2).$$

Next three saxpy operations are used to finish the computation of  $\tilde{v}_{j+1}$ :

$$\hat{s}_9 = fl(\hat{s}_8 \hat{a}_j^{-1} + \hat{w}_j) = \hat{s}_8 \hat{a}_j^{-1} + \hat{w}_j + \hat{e}_9,$$

where

$$|\hat{e}_9| \leq \mathbf{u} (2|\hat{s}_8 \hat{a}_j^{-1}| + |\hat{w}_j|) + \mathcal{O}(\mathbf{u}^2),$$

and

$$\hat{s}_{10} = fl(\hat{s}_9 - \hat{c}_j \hat{v}_j) = \hat{s}_9 - \hat{c}_j \hat{v}_j + \hat{e}_{10},$$

where

$$|\hat{e}_{10}| \leq \mathbf{u} (2|\hat{c}_j \hat{v}_j| + |\hat{s}_9|) + \mathcal{O}(\mathbf{u}^2),$$

and

$$\hat{s}_{11} = fl(\hat{s}_{10} - \hat{d}_j \hat{v}_{j-1}) = \hat{s}_{10} - \hat{d}_j \hat{v}_{j-1} + \hat{e}_{11},$$

where

$$|\hat{e}_{11}| \leq \mathbf{u} (2|\hat{d}_j \hat{v}_{j-1}| + |\hat{s}_{10}|) + \mathcal{O}(\mathbf{u}^2).$$

Overall, we have for  $\tilde{v}_{j+1}$

$$(33) \quad \hat{v}_{j+1} = -\hat{d}_j \hat{v}_{j-1} - \hat{c}_j \hat{v}_j + \hat{w}_j + \hat{a}_j^{-1} M_P^{-1} \hat{v}_j + \hat{f}_{j+1}^{[5]},$$

where the roundoff error vector  $\hat{f}_{j+1}^{[5]} = \hat{a}_j^{-1} \hat{e}_8 + \hat{e}_9 + \hat{e}_{10} + \hat{e}_{11}$  is bounded by

$$\begin{aligned} |\hat{f}_{j+1}^{[5]}| &\leq m\mathbf{u} |\hat{a}_j^{-1}| |M_P^{-1}| |\hat{v}_j| + \mathbf{u} (2|\hat{s}_8 \hat{a}_j^{-1}| + |\hat{w}_j|) + \mathbf{u} (2|\hat{c}_j \hat{v}_j| + |\hat{s}_9|) \\ &\quad + \mathbf{u} (2|\hat{d}_j \hat{v}_{j-1}| + |\hat{s}_{10}|) + \mathcal{O}(\mathbf{u}^2) \\ &\leq (m + 4)\mathbf{u} |\hat{a}_j^{-1}| |M_P^{-1}| |\hat{v}_j| + 3\mathbf{u} |\hat{w}_j| + 3\mathbf{u} |\hat{c}_j \hat{v}_j| + 2\mathbf{u} |\hat{d}_j \hat{v}_{j-1}| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

Next we compute  $d_{j+1} = \sqrt{\widehat{v}_{j+1}^T \widetilde{v}_{j+1}}$ .

$$\widehat{s}_{12} = fl(\widehat{v}_{j+1}^T \widehat{v}_{j+1}) = \widehat{v}_{j+1}^T \widehat{v}_{j+1} + \widehat{e}_{12},$$

with

$$|\widehat{e}_{12}| \leq 2n\mathbf{u} |\widehat{v}_{j+1}|^T |\widehat{v}_{j+1}| + \mathcal{O}(\mathbf{u}^2).$$

Hence,

$$\widehat{d}_{j+1} = fl(\sqrt{\widehat{s}_{12}}) = \sqrt{\widehat{v}_{j+1}^T \widehat{v}_{j+1}} + \widehat{f}_{j+1}^{[6]},$$

where the roundoff error  $\widehat{f}_{j+1}^{[6]}$  is bounded by

$$|\widehat{f}_{j+1}^{[6]}| \leq \mathbf{u} \sqrt{\widehat{s}_{12}} \leq \mathbf{u} \sqrt{\widehat{v}_{j+1}^T \widehat{v}_{j+1}} + \mathcal{O}(\mathbf{u}^2).$$

The symplectic Lanczos step is completed by computing  $v_{j+1} = \widetilde{v}_{j+1}/d_{j+1}$ :

$$(34) \quad \widehat{v}_{j+1} = fl(\widehat{v}_{j+1} \widehat{d}_{j+1}^{-1}) = \widehat{v}_{j+1} \widehat{d}_{j+1}^{-1} + \widehat{f}_{j+1}^{[7]},$$

with

$$|\widehat{f}_{j+1}^{[7]}| \leq \mathbf{u} |\widehat{v}_{j+1}| |\widehat{d}_{j+1}^{-1}| + \mathcal{O}(\mathbf{u}^2).$$

From (34) and (33) we know that

$$(35) \quad \widehat{d}_{j+1} \widehat{v}_{j+1} = -\widehat{d}_j \widehat{v}_{j-1} - \widehat{c}_j \widehat{v}_j + \widehat{w}_j + \widehat{a}_j^{-1} M_P^{-1} \widehat{v}_j + g_{j+1}$$

where  $g_{j+1}$  is the sum of roundoff errors in computing the intermediate vector  $\widetilde{v}_{j+1}$  and the symplectic Lanczos vector  $v_{j+1}$

$$g_{j+1} = \widehat{f}_{j+1}^{[5]} + \widehat{d}_{j+1} \widehat{f}_{j+1}^{[7]}.$$

Using the bounds for the rounding errors  $\widehat{f}_{j+1}^{[5]}$  and  $\widehat{f}_{j+1}^{[7]}$  we have

$$(36) \quad \begin{aligned} |g_{j+1}| &\leq (m+4)\mathbf{u} |\widehat{a}_j^{-1}| |M_P^{-1}| |\widehat{v}_j| + 3\mathbf{u} |\widehat{w}_j| + 3\mathbf{u} |\widehat{c}_j \widehat{v}_j| + 2\mathbf{u} |\widehat{d}_j \widehat{v}_{j-1}| \\ &\quad + \mathbf{u} |\widehat{v}_{j+1}| + \mathcal{O}(\mathbf{u}^2) \\ &\leq (m+5)\mathbf{u} |\widehat{a}_j^{-1}| |M_P^{-1}| |\widehat{v}_j| + 4\mathbf{u} |\widehat{w}_j| + 4\mathbf{u} |\widehat{c}_j \widehat{v}_j| + 3\mathbf{u} |\widehat{d}_j \widehat{v}_{j-1}| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

Similar, from (32) and (31) we know that

$$(37) \quad \widehat{a}_j \widehat{w}_j = M_P \widehat{v}_j - \widehat{v}_j + h_j,$$

where

$$h_j = \widehat{f}_j^{[2]} + \widehat{a}_j \widehat{f}_j^{[3]},$$

and

$$(38) \quad \begin{aligned} |h_j| &\leq (m+1)\mathbf{u} |M_P| |\widehat{v}_j| + \mathbf{u} |\widehat{v}_j| + \mathbf{u} |\widehat{w}_j| + \mathcal{O}(\mathbf{u}^2) \\ &\leq (m+2)\mathbf{u} |M_P| |\widehat{v}_j| + 2\mathbf{u} |\widehat{v}_j| + \mathcal{O}(\mathbf{u}^2). \end{aligned}$$

While the equation  $a_j w_j = M_P v_j - v_j$  is given by the  $(2j)$ th column of  $M_P S_P N_P^{-1} = S_P K_P^{-1}$ , the equation  $d_{j+1} v_{j+1} = -d_j v_{j-1} - c_j v_j + w_j + a_j^{-1} M_P^{-1} v_j$  corresponds to the  $(2j-1)$ th column of



This completes the proof of Theorem 3.1. ✓

**Proof of Lemma 3.3:**

Obviously,

$$k_{2j,2j} = k_{2j-1,2j-1} = 0$$

for  $j = 1, \dots, k$  as  $x^T J_P x = 0$  for any vector  $x$ . Moreover, as  $k_{2m,2j-1} = -k_{2j-1,2m}$ , we only need to examine  $k_{2j,2j-1}$  for  $j = 1, \dots, k$ , and  $k_{2j,2m-1}$ ,  $k_{2j-1,2m-1}$  and  $k_{2j,2m}$  for  $j, m = 1, \dots, k$ ,  $j < m$ .

Let us start with  $k_{2j,2j-1}$ . Using (32) and (34) we have

$$\begin{aligned} (42) \quad k_{2j,2j-1} &= \widehat{w}_j^T J_P \widehat{v}_j \\ &= \left( \frac{\widehat{w}_j}{\widehat{a}_j} + \widehat{f}_j^{[3]} \right)^T J_P \left( \frac{\widehat{v}_j}{\widehat{d}_j} + \widehat{f}_j^{[7]} \right) \\ &= \frac{\widehat{w}_j^T J_P \widehat{v}_j + \widehat{a}_j (\widehat{f}_j^{[3]})^T J_P \widehat{v}_j + \widehat{d}_j \widehat{w}_j^T J_P \widehat{f}_j^{[7]}}{\widehat{a}_j \widehat{d}_j} + \mathcal{O}(\mathbf{u}^2) \\ (43) \quad &=: \frac{\widehat{w}_j^T J_P \widehat{v}_j + \zeta_1}{\widehat{a}_j \widehat{d}_j} + \mathcal{O}(\mathbf{u}^2), \end{aligned}$$

where

$$\begin{aligned} |\zeta_1| &\leq |\widehat{a}_j (\widehat{f}_j^{[3]})^T J_P \widehat{v}_j| + |\widehat{d}_j \widehat{w}_j^T J_P \widehat{f}_j^{[7]}| \\ &\leq 2\mathbf{u} |\widehat{w}_j|^T |J_P| |\widehat{v}_j| \\ &\leq 2\mathbf{u} |\widehat{v}_j|^T |J_P| (|M_P| |\widehat{v}_j| - |\widehat{v}_j|). \end{aligned}$$

We would like to be able to rewrite  $k_{2j,2j-1} = \widehat{w}_j^T J_P \widehat{v}_j = -1 + \text{some small error}$ . In order to do so, we rewrite  $\widehat{a}_j \widehat{d}_j$  suitably. From (30) and (34) we have

$$\begin{aligned} \widehat{a}_j \widehat{d}_j &= (\widehat{v}_j^T J_P M_P \widehat{v}_j + \widehat{f}_j^{[1]}) \widehat{d}_j \\ &= \left[ \left( \frac{\widehat{v}_j}{\widehat{d}_j} + \widehat{f}_j^{[7]} \right)^T J_P M_P \widehat{v}_j + \widehat{f}_j^{[1]} \right] \widehat{d}_j \\ &= \widehat{v}_j^T J_P M_P \widehat{v}_j + \widehat{d}_j ((\widehat{f}_j^{[7]})^T J_P M_P \widehat{v}_j + \widehat{f}_j^{[1]}). \end{aligned}$$

Using (31) we obtain

$$\begin{aligned} \widehat{a}_j \widehat{d}_j &= \widehat{v}_j^T J_P (\widehat{w}_j + \widehat{v}_j - \widehat{f}_j^{[2]}) + \widehat{d}_j ((\widehat{f}_j^{[7]})^T J_P M_P \widehat{v}_j + \widehat{f}_j^{[1]}) \\ &= \widehat{v}_j^T J_P \widehat{w}_j + \widehat{v}_j^T J_P \widehat{v}_j - \widehat{v}_j^T J_P \widehat{f}_j^{[2]} + \widehat{d}_j ((\widehat{f}_j^{[7]})^T J_P M_P \widehat{v}_j + \widehat{f}_j^{[1]}) \\ &= \widehat{v}_j^T J_P \widehat{w}_j + \widehat{d}_j (\widehat{v}_j - \widehat{f}_j^{[2]})^T J_P \widehat{v}_j - \widehat{v}_j^T J_P \widehat{f}_j^{[2]} + \widehat{d}_j ((\widehat{f}_j^{[7]})^T J_P M_P \widehat{v}_j + \widehat{f}_j^{[1]}). \end{aligned}$$

For the last equation we used again (34). This rewriting allows us to make use of the fact that  $x^T J x = 0$  for any vector  $x$ . Thus

$$\begin{aligned} \widehat{a}_j \widehat{d}_j &= \widehat{v}_j^T J_P \widehat{w}_j - \widehat{v}_j^T J_P \widehat{f}_j^{[2]} + \widehat{d}_j (\widehat{f}_j^{[7]})^T J_P (M_P \widehat{v}_j - \widehat{v}_j) + \widehat{d}_j \widehat{f}_j^{[1]} \\ &= \widehat{v}_j^T J_P \widehat{w}_j - \widehat{v}_j^T J_P \widehat{f}_j^{[2]} + \widehat{d}_j (\widehat{f}_j^{[7]})^T J_P (\widehat{w}_j - \widehat{f}_j^{[2]}) + \widehat{d}_j \widehat{f}_j^{[1]} \\ (44) \quad &=: \widehat{v}_j^T J_P \widehat{w}_j + \zeta_2 \\ &= -\widehat{w}_j^T J_P \widehat{v}_j + \zeta_2, \end{aligned}$$

where we used (31). The roundoff error is bounded by

$$\begin{aligned}
|\zeta_2| &\leq |\widehat{v}_j^T J_P \widehat{f}_j^{[2]}| + |\widehat{d}_j| |\widehat{f}_j^{[7]}|^T |J_P| (|\widehat{w}_j| + |\widehat{f}_j^{[2]}|) + |\widehat{d}_j| |\widehat{f}_j^{[1]}| \\
&\leq (m+1)\mathbf{u} |\widehat{v}_j^T|^T |J_P| |M_P| |\widehat{v}_j| + 2\mathbf{u} |\widehat{v}_j^T|^T |J_P| |\widehat{v}_j| \\
&\quad + \mathbf{u} |\widehat{v}_j^T|^T |J_P| |\widehat{w}_j| + (m+2n)\mathbf{u} |\widehat{d}_j| |\widehat{v}_j^T|^T |J_P M_P| |\widehat{v}_j| + \mathcal{O}(\mathbf{u}^2) \\
&\leq (m+1)\mathbf{u} |\widehat{v}_j^T|^T |J_P| |M_P| |\widehat{v}_j| + 2\mathbf{u} |\widehat{v}_j^T|^T |J_P| |\widehat{v}_j| \\
&\quad + \mathbf{u} |\widehat{v}_j^T|^T |J_P| (|M_P| |\widehat{v}_j| + |\widehat{v}_j|) \\
&\quad + (m+2n)\mathbf{u} |\widehat{v}_j^T|^T |J_P| |M_P| |\widehat{v}_j| + \mathcal{O}(\mathbf{u}^2) \\
&\leq (2m+2n+2)\mathbf{u} |\widehat{v}_j^T|^T |J_P| |M_P| |\widehat{v}_j| + 3\mathbf{u} |\widehat{v}_j^T|^T |J_P| |\widehat{v}_j| + \mathcal{O}(\mathbf{u}^2).
\end{aligned}$$

Combining (43) and (44) we have

$$\begin{aligned}
k_{2j,2j-1} &= \frac{\widehat{w}_j^T J_P \widehat{v}_j + \zeta_1}{-\widehat{w}_j^T J_P \widehat{v}_j + \zeta_2} + \mathcal{O}(\mathbf{u}^2) \\
&= -1 + \frac{\zeta_1 + \zeta_2}{-\widehat{w}_j^T J_P \widehat{v}_j + \zeta_2} + \mathcal{O}(\mathbf{u}^2) \\
&=: -1 + \kappa_j + \mathcal{O}(\mathbf{u}^2).
\end{aligned}$$

Using the Taylor expansion of  $f(x) = \frac{\zeta_1 + \zeta_2}{x + \zeta_2}$  at  $t = x - \zeta_2$ ,

$$\begin{aligned}
f(x) &= f(t) + f'(t)(x-t) + \frac{f''(t)}{2}(x-t)^2 + \text{higher order terms} \\
&= \frac{\zeta_1 + \zeta_2}{x} - \frac{\zeta_1 + \zeta_2}{x^2} \zeta_2 + \frac{\zeta_1 + \zeta_2}{x^3} \zeta_2^2 + \text{higher order terms},
\end{aligned}$$

we obtain

$$\begin{aligned}
|\kappa_j| &\leq \frac{|\zeta_1| + |\zeta_2|}{|\widehat{w}_j^T J_P \widehat{v}_j|} + \mathcal{O}(\mathbf{u}^2) \\
(45) \quad &\leq \frac{2(m+n+2)\mathbf{u} |\widehat{v}_j^T|^T |J_P| |M_P| |\widehat{v}_j| + 5\mathbf{u} |\widehat{v}_j^T|^T |J_P| |\widehat{v}_j|}{|\widehat{w}_j^T J_P \widehat{v}_j|} + \mathcal{O}(\mathbf{u}^2).
\end{aligned}$$

Next we turn our attention to the terms  $k_{2j,2m-1}$ ,  $k_{2j-1,2m-1}$ , and  $k_{2j,2m}$ . The analysis of these three terms will be demonstrated by considering  $k_{2j,2m} = \widehat{w}_j^T J_P \widehat{w}_m$ . Let us assume that we have already analyzed all previous terms, that is, all the terms in the  $2m \times 2m$  leading principal submatrix of  $K$ , printed in bold face,

$$\begin{bmatrix}
\mathbf{k}_{11} & \cdots & \mathbf{k}_{1,2j-1} & \mathbf{k}_{1,2j} & \cdots & \mathbf{k}_{1,2m-1} & \mathbf{k}_{1,2m} \\
\vdots & & \vdots & \vdots & & \vdots & \vdots \\
\mathbf{k}_{2j-1,1} & \cdots & \mathbf{k}_{2j-1,2j-1} & \mathbf{k}_{2j-1,2j} & \cdots & \mathbf{k}_{2j-1,2m-1} & \mathbf{k}_{2j-1,2m} \\
\mathbf{k}_{2j,1} & \cdots & \mathbf{k}_{2j,2j-1} & \mathbf{k}_{2j,2j} & \cdots & \mathbf{k}_{2j,2m-1} & k_{2j,2m} \\
\mathbf{k}_{2j+1,1} & \cdots & \mathbf{k}_{2j+1,2j-1} & \mathbf{k}_{2j+1,2j} & \cdots & \mathbf{k}_{2j+1,2m-1} & k_{2j+1,2m} \\
\vdots & & \vdots & \vdots & & \vdots & \vdots \\
\mathbf{k}_{2m-1,1} & \cdots & \mathbf{k}_{2m-1,2j-1} & \mathbf{k}_{2m-1,2j} & \cdots & \mathbf{k}_{2m-1,2m-1} & k_{2m-1,2m} \\
\mathbf{k}_{2m,1} & \cdots & \mathbf{k}_{2m,2j-1} & k_{2m,2j} & \cdots & k_{2m,2m-1} & k_{2m,2m}
\end{bmatrix}.$$

Our goal is to rewrite  $k_{2j,2m}$  in terms of any of these already analyzed terms. First of all, note, that for  $j = m$  we have  $k_{2m,2m} = 0$ . Hence for the following discussion we assume  $j < m$ . From

(37) we have

$$\begin{aligned}\widehat{a}_m k_{2j,2m} &= \widehat{w}_j^T J_P (M_P \widehat{v}_m - \widehat{v}_m + h_m) \\ &= \widehat{w}_j^T J_P M_P \widehat{v}_m - k_{2j,2m-1} + \widehat{w}_j^T J_P h_m.\end{aligned}$$

Using (35) we obtain for  $\widehat{w}_j^T J_P M_P \widehat{v}_m$

$$\begin{aligned}\widehat{d}_m \widehat{w}_j^T J_P M_P \widehat{v}_m &= -\widehat{d}_{m-1} \widehat{w}_j^T J_P M_P \widehat{v}_{m-2} - \widehat{c}_{m-1} \widehat{w}_j^T J_P M_P \widehat{v}_{m-1} + \widehat{w}_j^T J_P M_P \widehat{w}_{m-1} \\ &\quad + \widehat{a}_{m-1}^{-1} \widehat{w}_j^T J_P \widehat{v}_{m-1} + \widehat{w}_j^T J_P M_P g_m.\end{aligned}$$

Using (37) twice yields

$$\begin{aligned}\widehat{d}_m \widehat{w}_j^T J_P M_P \widehat{v}_m &= -\widehat{d}_{m-1} \widehat{w}_j^T J_P (\widehat{a}_{m-2} \widehat{w}_{m-2} + \widehat{v}_{m-2} - h_{m-2}) \\ &\quad - \widehat{c}_{m-1} \widehat{w}_j^T J_P (\widehat{a}_{m-1} \widehat{w}_{m-1} + \widehat{v}_{m-1} - h_{m-1}) \\ &\quad + \widehat{w}_j^T J_P M_P \widehat{w}_{m-1} + \widehat{a}_{m-1}^{-1} k_{2j,2m-3} + \widehat{w}_j^T J_P M_P g_m \\ &= \widehat{w}_j^T J_P M_P \widehat{w}_{m-1} - \widehat{d}_{m-1} \widehat{a}_{m-2} k_{2j,2m-4} - \widehat{d}_{m-1} k_{2j,2m-5} \\ &\quad - \widehat{c}_{m-1} \widehat{a}_{m-1} k_{2j,2m-2} - \widehat{c}_{m-1} k_{2j,2m-3} + \widehat{a}_{m-1}^{-1} k_{2j,2m-3} \\ &\quad + \widehat{d}_{m-1} \widehat{w}_j^T J_P h_{m-2} + \widehat{c}_{m-1} \widehat{w}_j^T J_P h_{m-1} + \widehat{w}_j^T J_P M_P g_m.\end{aligned}$$

The last term that needs our attention here is  $\widehat{w}_j^T J_P M_P \widehat{w}_{m-1}$ . From (37) we have

$$\begin{aligned}\widehat{a}_j \widehat{w}_j^T J_P M_P \widehat{w}_{m-1} &= (M_P \widehat{v}_j - \widehat{v}_j + h_j)^T J_P M_P \widehat{w}_{m-1} \\ &= \widehat{v}_j^T M_P^T J_P M_P \widehat{w}_{m-1} - \widehat{v}_j^T J_P M_P \widehat{w}_{m-1} + h_j^T J_P M_P \widehat{w}_{m-1} \\ &= \widehat{v}_j^T J_P \widehat{w}_{m-1} + \widehat{w}_{m-1}^T J_P M_P^{-1} \widehat{v}_j + h_j^T J_P M_P \widehat{w}_{m-1}\end{aligned}$$

as  $M$  is symplectic. Using (35) yields

$$\begin{aligned}\widehat{a}_j \widehat{w}_j^T J_P M_P \widehat{w}_{m-1} &= k_{2j-1,2m-2} + h_j^T J_P M_P \widehat{w}_{m-1} \\ &\quad + \widehat{a}_j \widehat{w}_{m-1}^T J_P (\widehat{d}_{j+1} \widehat{v}_{j+1} + \widehat{d}_j \widehat{v}_{j-1} + \widehat{c}_j \widehat{v}_j - \widehat{w}_j - g_{j+1}) \\ &= \widehat{a}_j (\widehat{d}_{j+1} k_{2m-2,2j+1} + \widehat{d}_j k_{2m-2,2j-3} + \widehat{c}_j k_{2m-2,2j-1} + k_{2m-2,2m}) \\ &\quad + k_{2j-1,2m-2} + h_j^T J_P M_P \widehat{w}_{m-1} - \widehat{a}_j \widehat{w}_{m-1}^T J_P g_{j+1}.\end{aligned}$$

Therefore,

$$\begin{aligned}\widehat{d}_m \widehat{a}_m k_{2j,2m} &= \widehat{a}_j^{-1} k_{2j-1,2m-2} - \widehat{d}_m k_{2j,2m-1} - \widehat{c}_{m-1} \widehat{a}_{m-1} k_{2j,2m-2} - \widehat{c}_{m-1} k_{2j,2m-3} \\ &\quad + \widehat{a}_{m-1}^{-1} k_{2j,2m-3} - \widehat{d}_{m-1} \widehat{a}_{m-2} k_{2j,2m-4} - \widehat{d}_{m-1} k_{2j,2m-5} \\ &\quad + \widehat{d}_j k_{2m-2,2j-3} + \widehat{c}_j k_{2m-2,2j-1} + \widehat{d}_{j+1} k_{2m-2,2j+1} + k_{2m-2,2m} \\ &\quad + \widehat{d}_{m-1} \widehat{w}_j^T J_P h_{m-2} + \widehat{c}_{m-1} \widehat{w}_j^T J_P h_{m-1} + \widehat{d}_m \widehat{w}_j^T J_P h_m \\ &\quad + \widehat{a}_j^{-1} h_j^T J_P M_P \widehat{w}_{m-1} - \widehat{w}_{m-1}^T J_P g_{j+1} + \widehat{w}_j^T J_P M_P g_m.\end{aligned}$$

Hence, as  $|h_j| \approx \mathcal{O}(\mathbf{u})$  and  $|g_j| \approx \mathcal{O}(\mathbf{u})$  we obtain  $k_{2j,2m} \approx \mathcal{O}(\mathbf{u})$ . A similar analysis can be done for  $k_{2j,2m-1}$  and  $k_{2j-1,2m-1}$ .  $\checkmark$

**Proof** of Theorem 4.1:

Our goal is to derive expressions for  $z_j^T J_P \widehat{r}_{k+1}$  and  $x_j^T J_P \widehat{r}_{k+1}$  that describe the way in which  $J$ -orthogonality is lost. In exact arithmetic, these expressions are zero. Our approach follows Bai's derivations in [1, Proof of Theorem 4.1]. Pre-multiplying (22) by  $(\widehat{S}_P^{2k})^T$  and taking the transpose yields

$$(\widehat{W}_P^{2k})^T M_P \widehat{S}_P^{2k} = \widehat{B}_P^{2k,2k} (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} + (\widehat{K}_P^{2k,2k})^{-1} \left[ \widehat{d}_{k+1} J_P \widehat{v}_{k+1} e_{2k}^T + F_k \right]^T \widehat{S}_P^{2k}.$$

Pre-multiplying (21) by  $(\widehat{W}_P^{2k})^T$  we obtain

$$(\widehat{W}_P^{2k})^T M_P \widehat{S}_P^{2k} = (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} \widehat{B}_P^{2k,2k} - (\widehat{W}_P^{2k})^T \left[ \widehat{d}_{k+1} \widehat{r}_{k+1} e_{2k-1}^T - E_k \right] \widehat{N}_P^{2k,2k}.$$

Subtracting these two equations, we obtain

$$\begin{aligned} & (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} \\ &= \widehat{d}_{k+1} (\widehat{K}_P^{2k,2k})^{-1} e_{2k} \widehat{v}_{k+1}^T J_P^T \widehat{S}_P^{2k} + \widehat{d}_{k+1} (\widehat{W}_P^{2k})^T \widehat{r}_{k+1} e_{2k-1}^T \widehat{N}_P^{2k,2k} \\ &\quad + (\widehat{K}_P^{2k,2k})^{-1} F_k^T \widehat{S}_P^{2k} - (\widehat{W}_P^{2k})^T E_k \widehat{N}_P^{2k,2k} \\ &= \widehat{d}_{k+1} (\widehat{K}_P^{2k,2k})^{-1} e_{2k} \widehat{v}_{k+1}^T J_P^T \widehat{S}_P^{2k} - \widehat{d}_{k+1} (\widehat{W}_P^{2k})^T \widehat{r}_{k+1} e_{2k}^T \\ &\quad + (\widehat{K}_P^{2k,2k})^{-1} F_k^T \widehat{S}_P^{2k} - (\widehat{W}_P^{2k})^T E_k \widehat{N}_P^{2k,2k}. \end{aligned} \tag{46}$$

We are most interested in deriving an expression for

$$(\widehat{S}_P^{2k})^T J_P \widehat{r}_{k+1} e_{2k-1}^T \quad (\text{or } (\widehat{W}_P^{2k})^T J_P \widehat{r}_{k+1} e_{2k-1}^T)$$

from the above equation. From this we can easily obtain expressions for  $z_j^T J_P r_{k+1}$  or  $x_j^T J_P r_{k+1}$  as desired. In order to do so, we note that most of the matrices in (46) have a very special form. Let us start with the left-hand side. From (23) we have  $(\widehat{S}_P^{2k})^T J_P \widehat{S}_P^{2k} = K = J_P^{2k,2k} + C_k + \Delta_k - C_k^T$ . This implies

$$\begin{aligned} (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} &= J_P^{2k,2k} (\widehat{S}_P^{2k})^T J_P^{2n,2n} \widehat{S}_P^{2k} \\ &= J_P^{2k,2k} K \\ &= -I^{2k,2k} + J_P^{2k,2k} C_k + J_P^{2k,2k} \Delta_k - J_P^{2k,2k} C_k^T, \end{aligned}$$

where  $J_P^{2k,2k} C_k$  and  $(J_P^{2k,2k} C_k^T)^T$  have the same form as  $C_k$ , and  $J_P^{2k,2k} \Delta_k$  is a diagonal matrix,

$$J_P^{2k,2k} \Delta_k = \text{diag} \left( \begin{bmatrix} -\kappa_1 & 0 \\ 0 & -\kappa_1 \end{bmatrix}, \dots, \begin{bmatrix} -\kappa_k & 0 \\ 0 & -\kappa_k \end{bmatrix} \right).$$

Therefore, we can rewrite the left-hand side of (46) as

$$\begin{aligned} & (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} \\ &= \left[ -I^{2k,2k} + J_P^{2k,2k} C_k + J_P^{2k,2k} \Delta_k + J_P^{2k,2k} C_k^T \right] \widehat{B}_P^{2k,2k} \\ &\quad - \widehat{B}_P^{2k,2k} \left[ -I^{2k,2k} + J_P^{2k,2k} C_k + J_P^{2k,2k} \Delta_k + J_P^{2k,2k} C_k^T \right] \\ &= \left[ J_P^{2k,2k} C_k \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} C_k \right] + \left[ J_P^{2k,2k} \Delta_k \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} \Delta_k \right] \\ &\quad + \left[ J_P^{2k,2k} C_k^T \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} C_k^T \right]. \end{aligned}$$

By the local  $J$ -orthogonality assumption (and, therefore, by the special form of  $J_P^{2k,2k} C_k$ ), it follows that

$$L_k^{(1)} := J_P^{2k,2k} C_k \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} C_k$$

is a strictly lower block triangular matrix with block size 2. With the same argument we have that

$$U_k^{(1)} := J_P^{2k,2k} C_k^T \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} C_k^T$$

is a strictly upper block triangular matrix with block size 2. Since the  $2 \times 2$  diagonal blocks of  $J_P^{2k,2k} \Delta_k \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} \Delta_k$  are zero, we can write

$$J_P^{2k,2k} \Delta_k \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} \Delta_k = L_k^{(2)} + U_k^{(2)},$$

where  $L_k^{(2)}$  is strictly lower block triangular and  $U_k^{(2)}$  strictly upper block triangular. Hence,

$$(\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} (\widehat{W}_P^{2k})^T \widehat{S}_P^{2k} = L_k^{(1)} + L_k^{(2)} + U_k^{(1)} + U_k^{(2)}.$$

Now let us turn our attention to the right-hand side of (46). The row vector

$$\widehat{v}_{k+1}^T J_P^T \widehat{S}_P^{2k} = [* \dots * 0 0]$$

has nonzero elements in its first  $(2n - 2)$  positions. As  $(\widehat{K}_P^{2k,2k})^{-1} e_{2k} = b_k e_{2k-1} + a_k e_{2k}$  we have that

$$L_k^{(3)} := \widehat{d}_{k+1} (\widehat{K}_P^{2k,2k})^{-1} e_{2k} \widehat{v}_{k+1}^T J_P^T \widehat{S}_P^{2k}$$

is a strictly lower block triangular matrix with block size 2. Similarly we have that

$$U_k^{(3)} := \widehat{d}_{k+1} (\widehat{W}_P^{2k})^T \widehat{r}_{k+1} e_{2k}^T$$

is a strictly upper block triangular matrix with block size 2. Hence, we can rewrite (46) as

$$(47) \quad \begin{aligned} L_k^{(1)} + L_k^{(2)} - L_k^{(3)} + U_k^{(1)} + U_k^{(2)} - U_k^{(3)} \\ = (\widehat{K}_P^{2k,2k})^{-1} F_k^T \widehat{S}_P^{2k} - (\widehat{W}_P^{2k})^T E_k \widehat{N}_P^{2k,2k}. \end{aligned}$$

This implies that the diagonal blocks of  $(\widehat{K}_P^{2k,2k})^{-1} F_k^T \widehat{S}_P^{2k} - (\widehat{W}_P^{2k})^T E_k \widehat{N}_P^{2k,2k}$  must be zero. Therefore, we can write

$$(\widehat{K}_P^{2k,2k})^{-1} F_k^T \widehat{S}_P^{2k} - (\widehat{W}_P^{2k})^T E_k \widehat{N}_P^{2k,2k} = L_k^{(4)} + U_k^{(4)}$$

where  $L_k^{(4)}$  is strictly lower block triangular and  $U_k^{(4)}$  is strictly upper block triangular. By writing down only the strictly upper block triangular part of (47) we have

$$U_k^{(3)} = U_k^{(1)} + U_k^{(2)} - U_k^{(4)}$$

or

$$\widehat{d}_{k+1} (\widehat{W}_P^{2k})^T \widehat{r}_{k+1} e_{2k}^T = J_P^{2k,2k} C_k^T \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} C_k^T + U_k^{(2)} - U_k^{(4)}.$$

This is equivalent to

$$(48) \quad \begin{aligned} & \widehat{d}_{k+1} (\widehat{S}_P^{2k})^T J_P^{2n,2n} \widehat{r}_{k+1} e_{2k}^T \\ & = - J_P^{2k,2k} \left[ J_P^{2k,2k} C_k^T \widehat{B}_P^{2k,2k} - \widehat{B}_P^{2k,2k} J_P^{2k,2k} C_k^T + U_k^{(2)} - U_k^{(4)} \right] \\ & = C_k^T \widehat{B}_P^{2k,2k} - (\widehat{B}_P^{2k,2k})^{-T} C_k^T - J_P^{2k,2k} \left[ U_k^{(2)} - U_k^{(4)} \right], \end{aligned}$$

where we have used the fact that  $\widehat{B}_P^{2k,2k}$  is symplectic.

From (25) we get

$$\widehat{B}_P^{2k,2k} y_{2j-1} = \lambda_j y_{2j-1} \quad \text{and} \quad \widehat{B}_P^{2k,2k} y_{2j} = \lambda_j^{-1} y_{2j}.$$

This implies

$$y_{2j-1}^T (\widehat{B}_P^{2k,2k})^{-T} = \lambda_j^{-1} y_{2j-1}^T \quad \text{and} \quad y_{2j}^T (\widehat{B}_P^{2k,2k})^{-T} = \lambda_j y_{2j}^T.$$

Pre-multiplying (48) by  $y_{2j}^T$  and post-multiplying by  $y_{2j-1}$  yields

$$\begin{aligned} & \widehat{d}_{k+1} y_{2j}^T (\widehat{S}_P^{2k})^T J_P^{2n,2n} \widehat{r}_{k+1} (e_{2k}^T y_{2j-1}) \\ & = y_{2j}^T C_k^T \widehat{B}_P^{2k,2k} y_{2j-1} - y_{2j}^T (\widehat{B}_P^{2k,2k})^{-T} C_k^T y_{2j-1} - y_{2j}^T J_P^{2k,2k} \left[ U_k^{(2)} - U_k^{(4)} \right] y_{2j-1} \\ & = \lambda_j y_{2j}^T C_k^T y_{2j-1} - \lambda_j y_{2j}^T C_k^T y_{2j-1} - y_{2j}^T J_P^{2k,2k} \left[ U_k^{(2)} - U_k^{(4)} \right] y_{2j-1} \\ & = y_{2j}^T J_P^{2k,2k} \left[ U_k^{(4)} - U_k^{(2)} \right] y_{2j-1}. \end{aligned}$$

Similarly, pre-multiplying (48) by  $y_{2j-1}^T$  and post-multiplying by  $y_{2j}$  yields

$$\hat{d}_{k+1} y_{2j-1}^T (\hat{S}_P^{2k})^T J_P^{2n, 2n} \hat{r}_{k+1} e_{2k}^T y_{2j} = y_{2j-1}^T J_P^{2k, 2k} \left[ U_k^{(4)} - U_k^{(2)} \right] y_{2j}.$$

With the assumptions (24) and (25) this concludes the proof of Theorem 4.1.  $\checkmark$

## References

- [1] Z. BAI, *Error analysis of the Lanczos algorithm for the nonsymmetric eigenvalue problem*, Math. Comp., 62 (1994), pp. 209–226.
- [2] G. BANSE, *Symplektische Eigenwertverfahren zur Lösung zeitdiskreter optimaler Steuerungsprobleme*, PhD thesis, Universität Bremen, Fachbereich 3 - Mathematik und Informatik, Bremen, Germany, 1995.
- [3] P. BENNER AND H. FASSBENDER, *An implicitly restarted symplectic Lanczos method for the symplectic eigenvalue problem*, Berichte aus der Technomathematik, Report 98-01, Universität Bremen, Fachbereich 3 - Mathematik und Informatik, Bremen, Germany, 1998.
- [4] ———, *The symplectic eigenvalue problem, the butterfly form, the SR algorithm, and the Lanczos method*, Linear Algebra Appl., 275–276 (1998), pp. 19–47.
- [5] H. FASSBENDER, *Symplectic Methods for Symplectic Eigenproblems*, Habilitationsschrift, Universität Bremen, Fachbereich 3 - Mathematik und Informatik, Bremen, Germany, 1998.
- [6] R. FREUND, *Transpose-free quasi-minimal residual methods for non-Hermitian linear systems*, in Recent advances in iterative methods. Papers from the IMA workshop on iterative methods for sparse and structured problems, held in Minneapolis, MN, February 24-March 1, 1992., G. Golub et al., ed., vol. 60 of IMA Vol. Math. Appl., New York, NY, 1994, Springer-Verlag, pp. 69–94.
- [7] G. GOLUB AND C. VAN LOAN, *Matrix Computations*, Johns Hopkins University Press, Baltimore, 3rd ed., 1996.
- [8] N. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM Publications, Philadelphia, PA, 1996.
- [9] W. KAHAN, B. PARLETT, AND E. JIANG, *Residual bounds on approximate eigensystems of nonnormal matrices*, SIAM J. Numer. Anal., 19 (1982), pp. 470–484.
- [10] P. LANCASTER AND L. RODMAN, *The Algebraic Riccati Equation*, Oxford University Press, Oxford, 1995.
- [11] C. LANCZOS, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*, J. Res. Nat. Bur. Standards, 45 (1950), pp. 255–282.
- [12] V. MEHRMANN, *The Autonomous Linear Quadratic Control Problem, Theory and Numerical Solution*, no. 163 in Lecture Notes in Control and Information Sciences, Springer-Verlag, Heidelberg, July 1991.
- [13] C. PAIGE, *Error analysis of the Lanczos algorithm for tridiagonalizing a symmetric matrix*, J. Inst. Math. Applics., 18 (1976), pp. 341–349.
- [14] B. PARLETT, *A new look at the Lanczos algorithm for solving symmetric systems of linear equations*, Linear Algebra Appl., 29 (1980), pp. 323–346.
- [15] ———, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, New Jersey, 1980.

- [16] H. SIMON, *Analysis of the symmetric Lanczos algorithm with reorthogonalization methods*, Linear Algebra Appl., 61 (1984), pp. 101–132.
- [17] K. Zhou, J.C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, Upper Saddle River, NJ, 1996.