

Conversion Methods, Block Triangularization, and Structural Analysis of Differential-Algebraic Equation Systems

Guangning Tan · Nedialko S. Nedialkov ·
John D. Pryce

Received: date / Accepted: date

Abstract In a previous article, the authors developed two conversion methods to improve the Σ -method for structural analysis (SA) of differential-algebraic equations (DAEs). These methods reformulate a DAE on which the Σ -method fails into an equivalent problem on which this SA is more likely to succeed with a generically nonsingular Jacobian. The basic version of these methods processes the DAE as a whole. This article presents the block version that exploits block triangularization of a DAE. Using a block triangular form of a Jacobian sparsity pattern, we identify which diagonal blocks of the Jacobian are identically singular and then perform a conversion on each such block. This approach improves the efficiency of finding a suitable conversion for fixing SA's failures. All of our conversion methods can be implemented in a computer algebra system so that every conversion can be automated.

Keywords differential-algebraic equations · structural analysis · block triangular form · modeling · symbolic computation

Mathematics Subject Classification (2000) 34A09 · 65L80 · 41A58 · 68W30

1 Introduction.

This article is a continuation of [13], in which we presented two conversion methods for improving the Σ -method [9] for structural analysis (SA) of DAEs. When this SA fails on a DAE with an identically singular (but structurally nonsingular) System

G. Tan
School of Computational Science and Engineering, McMaster University,
1280 Main Street West, Hamilton, Ontario L8S 4K1, Canada, E-mail: tang4@mcmaster.ca

N. S. Nedialkov
Department of Computing and Software, McMaster University,
1280 Main Street West, , Hamilton, Ontario L8S 4K1, Canada, E-mail: nedialk@mcmaster.ca

J. D. Pryce
School of Mathematics, Cardiff University,
Senghennydd Road, Cardiff CF24 4AG, Wales, UK., E-mail: prycej1@cardiff.ac.uk

Jacobian, our conversion methods reformulate the DAE into an equivalent problem on which the SA is more likely to succeed with a generically nonsingular System Jacobian [12, 13].

These two conversion methods are the linear combination (LC) method and the expression substitution (ES) method. The former is based on replacing an existing equation by a linear combination of some equations and derivatives of them. The latter is based on replacing some existing derivatives¹ by expressions that contain newly introduced variables and derivatives of them. In the ES method, the equations that prescribe such replacements are also appended to the original DAE, so the resulting system is an enlarged one. The main result of a conversion using either method is a strict decrease in the value of the signature matrix [13]. Based on our experience, we conjecture that such a decrease tends to give a better problem formulation of a DAE from SA perspective.

Our works [6, 10, 11] show how to construct block triangular forms (BTFs) of a DAE using the structural data obtained from the Σ -method. A BTF indicates how each part of the DAE influences [resp. is influenced by] other parts. The interdependences between all pairs of blocks may be depicted by a fine-block graph [11]. Exploiting the underlying structure of a DAE, we can compute the derivatives of its solution in a blockwise fashion [7], or perform a dummy derivative index reduction algorithm [4, 5].

We refer the reader to the previous article [13] for a summary of the Σ -method, details of its failures, the basic conversion methods, and explanations of the equivalence of DAEs. By “basic” we mean that these methods do not exploit BTFs of a DAE. We shall follow the notation in [13].

In this article, we combine our conversion methods with a block triangularization of a DAE and derive our block conversion methods. When the System Jacobian is identically singular, and the DAE has a nontrivial BTF—that is, having at least two diagonal blocks—we can identify which blocks are identically singular and perform a conversion on each such block. Now that we only deal with equations and variables within a block, which is usually of a smaller size compared to the whole DAE, these block methods require fewer symbolic computations and hence are expected to be more efficient in finding a useful conversion for fixing SA’s failures.

Section 2 reviews BTFs of a sparsity pattern and BTFs of a DAE. Section 3 presents our block conversion methods and demonstrates their application on a DAE from [1]. Section 4 gives more examples, in which the two DAEs are obtained from electrical circuit analysis [3]. Section 5 gives concluding remarks.

2 Block triangularization of DAEs.

In §2.1, we introduce notation for a BTF of a sparsity pattern. In §2.2, we review how to derive a BTF of a DAE; more details are in [10, 11].

We do not repeat the definitions and formulas for the notation in the Σ -method theory, such as a *signature matrix* $\Sigma = (\sigma_{ij})$ and its *value* $\text{Val}(\Sigma)$, a *highest-value*

¹Throughout this article, “derivatives of a variable” include the variable itself as its 0th derivative.

transversal (HVT) T of Σ , a *valid offset pair* $(\mathbf{c}; \mathbf{d})$, a *System Jacobian* $\mathbf{J}(\mathbf{c}; \mathbf{d}) = (J_{ij})$, and so forth. We refer the reader to [13] for details.

Terms are in *slanted font* at their defining occurrence. We use bold font for matrices that may split into blocks, and for the sub-matrices. Individual entries of a matrix are in lowercase. For example, matrix \mathbf{A} has sub-matrices \mathbf{A}_{lm} and entries a_{ij} .

2.1 Block triangular forms of a sparsity pattern.

Let² $R = 1:n$ be the set of indices of n rows (equations), and let $C = 1:n$ be the set of indices of n columns (variables). A *sparsity pattern* \mathbf{A} is a subset of the Cartesian product $R \times C$ that contains row-column index pairs (i, j) . We can view \mathbf{A} as its incidence matrix (a_{ij}) , where a_{ij} equals 1 if $(i, j) \in \mathbf{A}$ and 0 otherwise. A *transversal* of \mathbf{A} is n positions in \mathbf{A} with exactly one position in each row and each column. If \mathbf{A} has some transversal, then it is *structurally nonsingular*. The union of all transversals of \mathbf{A} comprise its *essential sparsity pattern* \mathbf{A}_{ess} [11]. Obviously, \mathbf{A} is structurally nonsingular if and only if \mathbf{A}_{ess} is nonempty.

Assume henceforth that \mathbf{A} is structurally nonsingular. Let P and Q be two suitable permutation matrices for \mathbf{A} , such that the permuted incidence matrix $\mathbf{A}' = P\mathbf{A}Q$ can be written in a $p \times p$ block form

$$\mathbf{A}' = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \cdots & \mathbf{A}_{1p} \\ & \mathbf{A}_{22} & \cdots & \mathbf{A}_{2p} \\ & & \ddots & \vdots \\ & & & \mathbf{A}_{pp} \end{bmatrix}, \quad (2.1)$$

where each diagonal block \mathbf{A}_{qq} , $q = 1:p$, is square of positive size N_q . We say the block form (2.1) is a *BTF* of \mathbf{A} . Blanks in (2.1) mean that a sub-matrix \mathbf{A}_{kl} below the block diagonal with $k > l$ is empty.

A sparsity pattern is *irreducible*, if it cannot be permuted to the form (2.1) with $p > 1$ [2]; otherwise it is *reducible*. A BTF is *irreducible* if each diagonal block is *irreducible*; otherwise it is *reducible* [11]. Hence, if (2.1) is irreducible, then p is the largest number of diagonal blocks among all possible BTFs of \mathbf{A}' .

When we say block q of a matrix in a BTF, we shall refer to the q th diagonal block submatrix. For $q = 1:p$, we define for block q the index set

$$B_q = \text{the set of indices } i \text{ that belong to block } q.$$

Throughout this article, they are the indices of the permuted \mathbf{A}' , not those of the original \mathbf{A} .

Another useful notation is $\text{blockOf}(i)$ that denotes the block number q such that index $i \in B_q$. Since each diagonal block is square, both B_q and $\text{blockOf}(i)$ notation apply to rows and columns equally. To summarize, for $i \in 1:n$ and $q \in 1:p$,

$$\text{blockOf}(i) = q \iff i \in B_q \iff \sum_{m=1}^{q-1} N_m + 1 \leq i \leq \sum_{m=1}^q N_m.$$

²The colon notation $p:q$ for integers p, q denotes either the unordered set or the enumerated list of integers i with $p \leq i \leq q$, depending on context.

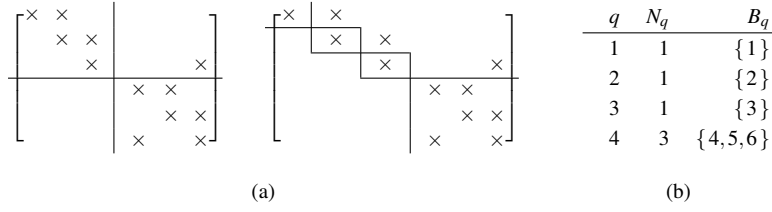


Fig. 2.1: (a) Two nontrivial BTFs of the same sparsity pattern. The left one is reducible with number of blocks $p = 2$. The right one is irreducible with $p = 4$. (b) Block information for the irreducible BTF.

Example 2.1 We illustrate in Figure 2.1 the above block notation with a sparsity pattern of two nontrivial BTFs.

The following lemma connects the transversals of a sparsity pattern \mathbf{A} and the transversals of its diagonal blocks in some BTF.

Lemma 2.1 [11, Lemma 2.4] *Any transversal T of a sparsity pattern \mathbf{A} is contained in the union of the diagonal blocks of any BTF of \mathbf{A} , that is, $T \subseteq \mathbf{A}_{11} \cup \dots \cup \mathbf{A}_{pp}$.*

Equivalently, the intersection of T with block q of \mathbf{A} is a transversal T_q of \mathbf{A}_{qq} .

2.2 Block triangular forms of a DAE.

The natural sparsity pattern of a DAE indicates if a variable x_j occurs in an equation f_i . Each such occurrence corresponds to a finite entry σ_{ij} in Σ , and hence we have

$$\mathbf{S} = \{ (i, j) \mid \sigma_{ij} > -\infty \} \quad (\text{the sparsity pattern of } \Sigma).$$

If \mathbf{S} has some transversal, then Σ has a transversal with finite σ_{ij} 's and a finite $\text{Val}(\Sigma)$ [9], so the DAE is structurally well posed (SWP) [10]; otherwise it is structurally ill posed. Here, we shall deal with the SWP case only.

A more informative BTF derives from the sparsity pattern $\mathbf{S}_0 = \mathbf{S}_0(\mathbf{c}; \mathbf{d})$ of a System Jacobian $\mathbf{J} = \mathbf{J}(\mathbf{c}; \mathbf{d})$ as defined in [13, (2.6)]:

$$\mathbf{S}_0 = \{ (i, j) \mid d_j - c_i = \sigma_{ij} \} \quad (\text{the sparsity pattern of } \mathbf{J}). \quad (2.2)$$

By [13, (2.2)], $d_j - c_i = \sigma_{ij}$ holds on a HVT T of Σ , so T is also a transversal of \mathbf{S}_0 .

A less obvious set contains the positions that contribute to $\det(\mathbf{J})$:

$$\mathbf{S}_{\text{ess}} = \text{the union of all HVTs of } \Sigma \quad (\text{the essential sparsity pattern of } \Sigma),$$

which is also the essential sparsity pattern of \mathbf{S}_0 for any valid offset pair $(\mathbf{c}; \mathbf{d})$ [11, Lemma 3.1].

Since $d_j - c_i = \sigma_{ij}$ holds on each HVT and hence implies $\sigma_{ij} > -\infty$, we have

$$\mathbf{S}_{\text{ess}} \subseteq \mathbf{S}_0 \subseteq \mathbf{S} \quad \text{for any offset pair } (\mathbf{c}; \mathbf{d}).$$

Our experience suggests that the (irreducible) BTF based on \mathbf{S}_0 can be significantly finer than that based on \mathbf{S} . We refer to the former BTF as *fine BTF*, and to the latter as *coarse BTF*. We refer to the diagonal blocks in the fine BTF as *fine blocks*, and refer to those in the coarse BTF as *coarse blocks*.

Assume that \mathbf{S}_0 is permuted into a $p \times p$ BTF. Following this BTF, we apply the same permutations on \mathbf{J} and $\mathbf{\Sigma}$, and write them in $p \times p$ block forms:

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_{11} & \mathbf{J}_{12} & \cdots & \mathbf{J}_{1p} \\ \mathbf{0} & \mathbf{J}_{22} & \cdots & \mathbf{J}_{2p} \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{J}_{pp} \end{bmatrix} \quad \text{and} \quad \mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_{11} & \mathbf{\Sigma}_{12} & \cdots & \mathbf{\Sigma}_{1p} \\ \mathbf{\Sigma}_{21} & \mathbf{\Sigma}_{22} & \cdots & \mathbf{\Sigma}_{2p} \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{\Sigma}_{p1} & \cdots & \cdots & \mathbf{\Sigma}_{pp} \end{bmatrix}. \quad (2.3)$$

We call this procedure a *block triangularization* of the DAE, and note that the sparsity pattern \mathbf{S} of $\mathbf{\Sigma}$ may not be in the same BTF as \mathbf{S}_0 of \mathbf{J} . That is, every σ_{ij} below the block diagonal of $\mathbf{\Sigma}$ is not necessarily $-\infty$, but must satisfy $\sigma_{ij} < d_j - c_i$ as $J_{ij} \equiv 0$. Hence,

$$d_j - c_i \begin{cases} > \sigma_{ij} & \text{if } \text{blockOf}(j) < \text{blockOf}(i) \\ \geq \sigma_{ij} & \text{if } \text{blockOf}(j) \geq \text{blockOf}(i). \end{cases} \quad (2.4)$$

We refer the reader to [6, 11] for more details on BTFs.

Example 2.2 We illustrate the coarse and fine BTFs with the (artificially) modified double pendula DAE in [7]. The state variables are x, y, λ, u, v, μ ; G is gravity, $L > 0$ is the length of both pendula, and α is a constant.

$$\begin{aligned} 0 = f_1 &= x'' + x\lambda & 0 = f_4 &= u'' + u\mu \\ 0 = f_2 &= y'' + y\lambda + (x')^3 - g & 0 = f_5 &= (v''')^3 + v\mu - G \\ 0 = f_3 &= x^2 + y^2 - L^2 & 0 = f_6 &= u^2 + v^2 - (L + \alpha\lambda)^2 + \lambda'' \end{aligned}$$

$$\mathbf{\Sigma} = \begin{array}{c} \begin{array}{cccccc} & v & \mu & u & x & y & \lambda & c_i \\ f_5 & 3^\bullet & 0 & & & & & 0 \\ f_4 & & 0^\bullet & 2 & & & & 0 \\ f_6 & 0 & & 0^\bullet & & & 2 & 2 \\ f_3 & & & & 0^\bullet & 0 & & 6 \\ f_2 & & & & 1 & 2^\bullet & 0 & 4 \\ f_1 & & & & 2 & & 0^\bullet & 4 \end{array} \\ \begin{array}{cccccc} d_j & 3 & 0 & 2 & 6 & 6 & 4 \end{array} \end{array} \quad \mathbf{J} = \begin{array}{c} \begin{array}{cccccc} & v''' & \mu & u'' & x^{(6)} & y^{(6)} & \lambda^{(4)} \\ f_5 & 2v''' & & & & & \\ f_4 & & v & & & & \\ f_6 & & u & 1 & & & \\ f_3^{(6)} & & & 2u & & & 1 \\ f_3^{(4)} & & & & 2x & 2y & \\ f_2^{(4)} & & & & & 1 & y \\ f_1^{(4)} & & & & & & 1 \\ f_1 & & & & & & x \end{array} \end{array}$$

The row and column labels in \mathbf{J} , showing equations and variables differentiated to order c_i and d_j , aim to remind the reader of the formula for \mathbf{J} in [13, (2.6)].

There are two 3×3 coarse blocks. The first one, comprising equations f_5, f_4, f_6 and variables v, μ, u , can further decompose into three 1×1 fine blocks, while the second coarse block, comprising equations f_3, f_2, f_1 and variables x, y, λ , is irreducible. Hence there are four blocks in the fine BTF.

The sparsity pattern \mathbf{S}_0 of \mathbf{J} is exactly the one in Figure 2.1(a), so the fine BTF information is in Figure 2.1(b).

If we state Lemma 2.1 in the context of a Jacobian sparsity pattern, then we have the following lemma.

Lemma 2.2 [11, Lemma 3.3] *Assume that a Jacobian sparsity pattern \mathbf{S}_0 is in some BTF. Let $(\boldsymbol{\Sigma}_{qm})_{q,m=1:p}$ be the corresponding sub-matrices of $\boldsymbol{\Sigma}$. Then a HVT T of $\boldsymbol{\Sigma}$ is the union of HVTs T_q of the diagonal blocks $\boldsymbol{\Sigma}_{qq}$: $T = T_1 \cup \dots \cup T_p$.*

This lemma is not difficult to prove, given that a transversal T of \mathbf{S}_0 is the union of transversals T_q of the diagonal blocks of \mathbf{S}_0 .

The following lemma is useful for proving the main Theorems 3.1 and 3.2 of the block conversion methods in §3.

Lemma 2.3 *Assume that $\boldsymbol{\Sigma}$ has a finite $\text{Val}(\boldsymbol{\Sigma})$ and is in a $p \times p$ block form as in (2.3). Let \mathbf{c} and \mathbf{d} be two nonnegative integer n -vectors. Assume also that*

- (a) $d_j - c_i > \sigma_{ij}$ holds for all entries below the diagonal blocks of $\boldsymbol{\Sigma}$,
- (b) $d_j - c_i \geq \sigma_{ij}$ holds elsewhere, and
- (c) $\text{Val}(\boldsymbol{\Sigma}) = \sum_j d_j - \sum_i c_i$.

Then

- (i) $(\mathbf{c}; \mathbf{d})$ is a valid offset pair of $\boldsymbol{\Sigma}$,
- (ii) the block form of $\boldsymbol{\Sigma}$ is a BTF of the Jacobian sparsity pattern \mathbf{S}_0 , and
- (iii) a HVT of $\boldsymbol{\Sigma}$ is the union of HVTs T_q of the diagonal blocks $\boldsymbol{\Sigma}_{qq}$, for all $q = 1:p$.

Proof (i) We let T denote a HVT of $\boldsymbol{\Sigma}$. Since $\text{Val}(\boldsymbol{\Sigma})$ is finite, $\sigma_{ij} \geq 0$ for all $(i, j) \in T$. For $(\mathbf{c}; \mathbf{d})$ to be a valid offset of $\boldsymbol{\Sigma}$, $d_j - c_i \geq \sigma_{ij}$ must hold for all $i, j = 1:n$, with equalities for all $(i, j) \in T$ [9].

By (a) and (b), $d_j - c_i \geq \sigma_{ij}$ holds everywhere. Summing these inequalities over T gives

$$\sum_{(i,j) \in T} (d_j - c_i) \geq \sum_{(i,j) \in T} \sigma_{ij}.$$

The left-hand side equals $\sum_j d_j - \sum_i c_i$, and the right-hand side equals $\text{Val}(\boldsymbol{\Sigma})$ by definition. By (c), these two values are equal, so $d_j - c_i = \sigma_{ij}$ holds for all $(i, j) \in T$, and $(\mathbf{c}; \mathbf{d})$ is valid for $\boldsymbol{\Sigma}$.

(ii) By (a), the blocks below the block diagonal in \mathbf{S}_0 , derived from $\boldsymbol{\Sigma}$ and $(\mathbf{c}; \mathbf{d})$ using (2.2), are empty. By the definition of a BTF of a Jacobian sparsity pattern, \mathbf{S}_0 is in a BTF as described by the $p \times p$ block form.

(iii) This follows immediately from (ii) and Lemma 2.2. \square

Following a $p \times p$ BTF based on \mathbf{S}_0 , we can write any valid offset pair $(\mathbf{c}; \mathbf{d})$ of $\boldsymbol{\Sigma}$ in a block form as

$$(\mathbf{c}_1; \mathbf{d}_1), (\mathbf{c}_2; \mathbf{d}_2), \dots, (\mathbf{c}_p; \mathbf{d}_p), \quad (2.5)$$

where each of the sub-vectors \mathbf{c}_q and \mathbf{d}_q is of length N_q , where $q = 1:p$.

Lemma 2.4 *Assume that a Jacobian pattern \mathbf{S}_0 , derived by Σ and a valid offset pair $(\mathbf{c}; \mathbf{d})$, is in some BTF. If we write $(\mathbf{c}; \mathbf{d})$ into block form as in (2.5), then $(\mathbf{c}_q; \mathbf{d}_q)$ is a valid offset pair of Σ_{qq} .*

Proof Let T be a HVT of Σ . By Lemma 2.2, the intersection of T with block q is a HVT T_q of Σ_{qq} . Then $d_j - c_i = \sigma_{ij}$ holds for all $(i, j) \in T_q \subseteq T$. Since $(\mathbf{c}; \mathbf{d})$ is valid for Σ , $d_j - c_i \geq \sigma_{ij}$ and $c_i \geq 0$ hold on Σ_{qq} , where $i, j \in B_q$. Thus the offset pair $(\mathbf{c}_q; \mathbf{d}_q)$ matched to block q satisfies the conditions [13, (2.2)] for being valid for Σ_{qq} . \square

From the view of Lemma 2.4, we can regard each diagonal block Σ_{qq} as a signature matrix in its own right. Equivalently, each block q , having N_q equations in N_q variables, can be viewed as a sub-DAE, with a signature matrix Σ_{qq} , a finite $\text{Val}(\Sigma_{qq})$, a local offset pair $(\mathbf{c}_q; \mathbf{d}_q)$, and a sub-Jacobian \mathbf{J}_{qq} . Expressions that contribute to entries in an off-diagonal block Σ_{qm} , where $q \neq m$, can be considered as driving terms, or equivalently, the influence of variables in block m on those in block q . We refer to $(\mathbf{c}; \mathbf{d})$ of Σ as a global offset pair. The reader is referred to [11] for more theoretical results about block triangularization and global/local offset pairs.

3 Block conversion methods.

They are suitable for improving the efficiency of finding a useful conversion for fixing SA's failures. If \mathbf{J} is identically singular, then by (2.3), $\det(\mathbf{J}) = \prod_{q=1}^p \det(\mathbf{J}_{qq}) \equiv 0$, so at least one \mathbf{J}_{qq} for some $q \in 1:p$ is identically singular. As discussed before, we can regard block q as a sub-DAE with a signature matrix Σ_{qq} . Then we wish to apply the basic conversion methods on this sub-DAE to achieve a strict decrease in $\text{Val}(\Sigma_{qq})$, provided the conditions for applying these methods are satisfied for those variables and equations within block q .

However, what we should ensure is a strict decrease in the value of the whole signature matrix, namely $\text{Val}(\bar{\Sigma}) < \text{Val}(\Sigma)$, where $\bar{\Sigma}$ is the signature matrix of the resulting DAE. Proving this inequality from a decrease in $\text{Val}(\Sigma_{qq})$ is nontrivial, because a conversion on block q may affect blocks Σ_{qm} for $m = 1, \dots, q-1, q+1, \dots, p$. Especially in the ES method, Σ_{qq} and these blocks are enlarged. Hence, the conditions and the conversion process need to be carefully modified, so that the conversion methods can adapt to a BTF based on \mathbf{S}_0 .

We give an introductory example in §3.1, present the block LC method in §3.2, and present the block ES method in §3.3.

Hereafter we use the fine BTF in the examples for demonstration, since each fine block contains an irreducible sub-Jacobian sparsity pattern. Our experience suggests that a useful conversion can usually be derived from the fine BTF of a DAE. However, we emphasize that the block conversion methods can be applied not only to the irreducible BTF of a Jacobian sparsity pattern \mathbf{S}_0 with some valid $(\mathbf{c}; \mathbf{d})$, but also to any BTF of \mathbf{S}_0 . For example, the basic conversion methods consider a DAE in a (trivial) BTF of one $n \times n$ block.

3.1 An introductory example.

We illustrate these block methods with the following DAE:

$$\begin{aligned} 0 &= f_1 = x_1 + x_2 + h_1(t) \\ 0 &= f_2 = x_1 + (x'_1 + x'_2)x'_3 + h_2(t) \\ 0 &= f_3 = x'_3 + h_3(t). \end{aligned} \quad (3.1)$$

$$\mathbf{\Sigma} = \begin{array}{c} f_1 \\ f_2 \\ f_3 \end{array} \left[\begin{array}{ccc|c} x_1 & x_2 & x_3 & c_i \\ \hline 0 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ \hline & & 1 & 0 \end{array} \right] \begin{array}{c} \\ \\ \\ d_j \\ 1 \\ 1 \\ 1 \end{array} \quad \mathbf{J} = \begin{array}{c} f'_1 \\ f'_2 \\ f'_3 \end{array} \left[\begin{array}{cc|c} x'_1 & x'_2 & x'_3 \\ \hline 1 & 1 & x'_1 + x'_2 \\ \hline x'_3 & x'_3 & 1 \end{array} \right]$$

(Here h_1, h_2, h_3 are driving functions.) The coarse BTF and the fine BTF are identical, both having two diagonal blocks.

In the basic LC method, we can choose $\mathbf{u} = [-x'_3, 1, -x'_1 - x'_2]^T \in \text{coker}(\mathbf{J})$. Using [13, (4.3)], we have

$$I = \{i \mid u_i \neq 0\} = \{1, 2, 3\}, \quad \underline{c} = \min_{i \in I} c_i = 0, \quad L = \{l \in I \mid c_l = \underline{c}\} = \{2, 3\}.$$

We let $\sigma(x_j, \mathbf{u})$ denote the order of the highest derivative to which x_j occurs in \mathbf{u} , or $-\infty$ if x_j does not occur in \mathbf{u} [13, (4.1)]. The LC condition [13, (4.4)] is violated since

$$\sigma(x_j, \mathbf{u}) = 1 \not\leq 1 = d_j - \underline{c} \quad \text{for all } j = 1:3.$$

Not surprisingly, replacing either f_2 or f_3 by

$$\bar{f} = \sum_{i \in I} u_i f_i^{(c_i - \underline{c})} = -x'_3 h'_1(t) + (x_1 + h_2(t)) - (x'_1 + x'_2)(x'_3 + h_3(t))$$

does not result in a decrease in $\text{Val}(\mathbf{\Sigma})$; verifying this is not difficult.

Notice that only the sub-Jacobian of block 1, $\mathbf{J}_{11} = \partial(f'_1, f_2)/\partial(x'_1, x'_2)$, is singular. Suppose we consider block 1, with $B_1 = \{1, 2\}$, as a sub-DAE, and choose $\mathbf{u} = [-x'_3, 1]^T \in \text{coker}(\mathbf{J}_{11})$. Within block 1, the LC method derives

$$I = \{i \in B_1 \mid u_i \neq 0\} = \{1, 2\}, \quad \underline{c} = \min_{i \in I} c_i = 0, \quad L = \{l \in I \mid c_l = \underline{c}\} = \{2\}.$$

Now the LC condition [13, (4.4)] is satisfied for the column indices in block 1:

$$\sigma(x_j, \mathbf{u}) = -\infty < d_j - \underline{c} \quad \text{for } j = 1, 2 \in B_1.$$

Replacing f_2 by $\bar{f}_2 = u_1 f'_1 + u_2 f_2 = x_1 + h_2(t) - x'_3 h'_1(t)$ results in the DAE with the following SA result.

$$\bar{\Sigma} = \begin{array}{c} \begin{array}{cccc|c} & x_2 & x_1 & x_3 & c_i \\ f_1 & \mathbf{0}^\bullet & 0 & & 0 \\ f_2 & & \mathbf{0}^\bullet & 1 & 0 \\ f_3 & & & \mathbf{1}^\bullet & 0 \end{array} \\ d_j \quad 0 \quad 0 \quad 1 \end{array} \quad \bar{\mathbf{J}} = \begin{array}{c} \begin{array}{ccc|c} & x_2 & x_1 & x'_3 \\ f_1 & 1 & 1 & \\ f_2 & & 1 & g'_1(t) \\ f_3 & & & 1 \end{array} \end{array}$$

The SA succeeds as $\bar{\mathbf{J}}$ is nonsingular. The conversion results in a decrease in the value of the signature matrix: $\text{Val}(\bar{\Sigma}) = 1 < 2 = \text{Val}(\Sigma)$.

The basic ES method can work on (3.1) by choosing $\mathbf{v} = [1, -1, 0]^T \in \ker(\mathbf{J})$. It is simpler—though trivial for this example—to work on block 1 only. We find $\mathbf{v} = [1, -1]^T \in \ker(\mathbf{J}_{11})$, and use [13, (4.10)] to derive

$$J = \{l \in B_1 \mid v_l \neq 0\} = \{1, 2\}, \quad s = |J| = 2, \quad M = \{1, 2\}, \quad \bar{c} = \max_{i \in M} c_i = 1.$$

Since \mathbf{v} is constant, it is not difficult to verify that the ES conditions [13, (4.11)] hold.

We choose $l = 2 \in J$ and introduce for x_1 a new variable

$$y_1 = x_1^{(d_1 - \bar{c})} - \frac{v_1}{v_2} \cdot x_2^{(d_2 - \bar{c})} = x_1 + x_2.$$

The ES method hence says: replace x_1 by $y_1 - x_2$ in f_1 , and replace x'_1 by $y'_1 - x'_2$ in f_2 . Finally we append the equation g_1 that prescribes such replacements, and obtain

$$\begin{array}{l} 0 = \bar{f}_1 = y_1 + h_1(t) \\ 0 = \bar{f}_2 = x_1 + y'_1 x'_3 + h_2(t) \end{array} \quad \begin{array}{l} 0 = \bar{f}_3 = x'_3 + h_3(t) \\ 0 = g_1 = -y_1 + x_1 + x_2. \end{array}$$

$$\bar{\Sigma} = \begin{array}{c} \begin{array}{cccc|c} & x_2 & x_1 & y_1 & x_3 & c_i \\ g_1 & \mathbf{0}^\bullet & 0 & 0 & & 0 \\ f_2 & & \mathbf{0}^\bullet & 1 & 1 & 0 \\ f_1 & & & \mathbf{0}^\bullet & & 1 \\ f_3 & & & & \mathbf{1}^\bullet & 0 \end{array} \\ d_j \quad 0 \quad 0 \quad 1 \quad 1 \end{array} \quad \bar{\mathbf{J}} = \begin{array}{c} \begin{array}{cccc|c} & x_2 & x_1 & y'_1 & x'_3 \\ g_1 & 1 & 1 & -1 & \\ f_2 & & 1 & x'_3 & y'_1 \\ f_1 & & & 1 & \\ f_3 & & & & 1 \end{array} \end{array}$$

Again $\text{Val}(\bar{\Sigma}) = 1 < 2 = \text{Val}(\Sigma)$, and the SA succeeds as $\det(\bar{\mathbf{J}}) = 1$.

3.2 Block linear combination method.

We first introduce some convenient notation for the block LC method. Let $\mathbf{0}_r$ denote the zero column vector of size r . Assume that a \mathbf{J}_{qq} is identically singular. Let $\hat{\mathbf{u}} \in \text{coker}(\mathbf{J}_{qq})$, where $\hat{\mathbf{u}} \neq \mathbf{0}_{N_q}$. Let also

$$\mathbf{u} = \begin{bmatrix} \mathbf{0}_{N_1 + \dots + N_{q-1}} \\ \hat{\mathbf{u}} \\ \mathbf{0}_{N_{q+1} + \dots + N_p} \end{bmatrix}.$$

Denote

$$\begin{aligned} I &= \{i \mid u_i \neq 0\} \subseteq B_q, \quad \underline{c} = \min_{i \in I} c_i, \\ L &= \{l \in I \mid c_l = \underline{c}\}, \quad \text{and} \quad \bar{L} = \{l \in L \mid u_l \text{ is (nonzero) constant}\}. \end{aligned} \quad (3.2)$$

The set \bar{L} is used to seek a conversion that guarantees equivalence between the original DAE and the converted one. The block LC method is based on the following theorem.

Theorem 3.1 *If*

$$\sigma(x_j, \mathbf{u}) < d_j - \underline{c} \quad \text{for all } j \in B_q \quad (3.3)$$

and we replace an equation $f_l, l \in L$, by

$$\bar{f} = \sum_{i \in I} u_i f_i^{(c_i - \underline{c})},$$

then $\text{Val}(\bar{\Sigma}) < \text{Val}(\Sigma)$, where $\bar{\Sigma} = (\bar{\sigma}_{ij})$ is the signature matrix of the resulting DAE.

Before proving this theorem, we show how to apply the block LC method and prove a related lemma.

Example 3.1 We illustrate the block LC method with the Campbell-Griepentrog two-link robot arm DAE [1]. We slightly simplify the problem formulation to (3.4), allowing the first-order derivatives x'_1, x'_2 , and x'_3 to occur implicitly in the equations. The two state variables u_1 and u_2 in the original formulation are renamed x_4 and x_5 , respectively (and not to be confused with entries in a vector \mathbf{u} in our notation).

The equations of this problem are

$$\begin{aligned} 0 = A &= x''_1 - \left[2c(x_3)(x'_1 + x'_3)^2 + x_1^2 d(x_3) + (2x_3 - x_2)(a(x_3) + 2b(x_3)) \right. \\ &\quad \left. + a(x_3)(x_4 - x_5) \right] \\ 0 = B &= x''_2 - \left[-2c(x_3)(x'_1 + x'_3)^2 - x_1^2 d(x_3) + (2x_3 - x_2)(1 - 3a(x_3) - 2b(x_3)) \right. \\ &\quad \left. - a(x_3)x_4 + (a(x_3) + 1)x_5 \right] \\ 0 = C &= x''_3 - \left[-2c(x_3)(x'_1 + x'_3)^2 - x_1^2 d(x_3) + (2x_3 - x_2)(a(x_3) - 9b(x_3)) \right. \\ &\quad \left. - 2x_1^2 c(x_3) - d(x_3)(x'_1 + x'_3)^2 - (a(x_3) + b(x_3))(x_4 - x_5) \right] \\ 0 = D &= \cos x_1 + \cos(x_1 + x_3) - p_1(t) \\ 0 = E &= \sin x_1 + \sin(x_1 + x_3) - p_2(t), \end{aligned} \quad (3.4)$$

where

$$\begin{aligned} a(\theta) &= 2/(2 - \cos^2 \theta) & b(\theta) &= \cos \theta / (2 - \cos^2 \theta) \\ c(\theta) &= \sin \theta / (2 - \cos^2 \theta) & d(\theta) &= \sin \theta \cos \theta / (2 - \cos^2 \theta) \\ p_1(t) &= \cos(1 - e^t) + \cos(1 - t) & p_2(t) &= \sin(1 - e^t) + \sin(1 - t). \end{aligned}$$

$$\mathbf{\Sigma} = \begin{array}{c} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ d_j \end{array} \begin{array}{c} B \\ C \\ A \\ D \\ E \\ 2 \end{array} \begin{array}{c} x_2 \\ x_4 \\ x_5 \\ x_1 \\ x_3 \\ c_i \end{array} \begin{array}{c} \left[\begin{array}{ccc|cc} 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 2 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 2 & 2 \end{array} \right] \end{array} \mathbf{J} = \begin{array}{c} B \\ C \\ A \\ D'' \\ E'' \end{array} \begin{array}{c} x_2'' \\ x_4 \\ x_5 \\ x_1'' \\ x_3'' \end{array} \begin{array}{c} \left[\begin{array}{ccc|cc} 1 & a_3 & -a_3 - 1 & & \\ a_3 + b_3 & -a_3 - b_3 & & & 1 \\ -a_3 & a_3 & & 1 & \\ \hline & & & \frac{\partial D}{\partial x_1} & \frac{\partial D}{\partial x_3} \\ & & & \frac{\partial E}{\partial x_1} & \frac{\partial E}{\partial x_3} \end{array} \right] \end{array}$$

Here in \mathbf{J} ,

$$\begin{aligned} a_3 &= a(x_3) = 2/(2 - \cos^2 x_3) & b_3 &= b(x_3) = \cos x_3/(2 - \cos^2 x_3) \\ \partial D/\partial x_1 &= -\sin x_1 - \sin(x_1 + x_3) & \partial D/\partial x_3 &= -\sin(x_1 + x_3) \\ \partial E/\partial x_1 &= \cos x_1 + \cos(x_1 + x_3) & \partial E/\partial x_3 &= \cos(x_1 + x_3). \end{aligned}$$

The DAE (3.4) is of differentiation index 5, while the SA reports structural index $v_S = 3$. Hence this must be a failure case, because v_S is an upper bound for the differentiation index when the SA succeeds [9]. We can see that the sub-Jacobian \mathbf{J}_{22} of block 2 is identically singular.

Our method first computes $\hat{\mathbf{u}} = [2, 2 + \cos x_3]^T \in \text{coker}(\mathbf{J}_{22})$. Then $\mathbf{u} = [0, 2, 2 + \cos x_3, 0, 0]^T$. Using (3.2), we have

$$I = \{i \mid u_i \neq 0\} = \{2, 3\}, \quad \underline{c} = \min_{i \in I} c_i = 0, \quad L = \{2, 3\}, \quad \text{and} \quad \bar{L} = \{2\}.$$

The variables x_4 and x_5 in block 2 do not occur in \mathbf{u} , so the condition (3.3) is satisfied.

Considering equivalence, we pick $l = 2 \in \bar{L}$ over $l = 3 \in L \setminus \bar{L}$, and replace $f_l = C$ by $\bar{C} = u_1 C + u_2 A = 2C + (2 + \cos x_3)A$. The SA results of the resulting DAE are as follows.

$$\bar{\mathbf{\Sigma}} = \begin{array}{c} A \\ B \\ \bar{C} \\ D \\ E \\ d_j \end{array} \begin{array}{c} x_4 \\ x_5 \\ x_2 \\ x_1 \\ x_3 \\ c_i \end{array} \begin{array}{c} \left[\begin{array}{ccc|cc} 0 & 0 & 0 & 2 & 1 \\ 0 & 0 & 2 & 1 & 1 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 4 & 4 \end{array} \right] \end{array} \bar{\mathbf{J}} = \begin{array}{c} A \\ B \\ \bar{C}'' \\ D^{(4)} \\ E^{(4)} \end{array} \begin{array}{c} x_4 \\ x_5 \\ x_2'' \\ x_1^{(4)} \\ x_3^{(4)} \end{array} \begin{array}{c} \left[\begin{array}{ccc|cc} a_3 & a_3 & & & \\ a_3 & -a_3 - 1 & & 1 & \\ \hline & & & \frac{\partial \bar{C}}{\partial x_2} & \\ & & & 2 + \cos x_3 & 2 \\ \hline & & & \frac{\partial D}{\partial x_1} & \frac{\partial D}{\partial x_3} \\ & & & \frac{\partial E}{\partial x_1} & \frac{\partial E}{\partial x_3} \end{array} \right] \end{array}$$

Here $\partial \bar{C}/\partial x_2 = 2(a_3^2 - 3a_3 b_3 + b_3^2)(2 - \cos^2 x_3)$. The SA reports $v_S = 5$ and succeeds at any point where

$$\det(\bar{\mathbf{J}}) = 4(a_3^2 - 3a_3 b_3 + b_3^2) \sin x_3 \neq 0.$$

Now $\text{Val}(\bar{\mathbf{\Sigma}}) = 0 < 2 = \text{Val}(\mathbf{\Sigma})$.

Lemma 3.1 Consider a BTF of a Jacobian pattern \mathbf{S}_0 derived from $\mathbf{\Sigma}$ and $(\mathbf{c}; \mathbf{d})$. If we perform the LC conversion as described in Theorem 3.1, then in the resulting $\bar{\mathbf{\Sigma}}$,

$$d_j - c_i \begin{cases} > \bar{\sigma}_{ij} & \text{if } \text{blockOf}(j) < \text{blockOf}(i) \\ \geq \bar{\sigma}_{ij} & \text{if } \text{blockOf}(j) \geq \text{blockOf}(i). \end{cases} \quad (3.5)$$

Proof We only replace f_l by $\bar{f}_l = \bar{f}$ in a conversion, so $\bar{\sigma}_{ij} = \sigma_{ij}$ for all $i \neq l$ and all j . By (2.4), (3.5) holds for all $i \neq l$.

When $i = l$, we consider two cases: (a) $\text{blockOf}(j) < q$ and (b) $\text{blockOf}(j) \geq q$.

(a) $\text{blockOf}(j) < q = \text{blockOf}(l)$. By (2.4), $\sigma_{lj} < d_j - c_l$. Then $\bar{\sigma}_{lj}$ is

$$\sigma(x_j, \bar{f}_l) = \sigma\left(x_j, \sum_{i \in I} u_i f_i^{(c_i - \underline{c})}\right) \leq \max\left\{\sigma(x_j, \mathbf{u}), \max_{i \in I} \sigma(x_j, f_i^{(c_i - \underline{c})})\right\}. \quad (3.6)$$

We use some simple derivations to obtain

$$\begin{aligned} \sigma(x_j, \mathbf{u}) &\leq \sigma(x_j, \mathbf{J}_{qq}) \leq \max_{i \in I} \sigma(x_j, f_i) = \max_{i \in I} \sigma_{ij} \\ &< \max_{i \in I} (d_j - c_i) = d_j - \min_{i \in I} c_i = d_j - c_l \quad \text{and} \end{aligned} \quad (3.7a)$$

$$\max_{i \in I} \sigma\left(x_j, f_i^{(c_i - \underline{c})}\right) = \max_{i \in I} (\sigma_{ij} + c_i - \underline{c}) < d_j - \underline{c} = d_j - c_l. \quad (3.7b)$$

Substituting (3.7a) and (3.7b) in (3.6), we obtain $\bar{\sigma}_{lj} = \sigma(x_j, \bar{f}_l) < d_j - c_l$.

(b) $\text{blockOf}(j) \geq q = \text{blockOf}(l)$. By (2.4), $\sigma_{lj} \leq d_j - c_l$. We can replace the two “<” in (3.7) by “ \leq ”, and using these inequalities in (3.6), we have $\bar{\sigma}_{lj} \leq d_j - c_l$. \square

Using Lemma 3.1, we can now prove Theorem 3.1.

Proof By Lemma 2.4, we can regard block q as a sub-DAE with $\mathbf{\Sigma}_{qq}$ and $(\mathbf{c}_q; \mathbf{d}_q)$. The conversion described in Theorem 3.1 can be considered as an application of the basic LC method to this sub-DAE. Since the block LC condition (3.3) holds, that is, $\sigma(x_j, \mathbf{u}) < d_j - \underline{c}$ for all $j \in B_q$ that belong to this sub-DAE, the basic LC condition [13, (4.4)] also holds for the sub-DAE. Hence $\text{Val}(\bar{\mathbf{\Sigma}}_{qq}) < \text{Val}(\mathbf{\Sigma}_{qq})$.

Let \bar{T} be a HVT of $\bar{\mathbf{\Sigma}}$. Using (3.5) in Lemma 3.1, we have $\text{Val}(\bar{\mathbf{\Sigma}}) = \sum_{(i,j) \in \bar{T}} \bar{\sigma}_{ij} \leq d_j - c_i = \text{Val}(\mathbf{\Sigma})$. Now we prove $\text{Val}(\bar{\mathbf{\Sigma}}) < \text{Val}(\mathbf{\Sigma})$ by contradiction. First assume that $\text{Val}(\bar{\mathbf{\Sigma}}) = \sum_j d_j - \sum_i c_i = \text{Val}(\mathbf{\Sigma}) \geq 0$. With (3.5), the three conditions in Lemma 2.3 are satisfied. From this lemma, it follows that the Jacobian patterns $\bar{\mathbf{S}}_0$, derived from $\bar{\mathbf{\Sigma}}$ and $(\mathbf{c}; \mathbf{d})$, and \mathbf{S}_0 , derived from $\mathbf{\Sigma}$ and $(\mathbf{c}; \mathbf{d})$, are in the same $p \times p$ BTF.

By Lemma 2.2, \bar{T} is the union of HVTs \bar{T}_m of all diagonal blocks $\bar{\mathbf{\Sigma}}_{mm}$, $m = 1 : p$. By the construction of $\bar{\mathbf{\Sigma}}$, $\text{Val}(\bar{\mathbf{\Sigma}}_{mm}) = \text{Val}(\mathbf{\Sigma}_{mm})$ for all $m \neq q$. Then a contradiction follows from

$$\begin{aligned} \text{Val}(\bar{\mathbf{\Sigma}}) &= \sum_{(i,j) \in \bar{T}} \bar{\sigma}_{ij} = \sum_{m=1}^p \sum_{(i,j) \in \bar{T}_m} \bar{\sigma}_{ij} = \sum_{m=1}^p \text{Val}(\bar{\mathbf{\Sigma}}_{mm}) \\ &= \sum_{m \neq q} \text{Val}(\bar{\mathbf{\Sigma}}_{mm}) + \text{Val}(\bar{\mathbf{\Sigma}}_{qq}) < \sum_{m \neq q} \text{Val}(\mathbf{\Sigma}_{mm}) + \text{Val}(\mathbf{\Sigma}_{qq}) \\ &= \sum_{m=1}^p \text{Val}(\mathbf{\Sigma}_{mm}) = \sum_{m=1}^p \sum_{(i,j) \in T_m} \sigma_{ij} = \sum_{(i,j) \in T} \sigma_{ij} = \text{Val}(\mathbf{\Sigma}), \end{aligned} \quad (3.8)$$

where T is a HVT of $\mathbf{\Sigma}$ and T_m are HVTs of its diagonal blocks $\mathbf{\Sigma}_{mm}$. The assumption $\text{Val}(\bar{\mathbf{\Sigma}}) = \text{Val}(\mathbf{\Sigma})$ is hence false, so $\text{Val}(\bar{\mathbf{\Sigma}}) < \text{Val}(\mathbf{\Sigma})$ holds.

3.3 Block expression substitution method.

Assume again that a \mathbf{J}_{qq} is identically singular. Let $\widehat{\mathbf{v}} \in \ker(\mathbf{J}_{qq})$, where $\widehat{\mathbf{v}} \neq \mathbf{0}_{N_q}$. Similarly, we construct the column n -vector

$$\mathbf{v} = \begin{bmatrix} \mathbf{0}_{N_1+\dots+N_{q-1}} \\ \widehat{\mathbf{v}} \\ \mathbf{0}_{N_{q+1}+\dots+N_p} \end{bmatrix}.$$

We use notation similar to that used in the basic ES method (see [13, §4.2]):

$$\begin{aligned} J &= \{j \mid v_j \neq 0\} \subseteq B_q, & M &= \{i \in B_q \mid d_j - c_i = \sigma_{ij} \text{ for some } j \in J\}, \\ s &= |J|, & \bar{c} &= \max_{i \in M} c_i & \text{ and } & \bar{J} &= \{l \mid v_l \text{ is (nonzero) constant}\}. \end{aligned} \quad (3.9)$$

The set \bar{J} is used to seek a conversion that guarantees equivalence between the original DAE and the converted one. The conditions for applying the block ES method are

$$\begin{aligned} \sigma(x_j, \mathbf{v}) &\begin{cases} < d_j - \bar{c} & \text{if } j \in J \text{ or } \text{blockOf}(j) < q \\ \leq d_j - \bar{c} & \text{if } j \in B_q \setminus J \text{ or } \text{blockOf}(j) > q, \end{cases} & \text{and} & \\ d_j - \bar{c} &\geq 0 & \text{for all } j \in J. \end{aligned} \quad (3.10)$$

We choose an $l \in J$, and introduce $s - 1$ new variables

$$y_j = x_j^{(d_j - \bar{c})} - \frac{v_j}{v_l} \cdot x_l^{(d_l - \bar{c})} \quad \text{for all } j \in J \setminus \{l\}. \quad (3.11)$$

In each f_i with $i \in B_q$, we

$$\begin{aligned} &\text{replace each } x_j^{(\sigma_{ij})} & \text{with } d_j - c_i = \sigma_{ij} \text{ and } j \in J \setminus \{l\} \\ &\text{by } \left(y_j + \frac{v_j}{v_l} \cdot x_l^{(d_l - \bar{c})}\right)^{(\bar{c} - c_i)}. \end{aligned} \quad (3.12)$$

Note that because of M in (3.9), we actually perform replacements (equivalently referred to as ‘‘expression substitutions’’) in only f_i ’s with $i \in M \subseteq B_q$. Denote each new f_i by \bar{f}_i , and let also $\bar{f}_i = f_i$ for the unchanged equations with $i \notin M$.

By (3.11), we append $s - 1$ equations that prescribe the substitutions in (3.12):

$$\mathbf{0} = \mathbf{g}_j = -y_j + x_j^{(d_j - \bar{c})} - \frac{v_j}{v_l} \cdot x_l^{(d_l - \bar{c})} \quad \text{for all } j \in J \setminus \{l\}. \quad (3.13)$$

The block ES method is based on the following theorem.

Theorem 3.2 *Let J , s , M , and \bar{c} be as defined in (3.9). Assume that the conditions (3.10) hold. For an $l \in J$, if we*

- 1) *introduce $s - 1$ new variables x_j , $j \in J \setminus \{l\}$, as defined in (3.11),*
- 2) *perform replacements in f_i , for all $i \in B_q$, as described in (3.12), and*
- 3) *append $s - 1$ equations g_j , $j \in J \setminus \{l\}$, as defined in (3.13),*

then $\text{Val}(\bar{\Sigma}) < \text{Val}(\Sigma)$, where $\bar{\Sigma}$ is the signature matrix of the resulting DAE.

Before proving this theorem, we illustrate the block ES method with the robot arm DAE (3.4) and prove two related lemmas.

Example 3.2 The method finds first $\widehat{\mathbf{v}} = [1, 1]^T \in \ker(\mathbf{J}_{22})$. Then $\mathbf{v} = [0, 0, 1, 1, 0]^T$. Using (3.9), we have

$$J = \bar{J} = \{j \mid v_j \neq 0\} = \{2, 3\}, \quad s = |J| = 2, \quad M = \{2, 3\}, \quad \bar{c} = \max_{i \in M} c_i = 0.$$

Since \mathbf{v} is constant, $J = \bar{J}$ and the first condition in (3.10) holds. The second condition in (3.10) holds also, as $d_4 - \bar{c} = d_5 - \bar{c} = 0$. We choose x_4 , whose column index in the permuted Σ is $l = 2 \in \bar{J}$. Then we introduce for x_5 , the other variable in block 2 with column index $j = 3$, a new variable

$$y_5 = x_5^{(d_5 - \bar{c})} - \frac{v_3}{v_2} \cdot x_4^{(d_4 - \bar{c})} = x_5 - x_4.$$

Correspondingly, we append $0 = g_5 = -y_5 + x_5 - x_4$ and replace x_5 by $y_5 + x_4$ in \mathbf{C} and \mathbf{A} , the equations in block 2.

The resulting DAE has the following new equations

$$\begin{aligned} 0 = \bar{A} &= x_1'' - \left[2c(x_3)(x_1' + x_3')^2 + x_1'^2 d(x_3) + (2x_3 - x_2)(a(x_3) + 2b(x_3)) + a(x_3)y_5 \right] \\ 0 = \bar{C} &= x_3'' - \left[-2c(x_3)(x_1' + x_3')^2 - x_1'^2 d(x_3) + (2x_3 - x_2)(a(x_3) - 9b(x_3)) \right. \\ &\quad \left. - 2x_1'^2 c(x_3) - d(x_3)(x_1' + x_3')^2 - (a(x_3) + b(x_3))y_5 \right] \\ 0 = g_5 &= -y_5 + x_4 - x_5. \end{aligned}$$

$$\bar{\Sigma} = \begin{array}{c} \begin{array}{cccccc} & x_4 & x_5 & x_2 & y_5 & x_1 & x_3 & c_i \\ g_5 & 0^\bullet & 0 & & 0 & & & 0 \\ B & 0 & 0^\bullet & 2 & & 1 & 1 & 0 \\ \bar{C} & & & 0^\bullet & 0 & 1 & 2 & 2 \\ \bar{A} & & & 0 & 0^\bullet & 2 & 1 & 2 \\ D & & & & & 0^\bullet & 0 & 4 \\ E & & & & & 0 & 0^\bullet & 4 \\ d_j & 0 & 0 & 2 & 2 & 4 & 4 & \end{array} \end{array} \quad \bar{\mathbf{J}} = \begin{array}{c} \begin{array}{cccccc} & x_4 & x_5 & x_2'' & y_5'' & x_1^{(4)} & x_3^{(4)} \\ g_5 & -1 & 1 & & & & \\ B & a_3 & -a_3 - 1 & 1 & & & \\ \bar{C}'' & & & a_3 - 9b_3 & -a_3 - b_3 & & 1 \\ \bar{A}'' & & & a_3 + 2b_3 & a_3 & & 1 \\ D^{(4)} & & & & & \frac{\partial D}{\partial x_1} & \frac{\partial D}{\partial x_3} \\ E^{(4)} & & & & & \frac{\partial E}{\partial x_1} & \frac{\partial E}{\partial x_3} \end{array} \end{array}$$

Now the System Jacobian $\bar{\mathbf{J}}$ is generically nonsingular. The SA reports the correct index 5, and succeeds at any point where $\det(\bar{\mathbf{J}}) = 2(a_3^2 - 3a_3b_3 + b_3^2) \sin x_3 \neq 0$. Again $\text{Val}(\bar{\Sigma}) = 0 < 2 = \text{Val}(\Sigma)$.

In [8], Pryce fixed the SA's failure on (3.4), and pointed out that *only* the introduction of $x_4 - x_5$ as a separate variable is essential to his fix. Example 3.2 verifies Pryce's argument and shows that the block ES method finds his reformulation in a systematic way.

To prove Theorem 3.2, we shall use the following two assumptions.

- (a) Without loss of generality, we assume the entries $\widehat{v}_j \neq 0$ are in the first s positions of $\widehat{\mathbf{v}}$, that is, $\widehat{\mathbf{v}} = [\widehat{v}_1, \dots, \widehat{v}_s, 0, \dots, 0]^T$. By (3.9), $J = \sum_{m=1}^{q-1} N_m + 1 : \sum_{m=1}^{q-1} N_m + s$.

- (b) We introduce one more variable $y_l = x_l^{(d_l - \bar{c})}$ for the chosen $l \in J$, and append correspondingly one more equation $0 = g_l = -y_l + x_l^{(d_l - \bar{c})}$.

We show first that the signature matrix $\bar{\Sigma}$ of the resulting DAE can be put in the block structure as shown in Figure 3.1. Then we construct two $(n+s)$ -vectors $\tilde{\mathbf{c}}$ and $\tilde{\mathbf{d}}$ in (3.14), and prove in Lemma 3.2 that $\tilde{d}_j - \tilde{c}_i > \bar{\sigma}_{ij}$ holds below the block diagonal, while $\tilde{d}_j - \tilde{c}_i \geq \bar{\sigma}_{ij}$ holds elsewhere. Lastly, we prove Theorem 3.2.

$$\begin{array}{c}
 \left. \begin{array}{ccc|ccc}
 \Sigma_{1,1} & \cdots & \Sigma_{1,q-1} & \Sigma_{1,q} & -\infty_{N_1 \times s} & \Sigma_{1,q+1} & \cdots & \Sigma_{1,p} \\
 \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \Sigma_{q-1,1} & \cdots & \Sigma_{q-1,q-1} & \Sigma_{q-1,q} & -\infty_{N_{q-1} \times s} & \Sigma_{q-1,q+1} & \cdots & \Sigma_{q-1,p}
 \end{array} \right\} f_i \text{ for } i \in B_{<q} \\
 \hline
 \left. \begin{array}{ccc|ccc}
 \bar{\Sigma}_{q,1} & \cdots & \bar{\Sigma}_{q,q-1} & \bar{\Sigma}_{qq,11} & \bar{\Sigma}_{qq,12} & \bar{\Sigma}_{qq,13} & \bar{\Sigma}_{q,q+1} & \cdots & \bar{\Sigma}_{q,p} \\
 & & & \bar{\Sigma}_{qq,21} & \bar{\Sigma}_{qq,22} & \bar{\Sigma}_{qq,23} & & & \\
 \Sigma_{q+1,1} & \cdots & \Sigma_{q+1,q-1} & \Sigma_{q+1,q} & -\infty_{N_{q+1} \times s} & \Sigma_{q+1,q+1} & \cdots & \Sigma_{q+1,p} \\
 \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \Sigma_{p,1} & \cdots & \Sigma_{p,q-1} & \Sigma_{p,q} & -\infty_{N_p \times s} & \Sigma_{p,q+1} & \cdots & \Sigma_{p,p}
 \end{array} \right\} \begin{array}{l} \bar{f}_i \text{ for } i \in B_q \\ g_r \text{ for } r \in J \end{array} \\
 \hline
 \left. \begin{array}{ccc|ccc}
 \Sigma_{q+1,1} & \cdots & \Sigma_{q+1,q-1} & \Sigma_{q+1,q} & -\infty_{N_{q+1} \times s} & \Sigma_{q+1,q+1} & \cdots & \Sigma_{q+1,p} \\
 \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \Sigma_{p,1} & \cdots & \Sigma_{p,q-1} & \Sigma_{p,q} & -\infty_{N_p \times s} & \Sigma_{p,q+1} & \cdots & \Sigma_{p,p}
 \end{array} \right\} f_i \text{ for } i \in B_{>q} \\
 \hline
 \underbrace{\hspace{1.5cm}}_{x_j \text{ for } j \in B_{<q}} \quad \underbrace{\hspace{1.5cm}}_{x_j \text{ for } j \in B_q} \quad \underbrace{\hspace{1.5cm}}_{y_j \text{ for } j \in J} \quad \underbrace{\hspace{1.5cm}}_{x_j \text{ for } j \in B_{>q}}
 \end{array}$$

Fig. 3.1: Block structure of $\bar{\Sigma}$ of the resulting DAE by the block ES method. The notation $B_{<q}$ is short for $\cup_{m=1}^{q-1} B_m$, and $B_{>q}$ is short for $\cup_{m=q+1}^p B_m$.

From the description of the conversion in Theorem 3.2, the substitutions (3.12) only occur in equations f_i with $i \in B_q$. Hence, in the resulting DAE, variables y_j for $j \in J$ only appear in \bar{f}_i for $i \in B_q$ and g_r for $r \in J$.

Considering the block structure of $\bar{\Sigma}$ in Figure 3.1, we distinguish between the four cases for a block submatrix $\bar{\Sigma}_{m_1 m_2}$: (a) $m_1 \neq q$ and $m_2 \neq q$, (b) $m_1 \neq q$ and $m_2 = q$, (c) $m_1 = q$ and $m_2 \neq q$, and (d) $m_1 = m_2 = q$.

- (a) $m_1 \neq q$ and $m_2 \neq q$. In $\bar{\Sigma}_{m_1 m_2}$, equations f_i are of indices $i \in B_{<q} \cup B_{>q}$. As noted in (3.9), the expression substitutions described in (3.12) only take place in $f_{i'}$ with $i' \in M \subseteq B_q$, so do not happen in such blocks $\bar{\Sigma}_{m_1 m_2}$. Hence, each $\Sigma_{m_1 m_2}$ remains unchanged in $\bar{\Sigma}$: $\bar{\Sigma}_{m_1 m_2} = \Sigma_{m_1 m_2}$ for $m_1 \neq q$ and $m_2 \neq q$.
- (b) $m_1 \neq q$ and $m_2 = q$. In $\bar{\Sigma}_{m_1 q}$, we include variables y_j for $j \in J$ as defined in (3.11). By the same arguments as in (a), the expression substitutions do not happen in these blocks. That is, y_j for $j \in J$ do not appear in equations f_i for $i \in B_{<q} \cup B_{>q}$. Hence, we can obtain $\bar{\Sigma}_{m_1 q}$ by concatenating horizontally $\Sigma_{m_1 q}$ with an $N_{m_1} \times s$ matrix of $-\infty$'s: $\bar{\Sigma}_{m_1 q} = [\Sigma_{m_1 q}, -\infty_{N_{m_1} \times s}]$ for $m_1 = 1, \dots, q-1, q+1, \dots, p$.
- (c) $m_1 = q$ and $m_2 \neq q$. In $\bar{\Sigma}_{qm_2}$, we include equations g_r for $r \in J$ as defined in (3.13). Also, due to the expression substitutions (3.12) occurring in f_i with $i \in M \subseteq B_q$, $\sigma(x_j, f_i)$ and $\sigma(x_j, \bar{f}_i)$ may not be the same for all $i \in B_q$ and all $j = 1:n$.

Hence, in contrast to cases (a) and (b), there are no obvious connections between $\Sigma_{m_1 m_2}$ and $\bar{\Sigma}_{m_1 m_2}$ for $m_1 = q$ and $m_2 \neq q$.

- (d) $m_1 = m_2 = q$. $\bar{\Sigma}_{qq}$ contains signature entries for equations \bar{f}_i and g_r , where $i \in B_q$ and $r \in J$, in variables x_j and y_r , where $j \in B_q$ and $r \in J$. Similar to the resulting signature matrix $\bar{\Sigma}$ by the basic ES method [13, §4.2], $\bar{\Sigma}_{qq}$ by the block ES method also has a (sub)block structure; c.f (A.11). We shall use it in the proof of Lemma 3.2 in Appendix A.

Let $Q = \sum_{m=1}^q N_m$ denote the total number of equations (or variables) in the first q blocks. Using a valid $(\mathbf{c}; \mathbf{d})$ of Σ , we construct an offset pair $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ defined as

$$\tilde{c}_i = \begin{cases} c_i & \text{if } i = 1:Q \\ \bar{c} & \text{if } i = Q+1:Q+s \\ c_{i-s} & \text{if } i = Q+s+1:n+s, \end{cases} \quad \tilde{d}_j = \begin{cases} d_j & \text{if } j = 1:Q \\ \bar{c} & \text{if } j = Q+1:Q+s \\ d_{j-s} & \text{if } j = Q+s+1:n+s. \end{cases} \quad (3.14)$$

Then we have the following lemma. Since its proof is rather technical, we present it in Appendix A.

Lemma 3.2 *Let $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ be as constructed in (3.14). In the block structure of $\bar{\Sigma}$ in Figure 3.1, if a position (i, j) in $\bar{\Sigma}$ is below the block diagonal, then $\tilde{d}_j - \tilde{c}_i > \bar{\sigma}_{ij}$; otherwise, $\tilde{d}_j - \tilde{c}_i \geq \bar{\sigma}_{ij}$.*

Using this lemma, we can now prove Theorem 3.2.

Proof Let \bar{T} be a transversal of $\bar{\Sigma}$. Using Lemma 3.2 and (3.14), we derive

$$\begin{aligned} \text{Val}(\bar{\Sigma}) &= \sum_{(i,j) \in \bar{T}} \bar{\sigma}_{ij} \leq \sum_{(i,j) \in \bar{T}} (\tilde{d}_j - \tilde{c}_i) = \sum_{j=1}^{n+s} \tilde{d}_j - \sum_{i=1}^{n+s} \tilde{c}_i \\ &= \left(\sum_{j=1}^Q d_j + s\bar{c} + \sum_{j=Q+s+1}^{n+s} d_{j-s} \right) - \left(\sum_{j=1}^Q c_j + s\bar{c} + \sum_{i=Q+s+1}^{n+s} c_{i-s} \right) \\ &= \sum_{j=1}^n d_j - \sum_{i=1}^n c_i = \text{Val}(\Sigma). \end{aligned}$$

As in the proof of Theorem 3.1, we prove $\text{Val}(\bar{\Sigma}) < \text{Val}(\Sigma)$ by contradiction. Assume $\text{Val}(\bar{\Sigma}) = \sum_{j=1}^{n+s} \tilde{d}_j - \sum_{i=1}^{n+s} \tilde{c}_i = \text{Val}(\Sigma)$, which is ≥ 0 . With Lemma 3.2, the three conditions in Lemma 2.3 are satisfied. Then it follows from Lemma 2.3 that

- (a) $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ is a valid offset pair of $\bar{\Sigma}$;
- (b) the Jacobian pattern \bar{S}_0 , derived from $\bar{\Sigma}$ and $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$, is in the $p \times p$ BTF shown in Figure 3.1;
- (c) \bar{T} is the union of HVTs \bar{T}_m of all diagonal blocks $\bar{\Sigma}_{11}, \dots, \bar{\Sigma}_{pp}$ of $\bar{\Sigma}$.

We can consider block q of the original DAE as a sub-DAE, with signature matrix Σ_{qq} and offset pair $(\mathbf{c}_q; \mathbf{d}_q)$ —this follows from Lemma 2.4. The conversion described in Theorem 3.2 can be regarded as an application of the basic ES method to this sub-DAE, given that the ES conditions [13, (4.11)] hold because of (3.10). By [13, Theorem 4.2] for the basic ES method, a conversion results in $\text{Val}(\bar{\Sigma}_{qq}) < \text{Val}(\Sigma_{qq})$. Also, since $\bar{\Sigma}_{mm} = \Sigma_{mm}$ for $m \neq q$, $\text{Val}(\bar{\Sigma}_{mm}) = \text{Val}(\Sigma_{mm})$. Then a contradiction $\text{Val}(\bar{\Sigma}) < \text{Val}(\Sigma)$ follows by the same derivations as in (3.8). The assumption $\text{Val}(\bar{\Sigma}) = \text{Val}(\Sigma)$ is hence false, so $\text{Val}(\bar{\Sigma}) < \text{Val}(\Sigma)$ holds. \square

$$\begin{array}{c}
\begin{array}{c} x_7 \ x_8 \ y_8 \ x_4 \ x_5 \ y_5 \ x_6 \ x_{11} \ x_{12} \ y_2 \ x_3 \ c_i \\
\begin{array}{c} h_8 \\ \bar{f}_8 \\ \bar{f}_7 \\ h_5 \\ \bar{f}_5 \\ \bar{f}_4 \\ f_6 \\ h_2 \\ \bar{f}_2 \\ \bar{f}_1 \\ f_3 \\ d_j \end{array} \\
\begin{array}{c} 1 \bullet \ 1 \ 0 \\ 0 \bullet \ 0 \bullet \\ 0 \bullet \ 0 \bullet \\ 1 \bullet \ 1 \ 0 \\ 0 \bullet \ 0 \bullet \\ 0 \bullet \ 0 \bullet \\ 1 \bullet \ 1 \ 0 \\ 0 \bullet \ 0 \bullet \\ 1 \bullet \ 1 \ 0 \\ 0 \bullet \ 0 \bullet \\ 0 \bullet \ 0 \bullet \\ 1 \bullet \ 1 \ 0 \\ 0 \bullet \ 0 \bullet \\ 0 \bullet \ 0 \bullet \\ 1 \bullet \ 1 \ 0 \\ 0 \bullet \ 0 \bullet \\ 1 \bullet \ 1 \ 0 \\ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \end{array} \\
\end{array}
\end{array}
\quad
\begin{array}{c}
\begin{array}{c} x'_7 \ x'_8 \ y'_8 \ x'_4 \ x'_5 \ y'_5 \ x'_6 \ x'_{11} \ x'_{12} \ y'_2 \ x'_3 \\
\begin{array}{c} h_8 \\ \bar{f}_8 \\ \bar{f}_7 \\ h_5 \\ \bar{f}_5 \\ \bar{f}_4 \\ f_6 \\ h_2 \\ \bar{f}_2 \\ \bar{f}_1 \\ f_3 \end{array} \\
\begin{array}{c} -1 \ 1 \\ R_9^{-1} \ C_5 \\ R_8^{-1} \ -C_5 \\ -1 \ 1 \\ \frac{\partial \bar{f}_7}{\partial x_5} \ \frac{\partial \bar{f}_7}{\partial x_6} \\ R_4^{-1} \ -C_3 \\ C_4 \\ -1 \ 1 \\ \frac{\partial \bar{f}_2}{\partial x_2} \ C_1 \\ R_0^{-1} \ -C_1 \\ C_3 \end{array} \\
\end{array}
\end{array}$$

Here in $\bar{\mathbf{J}}$, $\partial \bar{f}_5 / \partial x_5 = R_5^{-1} + R_6^{-1}$ and $\partial \bar{f}_2 / \partial x_2 = R_1^{-1} + R_2^{-1}$. The SA succeeds with a nonzero constant $\det(\bar{\mathbf{J}})$ and $\text{Val}(\bar{\Sigma}) = 5 < 8 = \text{Val}(\Sigma)$.

4.2 Ring modulator.

We study the ring modulator problem from [3]. When $C_s \neq 0$, it is a stiff ODE system of 15 nonlinear equations. Setting $C_s = 0$ gives a DAE of differentiation index 2 that consists of 11 differential and 4 algebraic equations:

$$\begin{aligned}
0 &= f_1 = -x'_1 + C^{-1}(x_8 - 0.5x_{10} + 0.5x_{11} + x_{14} - R^{-1}x_1) \\
0 &= f_2 = -x'_2 + C^{-1}(x_9 - 0.5x_{11} + 0.5x_{12} + x_{15} - R^{-1}x_2) \\
0 &= f_3 = x_{10} - q(U_{D1}) + q(U_{D4}) \\
0 &= f_4 = -x_{11} + q(U_{D2}) - q(U_{D3}) \\
0 &= f_5 = x_{12} + q(U_{D1}) - q(U_{D3}) \\
0 &= f_6 = -x_{13} - q(U_{D2}) + q(U_{D4}) \\
0 &= f_7 = -x'_7 + C_p^{-1}(-R_p^{-1}x_7 + q(U_{D1}) + q(U_{D2}) - q(U_{D3}) - q(U_{D4})) \\
0 &= f_8 = -x'_8 + -L_h^{-1}x_1 \\
0 &= f_9 = -x'_9 + -L_h^{-1}x_2 \\
0 &= f_{10} = -x'_{10} + L_{s2}^{-1}(0.5x_1 - x_3 - R_{g2}x_{10}) \\
0 &= f_{11} = -x'_{11} + L_{s3}^{-1}(-0.5x_1 + x_4 - R_{g3}x_{11}) \\
0 &= f_{12} = -x'_{12} + L_{s2}^{-1}(0.5x_2 - x_5 - R_{g2}x_{12}) \\
0 &= f_{13} = -x'_{13} + L_{s3}^{-1}(-0.5x_2 + x_6 - R_{g3}x_{13}) \\
0 &= f_{14} = -x'_{14} + L_{s1}^{-1}(-x_1 + U_{in1}(t) - (R_i + R_{g1})x_{14}) \\
0 &= f_{15} = -x'_{15} + L_{s1}^{-1}(-x_2 - (R_c + R_{g1})x_{15}) .
\end{aligned}$$

The functions are

$$\begin{aligned}
 U_{D1} &= x_3 - x_5 - x_7 - U_{in2}(t) & q(U) &= \gamma(e^{\delta U} - 1) \\
 U_{D2} &= -x_4 + x_6 - x_7 - U_{in2}(t) & U_{in1}(t) &= 0.5 \sin(2000\pi t) \\
 U_{D3} &= x_4 + x_5 + x_7 + U_{in2}(t) & U_{in2}(t) &= 2 \sin(20000\pi t) \\
 U_{D4} &= -x_3 - x_6 + x_7 + U_{in2}(t) .
 \end{aligned}$$

We refer the reader to [3] for the nonzero constants $C, C_p, R, R_p, R_c, \gamma, L_h, L_{s1}, L_{s2}, L_{s3}, R_{g1}, R_{g2}, R_{g3}, R_i$, and δ .

$$\Sigma = \begin{array}{c}
 \begin{array}{cccccccccccccccc}
 x_1 & x_2 & x_7 & x_{13} & x_{11} & x_{12} & x_{10} & x_3 & x_4 & x_5 & x_6 & x_8 & x_9 & x_{14} & x_{15} & c_i \\
 f_1 & 1^\bullet & & & 0 & & 0 & & & & & 0 & & 0 & & 0 \\
 f_2 & & 1^\bullet & & 0 & & 0 & & & & & & 0 & & 0 & 0 \\
 f_7 & & & 1^\bullet & & & & 0 & 0 & 0 & 0 & & & & & 0 \\
 f_{13} & & 0 & & 1^\bullet & & & & & & & 0 & & & & 0 \\
 f_{11} & 0 & & & & 1^\bullet & & & 0 & & & & & & & 0 \\
 f_{12} & & 0 & & & & 1^\bullet & & & 0 & & & & & & 0 \\
 f_{10} & 0 & & & & & & 1^\bullet & 0 & & & & & & & 0 \\
 f_3 & & & 0 & & & & 0^\bullet & 0 & 0 & & & & & & 0 \\
 f_4 & & & 0 & & 0 & & & 0^\bullet & 0 & 0 & & & & & 0 \\
 f_5 & & & 0 & & & 0 & & & 0^\bullet & & & & & & 0 \\
 f_6 & & & 0 & 0 & & & & & & 0^\bullet & & & & & 0 \\
 f_8 & 0 & & & & & & & & & & 1^\bullet & & & & 0 \\
 f_9 & & 0 & & & & & & & & & & 1^\bullet & & & 0 \\
 f_{14} & 0 & & & & & & & & & & & & 1^\bullet & & 0 \\
 f_{15} & & 0 & & & & & & & & & & & & 1^\bullet & 0
 \end{array} \\
 d_j & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0
 \end{array}$$

Each block of size 1 has a nonsingular sub-Jacobian -1 . Block 8 has an identically singular sub-Jacobian

$$\mathbf{J}_{88} = \begin{array}{c}
 \begin{array}{cccc}
 x_3 & x_4 & x_5 & x_6 \\
 f_3 & -s_1 - s_4 & & s_1 & -s_4 \\
 f_4 & & -s_2 - s_3 & -s_3 & s_2 \\
 f_5 & s_1 & -s_3 & -s_1 - s_3 & \\
 f_6 & -s_4 & s_2 & & -s_2 - s_4
 \end{array}
 \end{array}, \quad \text{where } s_i = \gamma \delta e^{\delta U_{Di}} .$$

This is a nonlinear block, since variables x_3, x_4, x_5, x_6 do not occur jointly linearly in equations f_3, f_4, f_5, f_6 . One can also see these variables appear in \mathbf{J}_{88} .

LC method. We find a constant vector $\hat{\mathbf{u}} = [1, -1, 1, -1]^T \in \text{coker}(\mathbf{J}_{88})$, which satisfies the block LC condition (3.3). Then $\mathbf{u} = [\mathbf{0}_7^T, 1, -1, 1, -1, \mathbf{0}_4^T]^T$. We use (3.2) to derive

$$I = \{i \mid u_i \neq 0\} = \{8, 9, 10, 11\}, \quad \underline{c} = 0, \quad \text{and} \quad L = \bar{L} = \{8, 9, 10, 11\}.$$

The row indices in \bar{L} correspond to the equations f_3, f_4, f_5, f_6 . We can pick any one of them and replace it by

$$\bar{f} = u_1 f_3 + u_2 f_4 + u_3 f_5 + u_4 f_6 = f_3 - f_4 + f_5 - f_6 = x_{10} + x_{11} + x_{12} + x_{13}.$$

We choose f_3 and replace it by $\bar{f}_3 = \bar{f}$. The resulting DAE has the following $\bar{\Sigma}$ with $\text{Val}(\bar{\Sigma}) = 10 < 11 = \text{Val}(\Sigma)$.

$$\bar{\Sigma} = \begin{array}{cccccccccccccccc} & x_1 & x_2 & x_7 & x_3 & x_4 & x_5 & x_6 & x_{10} & x_{11} & x_{12} & x_{13} & x_8 & x_9 & x_{14} & x_{15} & c_i \\ f_1 & 1^\bullet & & & & & & & 0 & 0 & & & 0 & & 0 & & 0 \\ f_2 & & 1^\bullet & & & & & & & & 0 & 0 & & 0 & & 0 & 0 \\ f_7 & & & 1^\bullet & 0 & 0 & 0 & 0 & & & & & & & 0 & & 0 \\ f_{10} & 0 & & & 0^\bullet & & & & 1 & & & & & & & & 0 \\ f_5 & & & 0 & 0 & 0^\bullet & 0 & & & & 0 & & & & & & 0 \\ f_4 & & & 0 & 0 & 0^\bullet & 0 & & & 0 & & & & & & & 0 \\ f_6 & & & 0 & 0 & 0 & 0^\bullet & & & & & 0 & & & & & 0 \\ \bar{f}_3 & & & & & & & & 0^\bullet & 0 & 0 & 0 & & & & & 1 \\ f_{11} & 0 & & & & & & & & 1^\bullet & & & & & & & 0 \\ f_{12} & & 0 & & & & & & & & 1^\bullet & & & & & & 0 \\ f_{13} & & & & & & & 0 & & & & 1^\bullet & & & & & 0 \\ f_8 & 0 & & & & & & & & & & & 1^\bullet & & & & 0 \\ f_9 & & 0 & & & & & & & & & & & 1^\bullet & & & 0 \\ f_{14} & 0 & & & & & & & & & & & & & 1^\bullet & & 0 \\ f_{15} & & 0 & & & & & & & & & & & & & 1^\bullet & 0 \\ d_j & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{array}$$

Again, each 1×1 block has a nonsingular Jacobian: $\partial f_i / \partial x'_i = -1$ for $i = 1, 2, 7, 8, 9, 14, 15$. The sub-Jacobian of block 4 in the resulting DAE is

$$\bar{J}_{44} = \begin{array}{cccccccc} & x_3 & x_4 & x_5 & x_6 & x'_{10} & x'_{11} & x'_{12} & x'_{13} \\ f_{10} & -L_{s_2}^{-1} & & & & & & & -1 \\ f_5 & s_1 & -s_3 & -s_1 - s_3 & & & & & \\ f_4 & & -s_2 - s_3 & -s_3 & s_2 & & & & \\ f_6 & -s_4 & s_2 & & -s_2 - s_4 & & & & \\ \bar{f}_3 & & & & & 1 & 1 & 1 & 1 \\ f_{11} & & L_{s_3}^{-1} & & & & -1 & & \\ f_{12} & & & -L_{s_2}^{-1} & & & & -1 & \\ f_{13} & & & & L_{s_3}^{-1} & & & & -1 \end{array},$$

whose determinant is $\det(\bar{J}_{44}) = 2s_1 s_2 s_3 s_4 (s_1^{-1} + s_2^{-1} + s_3^{-1} + s_4^{-1})(L_{s_2}^{-1} + L_{s_3}^{-1})$. The SA succeeds at any point where $\det(\bar{J}_{44}) \neq 0$, and the DAE is of index 2.

ES method. Find $\hat{\mathbf{v}} = [-1, 1, -1, 1]^T \in \ker(\mathbf{J}_{88})$. Then $\mathbf{v} = [\mathbf{0}_7^T, -1, 1, -1, 1, \mathbf{0}_4^T]^T$. We use (3.9) to derive

$$J = \bar{J} = \{j \mid v_j \neq 0\} = \{8, 9, 10, 11\}, \quad s = |J| = 4, \quad M = J, \quad \text{and} \quad \bar{c} = \mathbf{0}.$$

We choose column index $l = 8 \in \bar{J}$ in the permuted Σ . The variable of this column is x_3 . The other variables in block 8 are x_4, x_5, x_6 , so we introduce for them, respectively,

$$y_4 = x_4 - (v_9/v_8) \cdot x_3, \quad y_5 = x_5 - (v_{10}/v_8) \cdot x_3, \quad \text{and} \quad y_6 = x_6 - (v_{11}/v_8) \cdot x_3.$$

Then we append the equations corresponding to these variables

$$0 = g_4 = -y_4 + x_4 + x_3, \quad 0 = g_5 = -y_5 + x_5 - x_3, \quad \text{and} \quad 0 = g_6 = -y_6 + x_6 + x_3.$$

The equations in block 8 are f_3, f_4, f_5, f_6 . In these equations, we perform the following substitutions.

replace	by	in
x_4	$y_4 - x_3$	f_4, f_5, f_6
x_5	$y_5 + x_3$	f_3, f_4, f_5
x_6	$y_6 - x_3$	f_3, f_4, f_6

The resulting index-2 DAE is of size 18. (We do not display the results of SA here.) It has $\text{Val}(\bar{\Sigma}) = 10 < 11 = \text{Val}(\Sigma)$ and $\det(\bar{J}) = -2s_1s_2s_3s_4(s_1^{-1} + s_2^{-1} + s_3^{-1} + s_4^{-1})(L_{s_2}^{-1} + L_{s_3}^{-1})$. The largest fine block is of size 12, and the other six fine blocks are of size 1. The SA succeeds at any point where $\det(\bar{J}) \neq 0$.

5 Conclusions.

We combined block triangularization with the LC and ES conversion methods for improving the Σ -method. When J is identically singular and the DAE has a nontrivial BTF, we can locate each diagonal block whose corresponding sub-Jacobian is identically singular, and perform a conversion on it. We base this strategy on the view that each diagonal block can be regarded as a sub-DAE, while formulas contributing to the off diagonal blocks are regarded as driving terms.

Compared with the basic conversion methods that work on the whole DAE, the block methods only work on singular blocks, which are usually smaller than the DAE itself. Hence the block methods require fewer symbolic computations, and can generally find a useful conversion for reducing $\text{Val}(\Sigma)$ more efficiently. As in the basic case, a conversion applied on a singular block guarantees (a) a strict decrease in the value of the (whole) signature matrix, and (b) the equivalence between the original DAE and the resulting one. The rationale for choosing a desirable conversion method is in [13, Table 4.1].

We combine MATLAB's Symbolic Math Toolbox [14] with our structural analysis software DAESA [6, 10], and have built a prototype code that automates the conversion process. We aim to incorporate them in a future version of DAESA.

With our prototype code, we have applied our methods on numerous DAEs on which the Σ -method fails. They are either arbitrarily constructed to be SA-failure cases for our investigations, or borrowed from the existing literature. Our conversion methods succeed in fixing all these solvable DAEs. We believe that our assumptions and conditions are reasonable for practical problems, and that these methods can help make the Σ -method more reliable.

We end these two articles with our main conjecture regarding SA's failure. In all our experiments, when we successfully fix the failure using our conversion methods, the value of a signature matrix always decreases. As Pryce points out in [8], the solvability of a DAE lies within its inherent nature, not the way it is formulated or analyzed. Hence we conjecture that a DAE formulation friendly to SA should have a reasonable but never overestimated $\text{Val}(\mathbf{\Sigma})$ that can be interpreted as number of degrees of freedom (DOF) of the underlying problem. In other words, a DAE should not be formulated to exhibit more DOF than the underlying problem has. However, based on our current knowledge, it appears difficult to show why overestimating DOF can lead to an identically singular System Jacobian.

A Proof of Lemma 3.2.

For $\bar{\mathbf{\Sigma}} = (\bar{\sigma}_{ij})$ in the block structure in Figure 3.1, we write the block sizes in the array

$$\bar{N} = (N_1, N_2, \dots, N_{q-1}, N_q + s, N_{q+1}, \dots, N_p),$$

and also write the block sizes of $\mathbf{\Sigma}$ in the array $N = (N_1, N_2, \dots, N_{q-1}, N_q, N_{q+1}, \dots, N_p)$. Let $\overline{\text{blockOf}}(i)$ denote the block number of a row or column index i in $\bar{\mathbf{\Sigma}}$. From the construction of \bar{N} and N , it is not difficult to show that

$$\overline{\text{blockOf}}(j) < q \Leftrightarrow 1 \leq j \leq \sum_{w=1}^{q-1} N_w \Leftrightarrow \text{blockOf}(j) < q \quad \text{and} \quad (\text{A.1})$$

$$\overline{\text{blockOf}}(j+s) > q \Leftrightarrow \sum_{w=1}^q N_w + s + 1 \leq j+s \Leftrightarrow \text{blockOf}(j) > q. \quad (\text{A.2})$$

From the construction of $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ in (3.14), each variable x_j for $j = 1:n$ has the same "variable offset" in $\bar{\mathbf{\Sigma}}$ as x_j has in $\mathbf{\Sigma}$. Also, each equation \tilde{f}_i for $i = 1:n$ has the same "equation offset" in $\bar{\mathbf{\Sigma}}$ as f_i has in $\mathbf{\Sigma}$. Quotation marks are used here because $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ is *not* a valid offset pair of $\bar{\mathbf{\Sigma}}$; this vector pair is merely used for proving $\text{Val}(\bar{\mathbf{\Sigma}}) < \text{Val}(\mathbf{\Sigma})$ in Theorem 3.2.

We aim to show that

$$\tilde{d}_j - \tilde{c}_i \begin{cases} > \bar{\sigma}_{ij} & \text{if } \overline{\text{blockOf}}(j) < \overline{\text{blockOf}}(i) \\ \geq \bar{\sigma}_{ij} & \text{if } \overline{\text{blockOf}}(j) \geq \overline{\text{blockOf}}(i). \end{cases} \quad (\text{A.3})$$

For the block structure of $\bar{\mathbf{\Sigma}}$ in Figure 3.1, we have shown on page 15 that $\bar{\Sigma}_{m_1 m_2} = \Sigma_{m_1 m_2}$ if $m_1 \neq q$ and $m_2 \neq q$, and that $\bar{\Sigma}_{m_1 m_2} = [\Sigma_{m_1 q}, -\infty_{N_{m_1} \times s}]$ if $m_1 \neq q$ and $m_2 = q$. Hence, provided $m_1 \neq q$, $\bar{\Sigma}_{m_1 m_2}$ is below [resp. above] the block diagonal of $\bar{\mathbf{\Sigma}}$, if $\Sigma_{m_1 m_2}$ is below [resp. above] the block diagonal of $\mathbf{\Sigma}$. By (2.4), the inequalities in (A.3) hold for i with $\overline{\text{blockOf}}(i) \neq q$.

What remains to show is the inequalities in (A.3) for i with $\overline{\text{blockOf}}(i) = q$. These inequalities are for the signature entries in $\bar{\Sigma}_{qm_2}$, the blocks affected by the expression substitutions. We consider three cases for $\bar{\Sigma}_{qm_2}$: it is (a) below the block diagonal, with $m_2 < q$, (b) above the block diagonal, with $m_2 > q$, or (c) the diagonal block $\bar{\Sigma}_{qq}$, with $m_2 = q$.

(a) $\bar{\Sigma}_{qm_2}$ with $m_2 < q$. An entry (i, j) in this block satisfies $\overline{\text{blockOf}}(j) < \overline{\text{blockOf}}(i) = q$. By (A.1), $\overline{\text{blockOf}}(j) < q$ and hence $j \notin B_q$.

Recall from (3.12) that, in each f_i with $i \in M \subseteq B_q$, we substitute $(y_r + \frac{v_r}{v_l} \cdot x_l^{(d_l - \bar{c})})^{(\bar{c} - c_i)}$ for each $x_r^{(\sigma_{ir})}$ with $d_r - c_i = \sigma_{ir}$ and $r \in J \setminus \{l\} \subset B_q$. For a $j \notin B_q \supset J \setminus \{l\}$, the corresponding derivatives $x_j^{(d_j - c_i)}$ are not replaced in the ES conversion, and for $r \in J \setminus \{l\}$ (so j, r, l are distinct),

$$\sigma \left(x_j, \left(y_r + \frac{v_r}{v_l} \cdot x_l^{(d_l - \bar{c})} \right)^{(\bar{c} - c_i)} \right) = \sigma \left(x_j, \left(\frac{v_r}{v_l} \right)^{(\bar{c} - c_i)} \right) \leq \sigma \left(x_j, \mathbf{v}^{(\bar{c} - c_i)} \right) = \sigma(x_j, \mathbf{v}) + (\bar{c} - c_i). \quad (\text{A.4})$$

By (3.10), $\sigma(x_j, \mathbf{v}) < d_j - \bar{c}$. Using (2.4) and (A.4), we derive

$$\begin{aligned} \sigma(x_j, \bar{f}_i) &\leq \max \left\{ \sigma(x_j, f_i), \max_{r \in J \setminus \{l\}} \sigma \left(x_j, \left(y_r + \frac{y_r}{v_l} \cdot x_l^{(d_l - \bar{c})} \right)^{(\bar{c} - c_i)} \right) \right\} \\ &\leq \max \{ \bar{\sigma}_{ij}, \sigma(x_j, \mathbf{v}) + (\bar{c} - c_i) \} \\ &< \max \{ d_j - c_i, (d_j - \bar{c}) + (\bar{c} - c_i) \} = d_j - c_i \quad \text{for } i \in M \subseteq B_q. \end{aligned} \quad (\text{A.5})$$

From the ES conversion described in Theorem 3.2, we have

$$\sigma(x_j, \bar{f}_i) = \sigma(x_j, f_i) < d_j - c_i \quad \text{for } i \in B_q \setminus M \text{ and} \quad (\text{A.6})$$

$$\sigma(x_j, g_r) \leq \sigma(x_j, \mathbf{v}) < d_j - \bar{c} \quad \text{for } r \in J. \quad (\text{A.7})$$

Since blocks $\bar{\Sigma}_{qm_2}$ with $m_2 < q$ contain signature entries $\bar{\sigma}_{ij}$ for equations \bar{f}_i and g_r , where $i \in B_q$ and $r \in J$, in variables x_j with $\text{blockOf}(j) < q$, by taking together the inequalities in (A.5)-(A.7), we have

$$\bar{\sigma}_{ij} < \begin{cases} d_j - c_i & \text{if } \text{blockOf}(j) < q \text{ and } i \in B_q \\ d_j - \bar{c} & \text{if } \text{blockOf}(j) < q \text{ and } i \in Q+1:Q+s; \end{cases}$$

recall $Q = \sum_{w=1}^q N_w$. Using (A.1) and the construction of $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ in (3.14), we have

$$\bar{\sigma}_{ij} < \tilde{d}_j - \tilde{c}_i \quad \text{for } \overline{\text{blockOf}}(j) < \overline{\text{blockOf}}(i) = q.$$

(q is the block number of both the original and enlarged diagonal blocks.)

(b) $\bar{\Sigma}_{qm_2}$ with $m_2 > q$. An entry $(i, j+s)$ in this block satisfies $\overline{\text{blockOf}}(j+s) > \overline{\text{blockOf}}(i) = q$. By (A.2), $\text{blockOf}(j) > q$ and hence $j \notin B_q \supset J \setminus \{l\}$. By the same arguments as in (a), the corresponding derivatives $x_j^{(d_j - c_i)}$ are not replaced in the ES conversion.

By (3.10), $\sigma(x_j, \mathbf{v}) \leq d_j - \bar{c}$. Then by the same derivations as (A.5-A.7) in (a), we have

$$\sigma(x_j, \bar{f}_i) \leq d_j - c_i \quad \text{for } i \in M \subseteq B_q, \quad (\text{A.8})$$

$$\sigma(x_j, \bar{f}_i) = \sigma(x_j, f_i) \leq d_j - c_i \quad \text{for } i \in B_q \setminus M, \text{ and} \quad (\text{A.9})$$

$$\sigma(x_j, g_r) \leq \sigma(x_j, \mathbf{v}) \leq d_j - \bar{c} \quad \text{for } r \in J. \quad (\text{A.10})$$

Since blocks $\bar{\Sigma}_{qm_2}$ with $m_2 > q$ contain signature entries $\bar{\sigma}_{i,j+s}$ for equations \bar{f}_i and g_r , where $i \in B_q$ and $r \in J$, in variables x_j with $\text{blockOf}(j) > q$, the inequalities (A.8-A.10) yield

$$\bar{\sigma}_{i,j+s} \leq \begin{cases} d_j - c_i & \text{if } \text{blockOf}(j) > q \text{ and } i \in B_q \\ d_j - \bar{c} & \text{if } \text{blockOf}(j) > q \text{ and } i \in Q+1:Q+s. \end{cases}$$

Using again (A.2) and $(\tilde{\mathbf{c}}; \tilde{\mathbf{d}})$ in (3.14), we have $\bar{\sigma}_{i,j+s} \leq \tilde{d}_{j+s} - \tilde{c}_i$ for $\overline{\text{blockOf}}(j+s) > \overline{\text{blockOf}}(i) = q$, with $j = Q+1:n$. We can rewrite this inequality as

$$\bar{\sigma}_{ij} \leq \tilde{d}_j - \tilde{c}_i \quad \text{for } \overline{\text{blockOf}}(j) > \overline{\text{blockOf}}(i) = q,$$

with $j = Q+1+s:n+s$.

(c) $\bar{\Sigma}_{qm_2}$ is $\bar{\Sigma}_{qq}$, with $m_2 = q$. An entry (i, j) in $\bar{\Sigma}_{qq}$ satisfies $\overline{\text{blockOf}}(j) = \overline{\text{blockOf}}(i) = q$. We view block q of the original DAE as a sub-DAE, with a signature matrix Σ_{qq} of size N_q and an offset pair $(\mathbf{c}_q; \mathbf{d}_q)$. Given that the ES conditions are satisfied by (3.10), performing the ES conversion as described in Theorem 3.2 is equivalent to applying the basic ES method to this sub-DAE. After a conversion, the resulting enlarged signature matrix $\bar{\Sigma}_{qq}$ of size $N_q + s$ has the form

$$\bar{\Sigma}_{qq} = \begin{bmatrix} \bar{\Sigma}_{qq,11} & \bar{\Sigma}_{qq,12} & \bar{\Sigma}_{qq,13} \\ \bar{\Sigma}_{qq,21} & \bar{\Sigma}_{qq,22} & \bar{\Sigma}_{qq,23} \end{bmatrix}; \quad (\text{A.11})$$

cf. Figure 3.1 and the details of the basic ES method in [13]. The two block rows of $\bar{\Sigma}_{qq}$ correspond to \bar{f}_i for $i \in B_q$ and g_j for $j \in J$, respectively. The three block columns of $\bar{\Sigma}_{qq}$ correspond to x_j for $j \in J$, x_j for $j \in B_q \setminus J$, and y_j for $j \in J$, respectively. If we apply the same arguments in the proof of [13, Lemma 4.4] for the basic ES method, then we have $\tilde{d}_j - \tilde{c}_i \geq \bar{\sigma}_{ij}$ for all entries in $\bar{\Sigma}_{qq}$. \square

Acknowledgements The authors acknowledge with thanks the financial support for this research: GT is supported in part by the McMaster Centre for Software Certification through the Ontario Research Fund, Canada, NSN is supported in part by the Natural Sciences and Engineering Research Council of Canada, and JDP is supported in part by the Leverhulme Trust, the UK.

References

1. Campbell, S.L., Griepentrog, E.: Solvability of general differential-algebraic equations. *SIAM Journal on Scientific Computing* **16**(2), 257–270 (1995)
2. Duff, I., Erisman, A., Reid, J.: *Direct Methods for Sparse Matrices*. Oxford Science Publications. Clarendon Press, Oxford (1986)
3. Mazzia, F., Iavernaro, F.: Test set for initial value problem solvers. Tech. Rep. 40, Department of Mathematics, University of Bari, Italy (2003). <http://pitagora.dm.uniba.it/~testset/>
4. McKenzie, R.: Structural analysis based dummy derivative selection for differential-algebraic equations. Tech. rep., Cardiff University, UK (2015). Submitted to BIT Numerical Mathematics, Oct 2015.
5. McKenzie, R.: Reducing the index of differential-algebraic equations by exploiting underlying structures. PhD Thesis, School of Mathematics, Cardiff University, Senghennydd Road, Cardiff CF24 4AG, UK (2016)
6. Nedialkov, N.S., Pryce, J.D., Tan, G.: Algorithm 948: DAESA—a Matlab tool for structural analysis of differential-algebraic equations: Software. *ACM Trans. Math. Softw.* **41**(2), 12:1–12:14 (2015)
7. Nedialkov, N.S., Tan, G., Pryce, J.D.: Exploiting fine block triangularization and quasilinearity in differential-algebraic equation systems (2016). McMaster University, Cardiff University. In preparation
8. Pryce, J.D.: Solving high-index DAEs by Taylor Series. *Numerical Algorithms* **19**, 195–211 (1998)
9. Pryce, J.D.: A simple structural analysis method for DAEs. *BIT Numerical Mathematics* **41**(2), 364–394 (2001)
10. Pryce, J.D., Nedialkov, N.S., Tan, G.: DAESA—a Matlab tool for structural analysis of differential-algebraic equations: Theory. *ACM Trans. Math. Softw.* **41**(2), 9:1–9:20 (2015)
11. Pryce, J.D., Nedialkov, N.S., Tan, G.: Fine block triangular structure of DAEs and its uses (2016). Cardiff University, McMaster University. In preparation
12. Tan, G., Nedialkov, N., Pryce, J.: Symbolic-numeric methods for improving structural analysis of differential-algebraic equation systems. Tech. rep., Department of Computing and Software, McMaster University, 1280 Main St. W., Hamilton, ON, L8S4L8, Canada (2015). CAS-15-07-NN, 84 pages
13. Tan, G., Nedialkov, N.S., Pryce, J.D.: Conversion methods for improving structural analysis of differential-algebraic equation systems. *BIT Numerical Mathematics* (2016). Submitted
14. The MathWorks, Inc.: Matlab Symbolic Math Toolbox (2016). <http://www.mathworks.com/products/symbolic/>