

# Conditioning of two deck differential algebraic equations

***Citation for published version (APA):***

Mattheij, R. M. M., & Wijckmans, P. M. E. J. (1997). *Conditioning of two deck differential algebraic equations*. (RANA : reports on applied and numerical analysis; Vol. 9719). Technische Universiteit Eindhoven.

***Document status and date:***

Published: 01/01/1997

***Document Version:***

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

***Please check the document version of this publication:***

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

***General rights***

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

***Take down policy***

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

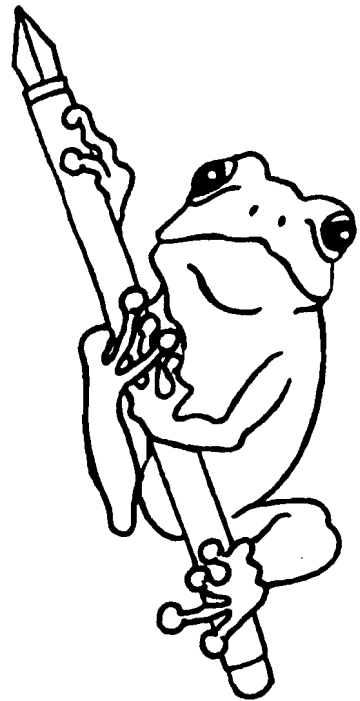
EINDHOVEN UNIVERSITY OF TECHNOLOGY  
Department of Mathematics and Computing Science

RANA 97-19  
November 1997

Conditioning of Two Deck  
Differential Algebraic Equations

by

R.M.M. Mattheij and P.M.E.J. Wijckmans



Reports on Applied and Numerical Analysis  
Department of Mathematics and Computing Science  
Eindhoven University of Technology  
P.O. Box 513  
5600 MB Eindhoven  
The Netherlands  
ISSN: 0926-4507

# Conditioning of Two Deck Differential Algebraic Equations

R.M.M. Mattheij<sup>†</sup>, P.M.E.J. Wijkmans<sup>‡</sup>

<sup>†</sup>*Department of Mathematics and Computing Science,  
Eindhoven University of Technology,  
P.O. Box 513, 5600 MB Eindhoven, The Netherlands*

<sup>‡</sup>*TNO Physics and Electronics Laboratory,  
P.O. Box 96864, 2509 JG 's-Gravenhage, The Netherlands*

## Abstract

For semi-explicit DAE consisting of an ODE system coupled with an algebraic equation an analysis is given of the sensitivity of the problem with respect to additive perturbations. In particular for index-1 and index-2 DAE conditioning constants are derived which show how an ill-conditioned problem may be "close" to some higher index problem. This knowledge might be employed to regularise problems. We also classify problems which are "nearly singular", thus leading to ill-conditioning beyond repair.

# 1 Introduction

Both from a theoretical and a practical point of view, it is interesting to know how sensitive a mathematical representation of a given problem is with respect to (small) perturbations. In purely analytical terms this question is usually referred to as well- (or ill-)posedness. In the numerical literature one often prefers the terminology well- (or ill-)conditioning. For ordinary differential equations (ODE) these concepts have been worked out extensively, both for initial value problems and for boundary value problems; one may consult general references like [2]. For ODE subject to constraints, i.e. differential algebraic equations (DAE), this is only partially true, see e.g. [12]. The problem at hand is also more difficult to describe. As is well known (cf. [6, 8, 9]) there are various notions of index, which describes the degree of complexity of a DAE. For our purposes the differential index (cf. [9]) is quite appropriate. Here, a solution is described in terms of dependence on its initial values and derivatives of a forcing function, up to the "index order". This is a straightforward generalisation of the standard concept of continuous dependence on parameters notion met in ODE.

The purpose of this paper is to find realistic estimates for the stability constants as are met in estimates for such solutions. As will become clear, the discrete (i.e. stepwise) notion of index is an analytical rather than a numerical tool. Indeed one wonders what may happen numerically for DAE with an "index somewhere in between", in particular whether a notion like "nearly index  $\nu$ " would help to clarify the situation (just like we often use the terminology "nearly singular matrix"). In order not to be overambitious we have settled for investigating linear two deck systems, i.e. DAE which typically read

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{Ax} + \mathbf{By} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{Cx} + \mathbf{Dy} + \mathbf{q},\end{aligned}\tag{1.1}$$

where all coefficients may depend on  $t$ . Here we let  $\mathbf{x}(t) \in \mathbb{R}^n$ ,  $\mathbf{y}(t) \in \mathbb{R}^m$ .

As is well known this DAE has index 1 if  $\mathbf{D}$  is nonsingular. On the other hand, if  $\mathbf{D} = \mathbf{O}$ , but  $\mathbf{CB}$  is nonsingular it has index 2. If  $\mathbf{D}$  is singular, some of the variables constituting  $\mathbf{y}$  may be considered as index-1 variables whereas others are (at least) index-2 variables. One may decompose further to find even higher order variables if  $\mathbf{CB}$  is singular.

In the next section we will give a straightforward stability estimation for the simplest index-1 case and likewise in Section 3 for the index-2 case. In Section 4 we first analyse DAE with scalars instead of matrices to illustrate what actually happens when an index one DAE is "close to" an index-2 DAE. This is then generalised to the matrix case in Section 5. In Section 6 we analyse the index-2 matrix case

which is close to a higher index problem. In Section 7 we shall illustrate a concept of effective index and we show in Section 8 why a regularisation method, even for a "formal index-1" problem can be beneficial.

Finally we derive a result for DAE of constant coefficients which have a singular pencil. This corresponds to a case where the effective index is infinite.

## 2 Stability constants for DAE of index-1

Consider the DAE

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{y} + \mathbf{q}.\end{aligned}\tag{2.1}$$

If  $\mathbf{D}$  is nonsingular for all  $t$ , we can express  $\mathbf{y}$  in terms of  $\mathbf{x}$  and  $\mathbf{q}$  through

$$\mathbf{x} = \mathbf{u}, \quad \mathbf{y} = -\mathbf{D}^{-1}\mathbf{C}\mathbf{u} - \mathbf{D}^{-1}\mathbf{q}.\tag{2.2}$$

Hence we find the *underlying ODE* (cf. [1])

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})\mathbf{x} + \mathbf{p} - \mathbf{B}\mathbf{D}^{-1}\mathbf{q} =: \hat{\mathbf{A}}\mathbf{x} + \mathbf{g}.\tag{2.3}$$

Clearly one only needs to specify  $\mathbf{x}(0)$ , say

$$\mathbf{x}(0) = \mathbf{x}_0.\tag{2.4}$$

Stability constants are induced by the fundamental solution of (2.3),  $\mathbf{X}$  say, i.e.

$$\begin{aligned}\dot{\mathbf{X}} &= \hat{\mathbf{A}}\mathbf{X}, \\ \mathbf{X}(0) &= \mathbf{I}.\end{aligned}\tag{2.5}$$

Hence we define

$$\begin{aligned}\kappa_1 &:= \sup_{t \geq 0} \|\mathbf{X}(t)\|, \\ \kappa_2 &:= \sup_{t \geq 0} \int_0^t \|\mathbf{X}(t)\mathbf{X}^{-1}(s)\| ds.\end{aligned}\tag{2.6}$$

Here and in the sequel  $\|\cdot\|$  for a matrix or vector is some Hölder norm. If we denote the  $t$ -dependence explicitly, a quantity like  $\mathbf{x}(t)$  will always denote a vector. A vector function is denoted without an argument (if no confusion can occur) but provided with  $(\cdot)$ . We shall denote an  $L_p$  function norm for  $\mathbf{x}(\cdot)$  by  $\|\mathbf{x}(\cdot)\|$  or  $\|\mathbf{x}\|_p$ . Clearly  $\kappa_1$  and  $\kappa_2$  in (2.6) are the result of taking some function norm (cf. [2]). We now have

$$\|\mathbf{x}\|_\infty \leq \kappa_1 \|\mathbf{x}_0\| + \kappa_2 \|\mathbf{p} - \mathbf{B}\mathbf{D}^{-1}\mathbf{q}\|_\infty,\tag{2.7a}$$

$$\begin{aligned}\|\mathbf{y}\|_\infty &\leq \kappa_1 \|\mathbf{D}^{-1}\mathbf{C}\|_\infty \|\mathbf{x}_0\| + \kappa_2 \|\mathbf{D}^{-1}\mathbf{C}\|_\infty \|\mathbf{p} - \mathbf{B}\mathbf{D}^{-1}\mathbf{q}\|_\infty \\ &\quad + \|\mathbf{D}^{-1}\mathbf{q}\|_\infty.\end{aligned}\tag{2.7b}$$

Note that this estimate will grow unbounded if  $\mathbf{D}$  is somewhere close to a singular matrix. Yet, if  $\mathbf{D}$  is zero we merely have a higher index problem, e.g. index 2 if  $\mathbf{CB}$  is nonsingular. In the next section we reiterate the well known fact that the latter type of DAE is still well-posed (in a Hadamard sense), so that (2.7) cannot be meaningful for nearly singular  $\mathbf{D}$ .

### 3 Stability constants for DAE of index-2

The next relatively straightforward situation we discuss is the DAE

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{p}, \quad (3.1a)$$

$$\mathbf{0} = \mathbf{C}\mathbf{x} + \mathbf{q}. \quad (3.1b)$$

We now assume that  $\mathbf{CB}$  is nonsingular, i.e. the index equals 2. Differentiating (3.1a) and using (3.1b) we obtain an expression of  $\mathbf{y}$  in terms of  $\mathbf{x}$  and the source functions, viz.

$$\mathbf{y} = -(\mathbf{CB})^{-1}((\dot{\mathbf{C}} + \mathbf{CA})\mathbf{x} + \mathbf{C}\mathbf{p} + \dot{\mathbf{q}}). \quad (3.2)$$

Substituting this in (3.1a) we hence obtain as an underlying ODE

$$\dot{\mathbf{x}} = ((\mathbf{I} - \mathbf{P})\mathbf{A} - \dot{\mathbf{P}})\mathbf{x} + (\mathbf{I} - \mathbf{P})(\mathbf{p} - (\mathbf{A}\mathbf{F} - \dot{\mathbf{F}})\mathbf{q}). \quad (3.3)$$

Here we have defined

$$\mathbf{F} := \mathbf{B}(\mathbf{CB})^{-1}, \quad (3.4a)$$

$$\mathbf{P} := \mathbf{F}\mathbf{C}. \quad (3.4b)$$

Note that  $\mathbf{P}$  is a projector (cf. [7]). Actually, a more meaningful underlying ODE can be found as follows. Define the state variable

$$\mathbf{z} := (\mathbf{I} - \mathbf{P})\mathbf{x}. \quad (3.5)$$

This implies that  $\mathbf{z} = \mathbf{x} + \mathbf{F}\mathbf{q}$ . Clearly  $\mathbf{C}\mathbf{z} = \mathbf{0}$ . Therefore, a consistent initial value  $\mathbf{z}(0)$  can be found as:

$$\mathbf{z}(0) = \mathbf{z}_0 = (\mathbf{I} - \mathbf{P}(0))\mathbf{u}_0, \quad \text{for arbitrary } \mathbf{u}_0 \in \mathbb{R}^n. \quad (3.6)$$

This means that

$$\mathbf{x}(0) = \mathbf{x}_0 = (\mathbf{I} - \mathbf{P}(0))\mathbf{u}_0 - \mathbf{F}(0)\mathbf{q}(0), \quad (3.7)$$

and therefore,

$$\mathbf{P}(0)\mathbf{x}_0 = -\mathbf{F}(0)\mathbf{q}(0) \quad (3.8)$$

and

$$(\mathbf{I} - \mathbf{P}(0))\mathbf{x}_0 = (\mathbf{I} - \mathbf{P}(0))\mathbf{u}_0. \quad (3.9)$$

Hence, only the components  $(\mathbf{I} - \mathbf{P}(0))\mathbf{x}(0)$  of the initial vector  $\mathbf{x}(0)$  can be chosen arbitrarily. The DAE (3.1a),(3.1b) therefore has the following equivalent state ordinary differential equation (ODE)

$$\dot{\mathbf{z}} = ((\mathbf{I} - \mathbf{P})\mathbf{A} - \dot{\mathbf{P}})\mathbf{z} + (\mathbf{I} - \mathbf{P})(\mathbf{p} - (\mathbf{A}\mathbf{F} - \dot{\mathbf{F}})\mathbf{q}) =: \hat{\mathbf{A}}\mathbf{z} + (\mathbf{I} - \mathbf{P})\mathbf{g}. \quad (3.10)$$

We remark that  $\mathbf{z}$  does not depend on the derivative of  $\mathbf{q}$ . Let the fundamental solution matrix  $\mathbf{Z} \in \mathbb{R}^{n \times n}$  of ODE (3.3) be defined by

$$\dot{\mathbf{Z}} = \hat{\mathbf{A}}\mathbf{Z}, \quad (3.11a)$$

$$\mathbf{Z}(0) = \mathbf{I}. \quad (3.11b)$$

With (3.11b), the solution of ODE (3.3) can be expressed as

$$\mathbf{z}(t) = \mathbf{Z}(t)\mathbf{z}(0) + \int_0^t \mathbf{Z}(t)\mathbf{Z}^{-1}(s)(\mathbf{I} - \mathbf{P}(s))\mathbf{g}(s)ds. \quad (3.12)$$

Since  $\mathbf{z} = (\mathbf{I} - \mathbf{P})\mathbf{z}$ , we find from (3.12) that

$$\mathbf{z}(t) = (\mathbf{I} - \mathbf{P}(t))\mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0))\mathbf{u}_0 + \int_0^t (\mathbf{I} - \mathbf{P}(t))\mathbf{Z}(t)\mathbf{Z}^{-1}(s)(\mathbf{I} - \mathbf{P}(s))\mathbf{g}(s)ds. \quad (3.13)$$

For stability we only have to consider the growth behaviour in the subspace defined by  $\text{range}(\mathbf{I} - \mathbf{P})$ . We can prove the following

**Lemma 3.1**

(i).  $\mathbf{Z}$  satisfies

$$\mathbf{P}(t)\mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0)) = \mathbf{O}.$$

(ii). The matrix function  $\mathbf{Z}^{-1}$  satisfies

$$\mathbf{P}(0)\mathbf{Z}^{-1}(t)(\mathbf{I} - \mathbf{P}(t)) = \mathbf{O}.$$



*Proof.*

- (i). Equation (3.11a) implies  $\frac{d}{dt}(\mathbf{PZ}) = (\dot{\mathbf{P}} - \mathbf{P}\dot{\mathbf{P}})\mathbf{Z}$ . Since  $\dot{\mathbf{P}} = \mathbf{P}\dot{\mathbf{P}} + \dot{\mathbf{P}}\mathbf{P}$ , it follows that  $\frac{d}{dt}(\mathbf{PZ}) = \dot{\mathbf{P}}\mathbf{Z}$ . Therefore, the matrix function  $\mathbf{PZ}(\mathbf{I} - \mathbf{P}(0))$  satisfies the differential equation  $\dot{\mathbf{U}} = \dot{\mathbf{P}}\mathbf{U}$ , with  $\mathbf{U}(0) = \mathbf{O}$ . Hence,  $\mathbf{U}(t) = \mathbf{O}$ .
- (ii). The matrix function  $\mathbf{Z}^{-1}$  satisfies the differential equation  $\frac{d}{dt}(\mathbf{Z}^{-1}) = -\mathbf{Z}^{-1}\hat{\mathbf{A}}$ . Moreover,

$$\begin{aligned} \frac{d}{dt}(\mathbf{Z}^{-1}(\mathbf{I} - \mathbf{P})) &= -\mathbf{Z}^{-1}\hat{\mathbf{A}}(\mathbf{I} - \mathbf{P}) - \mathbf{Z}^{-1}\dot{\mathbf{P}} \\ &= -\mathbf{Z}^{-1}((\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{P}) - \dot{\mathbf{P}}(\mathbf{I} - \mathbf{P}) + \dot{\mathbf{P}}) \\ &= -\mathbf{Z}^{-1}((\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{P}) + \dot{\mathbf{P}}\mathbf{P}) \\ &= -\mathbf{Z}^{-1}(\mathbf{I} - \mathbf{P})(\mathbf{A}(\mathbf{I} - \mathbf{P}) + \dot{\mathbf{P}}). \end{aligned}$$

Hence the matrix function  $\mathbf{P}(0)\mathbf{Z}^{-1}(t)(\mathbf{I} - \mathbf{P}(t))$  satisfies the homogeneous differential equation above, with  $\mathbf{P}(0)\mathbf{Z}^{-1}(0)(\mathbf{I} - \mathbf{P}(0)) = \mathbf{O}$ . As a consequence,  $\mathbf{P}(0)\mathbf{Z}^{-1}(t)(\mathbf{I} - \mathbf{P}(t)) = \mathbf{O}$ .

□

### Corollary 3.2

(i).

$$(\mathbf{I} - \mathbf{P}(t))\mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0)) = \mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0)),$$

(ii).

$$(\mathbf{I} - \mathbf{P}(0))\mathbf{Z}^{-1}(t)(\mathbf{I} - \mathbf{P}(t)) = \mathbf{Z}^{-1}(t)(\mathbf{I} - \mathbf{P}(t)).$$

From Corollary 3.2 we deduce

$$\begin{aligned} (\mathbf{I} - \mathbf{P}(t))\mathbf{Z}(t)\mathbf{Z}^{-1}(s)(\mathbf{I} - \mathbf{P}(s)) &= (\mathbf{I} - \mathbf{P}(t))\mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0))\mathbf{Z}^{-1}(s)(\mathbf{I} - \mathbf{P}(s)) \\ &= \mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0))\mathbf{Z}^{-1}(s)(\mathbf{I} - \mathbf{P}(s)) \\ &= \mathbf{Z}(t)\mathbf{Z}^{-1}(s)(\mathbf{I} - \mathbf{P}(s)). \end{aligned} \tag{3.14}$$

We can now define a "restricted" fundamental solution of (3.3) in:

**Definition 3.3**  $\mathbf{W}(t) = (\mathbf{I} - \mathbf{P}(t))\mathbf{Z}(t)(\mathbf{I} - \mathbf{P}(0))$ , with  $\mathbf{W}(0) = (\mathbf{I} - \mathbf{P}(0))$ .

For  $\mathbf{W}^+$  see [7];  $\mathbf{W}^+(t) = (\mathbf{I} - \mathbf{P}(0))\mathbf{Z}^{-1}(t)(\mathbf{I} - \mathbf{P}(t))$ . We finally obtain

**Theorem 3.4** *The solution  $\mathbf{z}$  of the state equation (3.3) can be written as*

$$\mathbf{z}(t) = \mathbf{W}(t)(\mathbf{I} - \mathbf{P}(0))\mathbf{u}_0 + \int_0^t \mathbf{W}(t)\mathbf{W}^+(s)(\mathbf{I} - \mathbf{P}(s))\mathbf{g}(s)ds. \quad (3.15)$$

*Proof.* Equation (3.11a) implies  $\frac{d}{dt}(\mathbf{P}\mathbf{Z}) = (\dot{\mathbf{P}} - \mathbf{P}\dot{\mathbf{P}})\mathbf{Z}$ . Therefore, the matrix function  $\mathbf{P}\mathbf{Z}(\mathbf{I} - \mathbf{P}(0))$  satisfies the differential equation  $\dot{\mathbf{U}} = \dot{\mathbf{P}}\mathbf{U}$ , with  $\mathbf{U}(0) = \mathbf{O}$ . Hence,  $\mathbf{U}(t) = \mathbf{O}$ . Since  $\mathbf{Z}$  satisfies  $\dot{\mathbf{Z}} = \hat{\mathbf{A}}\mathbf{Z}$ , the matrix function  $\mathbf{W}$  satisfies  $\dot{\mathbf{W}} = \hat{\mathbf{A}}\mathbf{W}$ .  $\square$

From (3.15), we deduce the following estimates for the growth behaviour of  $\mathbf{z}$  and  $\mathbf{x}$ , i.e.

$$\|\mathbf{z}\|_\infty \leq \hat{\kappa}_1 \|(\mathbf{I} - \mathbf{P}(0))\mathbf{x}_0\| \quad (3.16)$$

$$+ \hat{\kappa}_2 (\|(\mathbf{I} - \mathbf{P})\mathbf{p}\|_\infty + \|(\mathbf{I} - \mathbf{P})\mathbf{A}\mathbf{F}\mathbf{q}\|_\infty + \|(\mathbf{I} - \mathbf{P})\dot{\mathbf{F}}\mathbf{q}\|_\infty) \quad (3.17)$$

and

$$\|\mathbf{x}\|_\infty \leq \hat{\kappa}_1 (\|(\mathbf{I} - \mathbf{P}(0))\mathbf{x}_0\| + \|\mathbf{F}(0)\mathbf{q}(0)\|) \quad (3.18)$$

$$+ \hat{\kappa}_2 (\|(\mathbf{I} - \mathbf{P})\mathbf{p}\|_\infty + \|(\mathbf{I} - \mathbf{P})\mathbf{A}\mathbf{F}\mathbf{q}\|_\infty + \|(\mathbf{I} - \mathbf{P})\dot{\mathbf{F}}\mathbf{q}\|_\infty) \quad (3.19)$$

$$+ \|\mathbf{F}\mathbf{q}\|_\infty, \quad (3.20)$$

respectively, where the conditioning constants are defined as

$$\begin{aligned} \hat{\kappa}_1 &:= \sup\{\|\mathbf{W}(t)\|, t \in I\}, \quad \text{and,} \\ \hat{\kappa}_2 &:= \sup\{[\int_0^t \|\mathbf{W}(t)\mathbf{W}^+(s)\|ds], t \in I\}. \end{aligned} \quad (3.21)$$

Hence, from (3.2)

$$\|\mathbf{y}\|_\infty \leq \|(\mathbf{C}\mathbf{B})^{-1}\|_\infty (\|\dot{\mathbf{C}}\|_\infty + \|\mathbf{C}\mathbf{A}\|_\infty)\|\mathbf{x}\|_\infty + \|\mathbf{C}\mathbf{p}\|_\infty + \|\dot{\mathbf{q}}\|_\infty. \quad (3.22)$$

Again, if  $\mathbf{C}\mathbf{B}$  is nearly singular these estimates are no longer meaningful.

## 4 A scalar 'ill-conditioned' index-1 case

In Section 2 we saw that the index-1 estimates exhibit ill-conditioning if the matrix  $\mathbf{D}$  is nearly singular. However, this is not sustained for a limiting case where  $\mathbf{D} = \mathbf{O}$ . In the latter situation we now will show that the index-2 problem may still be regarded as well conditioned, provided some quantities like  $\|(\mathbf{C}\mathbf{B})^{-1}\|_\infty$  are not large.

In this section we shall give an explanation of this phenomenon for a DAE where all matrices involved are scalar and constant. So consider

$$\begin{aligned} \dot{x} &= ax + by + p, \\ 0 &= cx + dy + q, \end{aligned} \quad (4.1)$$

where  $d \neq 0$ . For this simple system the governing state ODE reads

$$\begin{aligned}\dot{x} &= (a - bd^{-1}c)x + p - bd^{-1}q, \\ &= \bar{a}x + \bar{p},\end{aligned}\tag{4.2}$$

where  $\bar{a}$  and  $\bar{p}$  are defined similarly to  $\bar{A}$  and  $\bar{p}$ , respectively. The corresponding solution reads

$$x(t) = e^{\bar{a}t}x(0) + \int_0^t e^{\bar{a}(t-\tau)} \bar{p}(\tau) d\tau.\tag{4.3}$$

Let  $d = \varepsilon$ , ( $\varepsilon \downarrow 0$ ), and note that

$$\begin{aligned}\int_0^t e^{\bar{a}(t-\tau)} b\varepsilon^{-1}q(\tau) d\tau &= -\bar{a}^{-1} (b\varepsilon^{-1}q(t) - e^{\bar{a}t}b\varepsilon^{-1}q(0)) \\ &\quad + \bar{a}^{-1} \int_0^t e^{\bar{a}(t-\tau)} b\varepsilon^{-1}\dot{q}(\tau) d\tau.\end{aligned}\tag{4.4}$$

Assume that  $|cb|$  is bounded away from zero. Using (4.4) in (4.3), we obtain

$$\begin{aligned}x(t) &= e^{\bar{a}t}x(0) + \int_0^t e^{\bar{a}(t-\tau)} p(\tau) d\tau \\ &\quad - \frac{1}{c} \left(1 + \frac{a}{bc}\varepsilon + \mathcal{O}(\varepsilon^2)\right) (q(t) - e^{\bar{a}t}q(0) - \int_0^t e^{\bar{a}(t-\tau)} \dot{q}(\tau) d\tau) \\ &\quad + \mathcal{O}(\varepsilon^2).\end{aligned}\tag{4.5}$$

The growth behaviour of  $x$  is therefore characterized by the following estimate

$$\begin{aligned}|x(t)| &\leq \kappa_1 |x(0)| + \varepsilon \hat{\kappa}_2 \max_{0 \leq \tau \leq t} |p(\tau)| \\ &\quad + \frac{1}{|c|} \left(1 + \left|\frac{a}{bc}\right| \varepsilon + \mathcal{O}(\varepsilon^2)\right) \left(\max_{0 \leq \tau \leq t} |q(\tau)| + \kappa_1 |q(0)| + \varepsilon \hat{\kappa}_2 \max_{0 \leq \tau \leq t} |\dot{q}(\tau)|\right)\end{aligned}\tag{4.6}$$

where the constant  $\kappa_1$  is defined like in (2.6), i.e.

$$\kappa_1 := \max_{0 \leq t \leq T} e^{(a - b\varepsilon^{-1}c)t}.\tag{4.7}$$

The conditioning constant  $\kappa_2$  (cf. (2.6)) can be seen to be of order of magnitude  $\varepsilon$ . For that reason a more useful quantity is  $\hat{\kappa}_2$ , defined by

$$\hat{\kappa}_2 := d^{-1} \max_{0 \leq t \leq T} \int_0^t e^{(a - bd^{-1}c)(t-\tau)} d\tau.\tag{4.8}$$

It can easily be seen that  $\hat{\kappa}_2$  is an  $\mathcal{O}(1)$  constant since  $d = \varepsilon$  ( $\varepsilon \downarrow 0$ ).

**Remark 4.1** (4.5) shows that  $x$  exhibits initial layer behaviour if  $d = \varepsilon$ ,  $\varepsilon$  small. Furthermore, it shows that  $x$  effectively depends on  $q$  only for  $\varepsilon$ ,  $\varepsilon$  small.

From equation (4.1) we find that  $y$  satisfies

$$y = -d^{-1}cx - d^{-1}q. \quad (4.9)$$

Substitution of (4.3) gives

$$y = -d^{-1}ce^{\bar{a}t}x(0) - d^{-1}c \int_0^t e^{\bar{a}(t-\tau)} \bar{p}(\tau) d\tau - d^{-1}q(t). \quad (4.10)$$

For  $d = \varepsilon$ ,  $\varepsilon$  small, partial integration yields

$$\begin{aligned} y(t) = & -\varepsilon^{-1}ce^{\bar{a}t}x(0) - \varepsilon^{-1}q(t) - \varepsilon^{-1}c \int_0^t e^{\bar{a}(t-\tau)} p(\tau) d\tau \\ & - \frac{\varepsilon^{-1}c}{\bar{a}} (b\varepsilon^{-1}q(t) - e^{\bar{a}t}b\varepsilon^{-1}q(0)) \\ & + \frac{\varepsilon^{-1}c}{\bar{a}} \int_0^t e^{\bar{a}(t-\tau)} b\varepsilon^{-1}\dot{q}(\tau) d\tau. \end{aligned} \quad (4.11)$$

First order expansion in  $\varepsilon$  results in

$$\begin{aligned} y(t) \doteq & e^{\bar{a}t} \left( y(0) - q(0) \left( \frac{a}{bc} + \mathcal{O}(\varepsilon) \right) \right) - \varepsilon^{-1}c \int_0^t e^{\bar{a}(t-\tau)} p(\tau) d\tau \\ & + \left( \frac{a}{bc} + \mathcal{O}(\varepsilon) \right) q(t) - \varepsilon^{-1} \left( 1 + \frac{a}{bc}\varepsilon + \mathcal{O}(\varepsilon^2) \right) \int_0^t e^{\bar{a}(t-\tau)} \dot{q}(\tau) d\tau. \end{aligned} \quad (4.12)$$

So, the algebraic variable  $y$  can be bounded as follows

$$\begin{aligned} |y(t)| \leq & \kappa_1 \left( |y(0)| + |q(0)| \left( \left| \frac{a}{bc} \right| + \mathcal{O}(\varepsilon) \right) \right) + |c|\hat{\kappa}_2 \max_{0 \leq \tau \leq t} |p(\tau)| \\ & + \left( \left| \frac{a}{bc} \right| + \mathcal{O}(\varepsilon) \right) \max_{0 \leq \tau \leq t} |q(\tau)| + \left( 1 + \left| \frac{a}{bc} \right| \varepsilon + \mathcal{O}(\varepsilon^2) \right) \hat{\kappa}_2 \max_{0 \leq \tau \leq t} |\dot{q}(\tau)| \\ & (\varepsilon \downarrow 0). \end{aligned} \quad (4.13)$$

For  $d$  bounded away from zero,  $\hat{\kappa}_2$  and  $\kappa_2$  are quite similar (as  $d\hat{\kappa}_2 = \kappa_2$ ), the solution  $x$  (cf. (4.3)) then can be estimated as

$$|x(t)| \leq \kappa_1|x(0)| + \kappa_2 \max_{0 \leq \tau \leq t} |p(\tau)| + \kappa_2|b|d^{-1} \max_{0 \leq \tau \leq t} |q(\tau)|. \quad (4.14)$$

For the algebraic variable  $y$ , we find from (4.10)

$$\begin{aligned} |y(t)| \leq & \kappa_1|d^{-1}c||x(0)| + \kappa_2|d^{-1}c| \max_{0 \leq \tau \leq t} |p(\tau)| \\ & + \kappa_2|d^{-1}c||bd^{-1}| \max_{0 \leq \tau \leq t} |q(\tau)| + |d^{-1}| \max_{0 \leq \tau \leq t} |q(\tau)|. \end{aligned} \quad (4.15)$$

## 5 Conditioning of index-1 problems where $\mathbf{D} = \varepsilon \mathbf{I}$

The estimate (4.13) shows that  $\mathbf{y}$  has a layer when  $d = \varepsilon$ , ( $\varepsilon \rightarrow 0$ ) (and  $|(\mathbf{cb})^{-1}|$  is not large). Moreover,  $\mathbf{y}$  effectively depends on  $\dot{\mathbf{q}}$  in such a situation; of course, this is reasonable since it resembles an index-2 case. We may also refer to papers like [13] where this sort of ill-conditioning was also noted. We shall now consider the more general matrix situation, where  $\mathbf{D} = \varepsilon \mathbf{I}$  (thus being an extreme form of a "nearly singular matrix"  $\mathbf{D}$ ). So let

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{y} + \mathbf{q}.\end{aligned}\tag{5.1}$$

Recalling (2.3) we see that the system (2.3), where

$$\hat{\mathbf{A}} := \mathbf{A} - \varepsilon^{-1}\mathbf{B}\mathbf{C},\tag{5.2}$$

induces a fundamental solution of the underlying equation, which is apparently stiff if  $\varepsilon \rightarrow 0$ . We now assume that  $\mathbf{C}\mathbf{B}$  is nonsingular. Then it is no restriction to also assume that  $\mathbf{B}$  has the form

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{B}_{11}(t) \in \mathbb{R}^{m^2} \text{ nonsingular}.\tag{5.3}$$

Hence the underlying equation itself is a two deck system with  $m$  fast and  $n - m$  slow components, say  $\mathbf{x}_1$  and  $\mathbf{x}_2$  respectively, where the homogeneous part reads

$$\begin{aligned}\varepsilon \dot{\mathbf{x}}_1 &= (\varepsilon \mathbf{A}_{11} - \mathbf{B}_{11} \mathbf{C}_{11}) \mathbf{x}_1 + (\varepsilon \mathbf{A}_{12} - \mathbf{B}_{11} \mathbf{C}_{12}) \mathbf{x}_2 \\ \mathbf{x}_2 &= \varepsilon \mathbf{A}_{21} \mathbf{x}_1 + \varepsilon \mathbf{A}_{22} \mathbf{x}_2.\end{aligned}\tag{5.4}$$

Here  $\hat{\mathbf{A}}$  has been partitioned in blocks commensurating with  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  ( $\hat{\mathbf{A}}_{11}(t) \in \mathbb{R}^{m^2}$ , etc.),  $\mathbf{C}_{11}(t) \in \mathbb{R}^{m^2}$ ,  $\mathbf{C}_{12}(t) \in \mathbb{R}^{m \times (n-m)}$ . For a relationship between DAE and singularly perturbed problems see e.g. [10]. Using a Riccati transformation, cf. [2], one can easily see that (5.4) can be transformed into a decoupled system with  $m$  fast and  $n - m$  slow modes. The latter actually do not play a role (they vanish in a regular way when  $\varepsilon \rightarrow 0$ ). Hence we conclude that it is not restrictive for our analysis to let  $m = n$ , which in turn is a straightforward generalisation of the scalar case dealt with in Section 4. Hence we shall consider the underlying ODE

$$\dot{\mathbf{x}} = (\mathbf{A} - \varepsilon^{-1}\mathbf{B}\mathbf{C})\mathbf{x} + \mathbf{p} - \varepsilon^{-1}\mathbf{B}\mathbf{q}.\tag{5.5a}$$

The variable  $\mathbf{y}$  then follows from

$$\mathbf{y} = -\varepsilon^{-1}(\mathbf{C}\mathbf{x} + \mathbf{q}).\tag{5.5b}$$

From (5.5a) we see that if  $\|(\mathbf{BC})^{-1}\|$  is uniformly bounded, the matrix  $\varepsilon^{-1}\mathbf{BC}$  essentially determines the growth properties of the modes. Now assume, for some  $\tilde{\kappa}_2, \sigma > 0$

$$\exp\left[-\int_s^t (-\varepsilon^{-1}\mathbf{B}(\tau)\mathbf{C}(\tau)) d\tau\right] \leq \tilde{\kappa}_2 \exp\left[\frac{\sigma}{\varepsilon}(s-t)\right]. \quad (5.6)$$

This merely expresses exponential decay. We now have

**Property 5.1** *Let  $\|(\mathbf{CB})^{-1}\|, \|\mathbf{A}\|$  be uniformly bounded and let (5.6) hold. Then  $\hat{\kappa}_2 := \varepsilon^{-1}\kappa_2$  is uniformly bounded in  $\varepsilon$ .*

*Proof.* For  $\varepsilon$  small enough we see that

$$\|\mathbf{X}(t)\mathbf{X}^{-1}(s)\| = \left\| \int_s^t (\mathbf{A}(\tau) - \varepsilon^{-1}\mathbf{B}(\tau)\mathbf{C}(\tau)) d\tau \right\| \leq \tilde{\kappa}_2 \exp\left[\frac{\sigma}{\varepsilon}(s-t)\right].$$

One should note that  $\|(\mathbf{CB})^{-1}\|$  bounded, implies the same for  $\|(\mathbf{BC})^{-1}\|$ . Hence  $\varepsilon^{-1}\|\mathbf{X}(t)\mathbf{X}^{-1}(s)\| \leq \tilde{\kappa}_2[1 - \exp(-\frac{\sigma}{\varepsilon}t)]$ .  $\square$

**Theorem 5.2** *Under the assumptions of Property 5.1 the following estimates hold to first order in  $\varepsilon$*

$$\begin{aligned} \|\mathbf{x}\| &\leq \kappa_1 \|\mathbf{x}(0)\| + \kappa_2 \|\mathbf{p}\| + (1 + \kappa_1) \|\mathbf{B}\| \|(\mathbf{CB})^{-1}\| \|\mathbf{q}\| + \kappa_2 (\kappa_3 \|\mathbf{q}\| + \|\dot{\mathbf{q}}\|), \\ \|\mathbf{y}\| &\leq \kappa_1 \|\mathbf{C}\| \|(\mathbf{C}(0))^{-1}\| \|\mathbf{y}(0)\| + \hat{\kappa}_2 \|\mathbf{C}\| \|\mathbf{p}\| + \hat{\kappa}_2 \|\mathbf{C}\| (\kappa_3 \|\mathbf{q}\| + \|\mathbf{B}\| \|(\mathbf{CB})^{-1}\| \|\dot{\mathbf{q}}\|). \end{aligned}$$

Here  $\hat{\kappa}_2 := \varepsilon^{-1}\kappa_2, \kappa_3 := \sup\left\|\frac{d}{dt}(\mathbf{B}(t)\mathbf{C}(t))^{-1}\mathbf{B}(t)\right\|$ .

*Proof.* From (5.1), (5.2) we find

$$\begin{aligned} \mathbf{x}(t) &= \mathbf{X}(t)\mathbf{x}(0) + \int_0^t \mathbf{X}(t)\mathbf{X}^{-1}(s) [\mathbf{p}(s) - \varepsilon^{-1}\mathbf{B}(s)\mathbf{q}(s)] ds \\ &=: \mathbf{X}(t)\mathbf{x}(0) + \int_0^t \mathbf{X}(t)\mathbf{X}^{-1}(s)\mathbf{p}(s) ds + I(t). \end{aligned}$$

Now we remark that

$$\mathbf{X}^{-1}(s) = -\frac{d}{dt}[\mathbf{X}^{-1}(s)]\hat{\mathbf{A}}^{-1}(s).$$

Hence,

$$\begin{aligned}
I(t) &= \varepsilon^{-1} \mathbf{X}(t) \int_0^t \frac{d}{ds} [\mathbf{X}^{-1}(s)] \hat{\mathbf{A}}^{-1}(s) \mathbf{B}(s) \mathbf{q}(s) ds \\
&= \varepsilon^{-1} \hat{\mathbf{A}}^{-1}(t) \mathbf{B}(t) \mathbf{q}(t) - \varepsilon^{-1} \mathbf{X}(t) \hat{\mathbf{A}}^{-1}(0) \mathbf{B}(0) \mathbf{q}(0) \\
&\quad - \varepsilon^{-1} \mathbf{X}(t) \int_0^t \mathbf{X}^{-1}(s) \left\{ \left[ \frac{d}{ds} (\hat{\mathbf{A}}^{-1}(s) \mathbf{B}(s)) \right] \mathbf{q}(s) + \hat{\mathbf{A}}^{-1}(s) \mathbf{B}(s) \dot{\mathbf{q}}(s) \right\} ds.
\end{aligned}$$

Since  $\varepsilon^{-1} \hat{\mathbf{A}}^{-1} = (\mathbf{BC})^{-1} (\mathbf{I} - \varepsilon \mathbf{A}(\mathbf{BC})^{-1})$  and  $(\mathbf{BC})^{-1} \mathbf{B} = \mathbf{B}(\mathbf{CB})^{-1}$ , we find neglecting  $\mathcal{O}(\varepsilon)$  terms

$$\begin{aligned}
\|I\| &\leq \|\mathbf{B}(\mathbf{CB})^{-1}\| \|\mathbf{q}\| + \kappa_1 \|\mathbf{B}(0)(\mathbf{C}(0)\mathbf{B}(0))^{-1}\| \|\mathbf{q}(0)\| \\
&\quad + \kappa_2 (\kappa_3 \|\mathbf{q}\| + \|\mathbf{B}(\mathbf{CB})^{-1}\| \|\dot{\mathbf{q}}\|).
\end{aligned}$$

From this the estimate for  $\mathbf{x}$  follows.

For  $\mathbf{y}$  we obtain (cf. (2.2))

$$\mathbf{y}(t) = -\varepsilon^{-1} \mathbf{C}(t) \mathbf{x}(t) - \varepsilon^{-1} \mathbf{q}(t).$$

Now we remark that

$$\begin{aligned}
-\varepsilon^{-1} \mathbf{C} \hat{\mathbf{A}}^{-1} \mathbf{B} &\doteq \mathbf{C}(\mathbf{BC})^{-1} (\mathbf{I} + \varepsilon (\mathbf{BC})^{-1} \mathbf{A}) \mathbf{B} \\
&\doteq \mathbf{C}(\mathbf{BC})^{-1} \mathbf{B} = \mathbf{CB}(\mathbf{CB})^{-1} = \mathbf{I}.
\end{aligned}$$

(Here we have neglected terms of order  $\varepsilon$ ). Hence we see that the second term in  $\mathbf{y}(t)$  is approximately cancelled by the first term arising from the contribution of  $I$  in  $\mathbf{x}(t)$ . In summary we thus have, neglecting higher order terms,

$$\begin{aligned}
\mathbf{y}(t) &\doteq -\varepsilon^{-1} [\mathbf{C}(t) \mathbf{X}(t) \mathbf{x}(0) + \mathbf{C}(t) \int_0^t \mathbf{X}(t) \mathbf{X}^{-1}(s) \mathbf{p}(s) ds] \\
&\quad - \varepsilon^{-1} \mathbf{C}(t) \mathbf{X}(t) (\mathbf{B}(0) \mathbf{C}(0))^{-1} \mathbf{B}(0) \mathbf{q}(0) \\
&\quad - \varepsilon^{-2} \mathbf{C}(t) \mathbf{X}(t) \int_0^t \mathbf{X}^{-1}(s) \left\{ \left[ \frac{d}{ds} (\hat{\mathbf{A}}^{-1}(s) \mathbf{B}(s)) \right] \mathbf{q}(s) + \hat{\mathbf{A}}^{-1}(s) \mathbf{B}(s) \dot{\mathbf{q}}(s) \right\} ds
\end{aligned}$$

One should note that the  $\varepsilon^{-2}$  factor for the last integral is in fact moderated by a factor from  $\hat{\mathbf{A}}^{-1}$ . Its actual value is controlled by the fast fundamental solutions, giving an  $\mathcal{O}(1)$  contribution after integration. Using the fact that

$$\begin{aligned}
\mathbf{y}(0) &= -\varepsilon^{-1} [\mathbf{C}(0) \mathbf{x}(0) + \mathbf{q}(0)] \\
&= -\varepsilon^{-1} \mathbf{C}(0) [\mathbf{x}(0) + (\mathbf{B}(0) \mathbf{C}(0))^{-1} \mathbf{B}(0) \mathbf{q}(0)]
\end{aligned}$$

the result trivially follows. Note that  $\mathbf{CB}$  and  $\mathbf{BC}$  have the same eigenvalues. Hence, the condition that  $\|(\mathbf{CB})^{-1}\|$  is reasonably bounded assures that all solutions of the underlying ODE are fast indeed.  $\square$

## 6 Conditioning of index-2 problems where $\mathbf{CB} = \varepsilon \mathbf{I}$

Like in the index-1 case we can study the conditioning constants for the index-2 problems. In Section 3 we have derived an underlying ODE (essential underlying ODE, cf. [1]) for  $\mathbf{z}$  rather than for  $\mathbf{x}$ . We first remark that we now have to require  $m < n$  if  $\mathbf{CB} = \varepsilon \mathbf{I}$ . Indeed for if  $m = n$  we have  $\mathbf{BC} = \mathbf{CB}$  and so  $\mathbf{P} = \mathbf{I}$  (cf. (3.4b)) for constant  $\mathbf{C}, \mathbf{B}$ . Hence  $\mathbf{z} = \mathbf{x} = -\mathbf{F}(0)\mathbf{q}(0)$  is constant and we find

$$\mathbf{y} = \mathbf{B}^{-1}\mathbf{A}\mathbf{B}(0)(\mathbf{C}(0)\mathbf{B}(0))^{-1}\mathbf{q}(0) - \mathbf{B}^{-1}\mathbf{p} - (\mathbf{CB})\dot{\mathbf{q}}. \quad (6.1)$$

If e.g.  $\mathbf{B} = \varepsilon \mathbf{I}$  and  $\mathbf{C} = \mathbf{I}$  we see that  $\mathbf{y}$  is not uniformly bounded in  $\varepsilon$  and instability would occur. We return to this problem in Section 9. Hence  $m < n$

Turning now to (3.3) we derive

$$\begin{aligned} \dot{\mathbf{x}} &= [(\mathbf{I} - \mathbf{P})\mathbf{A} - \mathbf{F}\dot{\mathbf{C}}]\mathbf{x} + (\mathbf{I} - \mathbf{P})\mathbf{p} - \mathbf{F}\dot{\mathbf{q}} \\ &= \hat{\mathbf{A}}\mathbf{x} + \hat{\mathbf{g}}. \end{aligned} \quad (6.2)$$

Here we have used  $\varepsilon^{-1}\dot{\mathbf{B}}\mathbf{C}\mathbf{x} = -\varepsilon^{-1}\dot{\mathbf{B}}\dot{\mathbf{q}}$ .

One should note that the system matrix  $\hat{\mathbf{A}}$  is precisely the one encountered in (3.10) for the restricted variable  $\mathbf{z}$ . As may be deduced from the analysis of the previous section the fact that the matrices  $\mathbf{A}, \mathbf{B}$ , etc., are varying is of minor importance. The same applies here too, unless  $\|\dot{\mathbf{P}}\|$  is large compared to  $\|(\mathbf{I} - \mathbf{P})\mathbf{A}\|$ . This case needs further research which is not carried out here. Hence we shall restrict ourselves to constant coefficient problems, implying  $\dot{\mathbf{P}} = \mathbf{O}$ ,  $\frac{d}{dt}(\mathbf{F}\mathbf{q}) = \varepsilon^{-1}\mathbf{B}\dot{\mathbf{q}}$ . Since we may transform the variable  $\mathbf{x}$  in (6.2) such that  $\mathbf{P}$  (in fact  $\mathbf{B}$ ) has zeros in the last  $n - m$  rows, we see that (6.2) is in fact a two deck system with  $m$  fast and  $n - m$  slow modes, not much different from (5.4) with  $\mathbf{C}$  replaced by  $\mathbf{CA}$ . In principle we can decouple (6.2) by e.g. Riccati transformation. As a consequence the slow modes can be found separately, after which we end up with an  $m$ -th order system having fast modes only. The essential part of the analysis is therefore greatly simplified by considering the  $m$ -th order "underlying ODE"

$$\dot{\mathbf{x}} = (\mathbf{A} - \varepsilon^{-1}\mathbf{BCA})\mathbf{x} + (\mathbf{I} - \varepsilon^{-1}\mathbf{BC})\mathbf{p} - \varepsilon^{-1}\mathbf{B}\dot{\mathbf{q}}, \quad (6.3a)$$

where  $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{m^2}$ . For the variable  $\mathbf{y}$  we moreover have (cf. (3.2))

$$\mathbf{y} = -\varepsilon^{-1}(\mathbf{CA}\mathbf{x} + \mathbf{C}\mathbf{p} + \dot{\mathbf{q}}). \quad (6.3b)$$

We can now simply apply Theorem 5.2 to obtain



**Theorem 6.1** For (6.3a), (6.3b) the following estimates hold

$$\begin{aligned} \|\mathbf{x}\| &\leq \kappa_1 \|\mathbf{x}(0)\| + ((1 + \kappa_1) \|\mathbf{B}\| \|(\mathbf{CAB})^{-1}\| \|\mathbf{C}\| + \kappa_2 (1 + \kappa_3 \|\mathbf{C}\|)) \|\mathbf{p}\| \\ &\quad + ((1 + \kappa_1) \|\mathbf{B}\| \|(\mathbf{CAB})^{-1}\| + \kappa_2 \kappa_3) \|\dot{\mathbf{q}}\| + \kappa_2 \left\| \frac{d}{dt}(\mathbf{Cp}) \right\| + \kappa_2 \|\ddot{\mathbf{q}}\|, \\ \|\mathbf{y}\| &\leq \kappa_1 \|\mathbf{CA}\| \|(\mathbf{CA})^{-1}\| \|\mathbf{y}(0)\| + \|\mathbf{CA}\| \{\hat{\kappa}_2 + \kappa_2 \kappa_3\} (\|\mathbf{C}\| \|\mathbf{p}\| + \|\dot{\mathbf{q}}\|) \\ &\quad + \kappa_2 \|\mathbf{CA}\| \|(\mathbf{CAB})^{-1}\| \left\{ \left\| \frac{d}{dt}(\mathbf{Cp}) \right\| + \|\ddot{\mathbf{q}}\| \right\}. \end{aligned}$$

*Proof.* For deriving the estimate one should replace the matrix  $\mathbf{C}(t)$  in the proof of Theorem 5.2 by  $\mathbf{CA}$  and  $\mathbf{q}$  by  $\mathbf{Cp} + \dot{\mathbf{q}}$ . The estimate for  $\mathbf{x}$  follows from identifying (6.3a) with (5.4) noting that  $(\mathbf{BCA})^{-1}\mathbf{B} = \mathbf{B}(\mathbf{CAB})^{-1}$ . For  $\mathbf{y}$  we obtain in particular

$$\begin{aligned} \mathbf{y}(t) &\doteq -\varepsilon^{-1} \mathbf{CA} \mathbf{X}(t) (\mathbf{CA})^{-1} \mathbf{y}(0) - \varepsilon^{-1} \mathbf{CA} \int_0^t \mathbf{X}(t) \mathbf{X}^{-1}(s) (\mathbf{Cp}(s) + \dot{\mathbf{q}}(s)) ds \\ &\quad + \mathbf{CA} \int_0^t \mathbf{X}(t) \mathbf{X}^{-1}(s) \left[ \frac{d}{dt} (\varepsilon \mathbf{B} (\mathbf{CAB})^{-1}) \{ \mathbf{Cp} + \dot{\mathbf{q}} \} \right. \\ &\quad \left. + \mathbf{B} (\mathbf{CAB})^{-1} \left( \frac{d}{dt} (\mathbf{Cp}(s)) + \ddot{\mathbf{q}}(s) \right) \right] ds. \end{aligned}$$

From which the second estimate more specifically follows.  $\square$

## 7 The effective index; effect of errors

As we already saw in the analysis of Sections 5 and 6, an underlying state equation may be a two deck system itself, thus inducing a further reduction to isolate the fast solutions. We may continue this process for nearly singular  $\mathbf{CAB}$ , etc., which we omit, however. From such a decoupling into fast and slow modes we see that the notion of index belongs to a *variable*, rather than to a *problem* (the index of the problem is then the highest index a variable may have). More than that, since this index notion is strongly depending on invertibility of some matrix (function) it is obvious that the conditioning of the latter matrix must determine whether the *effective index* is at least one higher.

A practical quantification of this notion should be related to a natural threshold for the problem, like inherent errors (data errors or rounding errors) or discretisation errors. We will illustrate this problem by analyzing the iteration matrix of e.g. Euler Backward (or any other BDF method, cf. also [3]). Consider the DAE

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{Ax} + \mathbf{By} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{Cx} + \mathbf{Dy} + \mathbf{q}, \end{aligned} \tag{7.1}$$

The iteration matrix will be of the following form

$$\mathbf{J}_1 := \begin{bmatrix} \mathbf{I} - h\mathbf{A} & -h\mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}, \quad (7.2)$$

where  $h$  is the time step. The inverse of  $\mathbf{J}_1$  can be computed through LU-decomposition, i.e.

$$\mathbf{J}_1 = \mathbf{L}\mathbf{U} := \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{I} - h\mathbf{A} & -h\mathbf{B} \\ \mathbf{O} & h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B} + \mathbf{D} \end{bmatrix}. \quad (7.3)$$

For  $\mathbf{U}^{-1}$  we need the inverse of the matrix  $h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$ .

Let us first assume that the matrix  $\mathbf{D}$  is such that  $\|\mathbf{D}^{-1}\|$  is moderate. Then we find

$$(h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B} + \mathbf{D})^{-1} = \mathbf{D}^{-1} - h\mathbf{D}^{-1}\mathbf{C}\mathbf{B}\mathbf{D}^{-1} + \mathcal{O}(h^2). \quad (7.4)$$

This implies that

$$\mathbf{J}_1^{-1} = \begin{bmatrix} \mathbf{I} + h(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}) & h\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{I} + h(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})) & (\mathbf{I} - h\mathbf{D}^{-1}\mathbf{C}\mathbf{B})\mathbf{D}^{-1} \end{bmatrix} + \mathcal{O}(h^2). \quad (7.5)$$

Notice the important role of the matrices  $\mathbf{D}^{-1}\mathbf{C}$  and  $\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$  in the equation above. Hence, rounding errors proportional to the *machine constant*  $\eta$  will be introduced both in  $\mathbf{x}$  and in  $\mathbf{y}$ , while solving this linear system. This means that such a DAE behaves as well conditioned as an ODE. Here and in the sequel the latter notion has to be understood in an absolute (i.e. not relative) sense. If, on the other hand,  $\|\mathbf{D}^{-1}\|$  is not small, problems arise. If we take  $\mathbf{D} = \varepsilon\mathbf{I}$ , ( $\varepsilon \rightarrow 0$ ), then we find for small fixed  $h$  that

$$\begin{aligned} (h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B} + \mathbf{D})^{-1} &= h^{-1}(\mathbf{C}\mathbf{B})^{-1}(\mathbf{I} - h\mathbf{C}\mathbf{A}\mathbf{B}(\mathbf{C}\mathbf{B})^{-1}) \\ &+ \mathcal{O}(h) + \mathcal{O}(\varepsilon/h^2), (\varepsilon \rightarrow 0), \end{aligned} \quad (7.6)$$

assuming  $(\mathbf{C}\mathbf{B})^{-1}$  is bounded. We thus find to first order in  $\varepsilon$

$$\mathbf{J}_1^{-1} \doteq \begin{bmatrix} (\mathbf{I} - \mathbf{P})(\mathbf{I} + h\mathbf{A}(\mathbf{I} - \mathbf{P})) & (\mathbf{I} + h(\mathbf{I} - \mathbf{P})\mathbf{A})\mathbf{B}(\mathbf{C}\mathbf{B})^{-1} \\ -h^{-1}(\mathbf{C}\mathbf{B})^{-1}\mathbf{C}(\mathbf{I} + h\mathbf{A}(\mathbf{I} - \mathbf{P})) & h^{-1}(\mathbf{C}\mathbf{B})^{-1}(\mathbf{I} - h\mathbf{C}\mathbf{A}\mathbf{B}(\mathbf{C}\mathbf{B})^{-1}) \end{bmatrix}, \quad (7.7)$$

where the projector  $\mathbf{P}$  is defined by  $\mathbf{P} := \mathbf{B}(\mathbf{C}\mathbf{B})^{-1}\mathbf{C}$ . This implies that rounding errors proportional to  $\eta$  and  $\eta h^{-1}$  may be introduced in the variables  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. Let us now assume that the (approximate) variables have a moderate norm,

so that absolute and relative errors are qualitatively the same. Hence, we shall consider absolute "machine" rounding errors  $\delta$  rather than relative rounding errors  $\eta$  in our examples. The matrix  $\mathbf{J}_1^{-1}$  (cf. (7.7)) shows of course the ill-posedness, due to differentiation of the forcing function, of an index one DAE which is close to a DAE of higher index.

For the index two system

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{C}\mathbf{x} + \mathbf{q},\end{aligned}\tag{7.8}$$

where  $\mathbf{CB}$  is nonsingular, the resulting iteration matrix equals

$$\mathbf{J}_2 := \begin{bmatrix} \mathbf{I} - h\mathbf{A} & -h\mathbf{B} \\ \mathbf{C} & \mathbf{O} \end{bmatrix}.\tag{7.9}$$

Performing LU-decomposition we find a similar expression as in (7.3), but with  $\mathbf{D} = \mathbf{O}$ . When the matrix  $\mathbf{CB}$  is well-conditioned and bounded away from zero the inverse of  $\mathbf{J}_2$  is equal to  $\mathbf{J}_1^{-1}$  ( $\varepsilon \rightarrow 0$ ) (cf. (7.7)). This means that the numerical solution of the index one DAE (7.1) behaves like a solution of the associated index two DAE (7.8). Rounding errors proportional to  $\eta$  and  $\eta h^{-1}$  are introduced in the variables  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. This shows the ill-conditioning of a higher index DAE.

However, for  $\mathbf{CB} = \varepsilon\mathbf{I}$  ( $\varepsilon \rightarrow 0$ ),  $(h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B})^{-1}$  has to be computed. For small fixed  $h$  we find that

$$\begin{aligned}(h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B})^{-1} &= h^{-2}(\mathbf{CAB})^{-1}(\mathbf{I} - h\mathbf{CA}^2\mathbf{B}(\mathbf{CAB})^{-1}) \\ &+ \mathcal{O}(1) + \mathcal{O}(\varepsilon/h^3), \quad (\varepsilon \rightarrow 0),\end{aligned}\tag{7.10}$$

when the matrix  $(\mathbf{CAB})^{-1}$  is bounded. Hence, to first order

$$\mathbf{J}_2^{-1} \doteq \begin{bmatrix} -h^{-1}\mathbf{B}(\mathbf{CAB})^{-1}\mathbf{C} & h^{-1}\mathbf{B}(\mathbf{CAB})^{-1} \\ -h^{-2}(\mathbf{CAB})^{-1}\mathbf{C} & h^{-2}(\mathbf{CAB})^{-1} \end{bmatrix},\tag{7.11}$$

implying that rounding errors proportional to  $\eta h^{-1}$  and  $\eta h^{-2}$  are introduced in the variables  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. Note that when  $\mathbf{CB} = \mathbf{O}$  in the DAE (7.8), discretizing this DAE results in the same iteration matrix (cf. (7.11)). This implies that, also numerically, the solution of the index two DAE is close to the solution of the associated index three DAE, i.e.  $\mathbf{CB} = \mathbf{O}$  in (7.8) (see also Section 6, where the analogue has been shown for the exact solution). If the matrix  $\mathbf{D}$  is ill-conditioned and  $\mathbf{CB} = \mathbf{O}$  in the DAE (7.1) then the inverse of the resulting iteration matrix  $h\mathbf{J}_1$  also results in (7.11) and the index one DAE (7.1) is numerically close to an index three DAE.

Now consider  $\mathbf{D} + h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B}$ . For small  $h$  this matrix can be written as

$$\mathbf{D} + h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B} = \mathbf{D} + h\mathbf{CB} + h^2\mathbf{CAB} + h^3\mathbf{CA}^2\mathbf{B} + \dots$$

Suppose that both  $\mathbf{D}$  and its inverse have a norm of moderate size, then

$$(\mathbf{D} + h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B})^{-1} = \mathbf{D}^{-1} + \mathcal{O}(h).$$

On the other hand if the matrix  $\mathbf{D}$  is almost singular, say  $\mathbf{D} = \varepsilon\hat{\mathbf{D}}$ , for simplicity, with  $\|\hat{\mathbf{D}}^{-1}\| \approx 1$ , and  $\|(\mathbf{CB})^{-1}\| \leq M$  for a constant  $M$  of moderate size, then

$$(\mathbf{D} + h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B})^{-1} \approx h^{-1}(\mathbf{CB})^{-1}.$$

Then it is immediately obvious that

$$\|(\mathbf{D} + h\mathbf{C}(\mathbf{I} - h\mathbf{A})^{-1}\mathbf{B})^{-1}\| = \mathcal{O}(\min\{\varepsilon^{-1}, h^{-m+1}\}),$$

where  $m$  is the highest effective index of any of the variables. So, the conditioning of the iteration matrix is practically determined by the highest effective index of any variable.

**Example 7.1** Consider the index one DAE

$$\begin{aligned} \dot{x} &= -y + \sin t, \\ 0 &= x - \varepsilon y - \cos t, \end{aligned}$$

subject to the initial condition  $x(0) = 2$ . On the interval  $[0, 10^{-3}]$  we obtain from discretization with Euler backward Table 7.1, where the errors  $e_x := |x(t_n) - x_n|$ ,  $e_y := |y(t_n) - y_n|$  and drift  $:= |x_n - \varepsilon y_n - \cos t_n|$ , i.e. the deviation from the constraint, for  $nh = T = 10^{-3}$ . In order to show the influence of the rounding errors we introduce artificial absolute rounding errors  $\delta$ , say, with  $\delta \in [4 \cdot 10^{-5}, 6 \cdot 10^{-5}]$ , into the system:  $\varepsilon$  is taken equal to  $10^{-6}$ . It is obvious that errors proportional to  $\delta$  are introduced into the state variable  $x$ . Note that from the first three rows of Table 7.1 it appears that the numerical approximation of the algebraic variable  $y$  has errors proportional to  $\delta h^{-1}$ . In the last three rows, however, the resulting errors are proportional to  $\delta \varepsilon^{-1}$  as is explained above.  $\square$

**Example 7.2** Consider the DAE of index two, which is the DAE of Example 7.1 with  $\varepsilon = 0$

$$\begin{aligned} \dot{x} &= -y + \sin t, \\ 0 &= x - \cos t, \end{aligned}$$

$h$	$e_x$	$e_y$	drift
$10^{-4}$	$0.49 \cdot 10^{-4}$	$0.75 \cdot 10^{-1}$	$0.49 \cdot 10^{-4}$
$10^{-5}$	$0.55 \cdot 10^{-4}$	$0.13 \cdot 10^0$	$0.55 \cdot 10^{-4}$
$10^{-6}$	$0.54 \cdot 10^{-4}$	$0.14 \cdot 10^1$	$0.54 \cdot 10^{-4}$
$10^{-7}$	$0.55 \cdot 10^{-4}$	$0.30 \cdot 10^1$	$0.58 \cdot 10^{-4}$
$10^{-8}$	$0.55 \cdot 10^{-4}$	$0.10 \cdot 10^0$	$0.56 \cdot 10^{-4}$
$10^{-9}$	$0.55 \cdot 10^{-4}$	$0.35 \cdot 10^1$	$0.58 \cdot 10^{-4}$

Table 7.1: Results for the nearly index two DAE of Example 7.2 showing the  $h^{-1}$  effect in the algebraic variable  $y$ .

subject to the initial condition  $x(0) = 1$ ,  $y(0) = 0$ . This DAE has the solution  $x(t) = \cos t$ ,  $y(t) = 2 \sin t$ . Introducing artificial errors  $\delta \in [4 \cdot 10^{-5}, 6 \cdot 10^{-5}]$ , into the system shows the influence of the rounding errors. On the interval  $[0, 10^{-3}]$  we obtain from discretization with Euler backward the following table, where we use the following definitions for the errors:  $e_x := |x(t_n) - x_n|$ ,  $e_y := |y(t_n) - y_n|$  and drift  $:= |x_n - \cos t_n|$ , i.e. the deviation from the constraint, for  $nh = T = 10^{-3}$ . Table 7.2 shows that errors proportional to  $\delta$  appear in the state variable  $x$  whereas

$h$	$e_x$	$e_y$	drift
$10^{-4}$	$0.54 \cdot 10^{-4}$	$0.57 \cdot 10^{-2}$	$0.54 \cdot 10^{-4}$
$10^{-5}$	$0.60 \cdot 10^{-4}$	$0.62 \cdot 10^0$	$0.60 \cdot 10^{-4}$
$10^{-6}$	$0.51 \cdot 10^{-4}$	$0.44 \cdot 10^0$	$0.51 \cdot 10^{-4}$
$10^{-7}$	$0.59 \cdot 10^{-4}$	$0.72 \cdot 10^1$	$0.59 \cdot 10^{-4}$
$10^{-8}$	$0.54 \cdot 10^{-4}$	$0.20 \cdot 10^3$	$0.54 \cdot 10^{-4}$
$10^{-9}$	$0.58 \cdot 10^{-4}$	$0.11 \cdot 10^5$	$0.58 \cdot 10^{-4}$

Table 7.2: Results for the index two index one DAE of Example 7.1 showing the  $h^{-1}$  effect in the algebraic variable  $y$ .

errors proportional to  $\delta h^{-1}$  are introduced into the algebraic variable  $y$ .  $\square$

These examples show that effective index depends on thresholds, like accuracy (or timestep). The index 1 problem in Example 7.1 behaves like an index 2 problem until  $h \approx \varepsilon$ , after which it is just an index 1 problem with stability constant  $\varepsilon^{-1}$  ( $\varepsilon$  fixed).

## 8 Stabilization of DAE

There exists a number of methods to stabilize higher index problems. One is based on the introduction of additional Lagrange multipliers, cf. [5, 14]. For instance the DAE

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{y} + \mathbf{q},\end{aligned}\tag{8.1}$$

is replaced by

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{L}\boldsymbol{\mu} + \mathbf{p}, \\ \mathbf{0} &= \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{y} + \mathbf{q}, \\ \mathbf{0} &= \mathbf{C}\mathbf{A}\mathbf{x} + \mathbf{C}\mathbf{B}\mathbf{y} + \mathbf{C}\mathbf{p} + \dot{\mathbf{q}},\end{aligned}\tag{8.2}$$

where  $\mathbf{C}\mathbf{L}$  should be well conditioned. For the iteration matrix we obtain

$$\mathbf{J} = \begin{bmatrix} \mathbf{I} - h\mathbf{A} & -h\mathbf{B} & -h\mathbf{L} \\ \mathbf{C} & \mathbf{D} & \mathbf{0} \\ h\mathbf{C}\mathbf{A} & h\mathbf{C}\mathbf{B} + \mathbf{D} & \mathbf{0} \end{bmatrix},$$

and hence to first order in  $\varepsilon$

$$\mathbf{J}^{-1} \doteq \begin{bmatrix} \mathbf{I} - \mathbf{Q} & \mathbf{L}(\mathbf{C}\mathbf{L})^{-1} & \mathbf{B}(\mathbf{C}\mathbf{B})^{-1} - \mathbf{L}(\mathbf{C}\mathbf{L})^{-1} \\ (\mathbf{C}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}(\mathbf{Q} - \mathbf{I}) & -(\mathbf{C}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}\mathbf{L}(\mathbf{C}\mathbf{L})^{-1} & h^{-1}(\mathbf{C}\mathbf{B})^{-1} \\ -h^{-1}(\mathbf{C}\mathbf{L})^{-1}\mathbf{C} & h^{-1}(\mathbf{C}\mathbf{L})^{-1} & -h^{-1}(\mathbf{C}\mathbf{L})^{-1} \end{bmatrix}.$$

Note that the third block column of the above matrix is scaled by a factor  $h^{-1}$ , which arises from differentiating  $\mathbf{y}$  in (8.2); this factor will not cause any problems, since the matrix  $\mathbf{J}^{-1}$  operates on a right-hand side vector with a factor  $h$  appearing in the corresponding third row. Hence, in the thus stabilized system errors proportional to  $\eta$  will appear in the approximation of the variables  $\mathbf{x}$  and  $\mathbf{y}$ , whereas errors proportional to  $\eta h^{-1}$  will appear in the newly defined variable  $\boldsymbol{\mu}$ .

If the matrix  $\mathbf{D}$  is well-conditioned, one can show that

$$\mathbf{J}^{-1} \doteq \begin{bmatrix} \mathbf{I} - \mathbf{Q} & \mathbf{L}(\mathbf{C}\mathbf{L})^{-1} & -\mathbf{L}(\mathbf{C}\mathbf{L})^{-1} \\ h\mathbf{D}^{-1}\mathbf{C}\mathbf{A}(\mathbf{Q} - \mathbf{I}) & -h\mathbf{D}^{-1}\mathbf{C}\mathbf{A}\mathbf{L}(\mathbf{C}\mathbf{L})^{-1} & \mathbf{D}^{-1} \\ -h^{-1}(\mathbf{C}\mathbf{L})^{-1}\mathbf{C} & h^{-1}(\mathbf{C}\mathbf{L})^{-1} & -h^{-1}(\mathbf{C}\mathbf{L})^{-1} \end{bmatrix},$$

and the conditioning of the DAE will be the same as before for the index one variables, whereas the index two variable  $\boldsymbol{\mu}$  has a conditioning constant  $\mathcal{O}(h^{-1})$ . The second row of  $\mathbf{J}^{-1}$  shows that errors proportional to  $\eta h$  are introduced in well-conditioned algebraic variables. Combining these two results shows that, whenever the matrix

$\mathbf{D}$  is ill-conditioned, one may stabilize without making any factorizations of  $\mathbf{D}$  in order to determine the effective index of any of the variables. Hence, we have constructed a method to cheaply improve the conditioning of an index one DAE, which contains variables that are effectively of higher index.

**Example 8.1** Consider an index one DAE (8.1), with

$$\mathbf{A} = \begin{bmatrix} -2 & 1 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{B} = [\mathbf{B}_1 \quad \mathbf{B}_2] = \begin{bmatrix} 1 & 1 \\ 2 & 0 \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} \mathbf{C}_1 \\ \mathbf{C}_2 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix},$$

$$\text{and } \mathbf{D} = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \end{bmatrix},$$

where  $p_1(t) = q_1(t) = t/2 + \sin t$ ,  $p_2(t) = t + 2 \cos t$ , and  $q_2 = \cos t - \sin t$ .

The reduced DAE (i.e.  $\varepsilon = 0$ ) has the solution

$$x_1(t) = -q_2(t), \quad x_2(t) = x_2(0) \exp(-t), \quad (8.3)$$

$$y_1(t) = x_2(0) \exp(-t) - t/2 - \cos t, \quad \text{and } y_2(t) = 2 \sin t. \quad (8.4)$$

For  $\varepsilon$  approaching to zero, the algebraic variable  $y_2$  behaves effectively like an index two variable, whereas  $y_1$  is an index one variable. According to Section 7 rounding errors proportional to  $h^{-1}$  are introduced in the numerical approximation of  $\mathbf{y}$  due to the ill-conditioning of the iteration matrix. Let us choose initial values  $x_1 = 2$  and  $x_2 = 0$ . Using the stabilized form (8.2) we obtain from discretization with Euler backward on the interval  $I = [0, 10^{-3}]$  the result of Table 8.3. Here the errors are  $e_{x_i} := |x_i(t_n) - x_{i,n}|$ ,  $e_{y_i} := |y_i(t_n) - y_{i,n}|$ ,  $i = (1, 2)$ , and the drift is defined as  $\text{drift} := \|(\mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{y} + \mathbf{q})_n\|_\infty$ , i.e. the deviation from the constraint, for  $nh = T = 10^{-3}$ . Introducing artificial errors  $\delta \in [4 \cdot 10^{-5}, 6 \cdot 10^{-5}]$  in the system shows the influence of the rounding errors, when  $\varepsilon = 10^{-6}$ .  $\square$

## 9 (Nearly) singular pencils

The problem of the effective index of a DAE or a variable is closely related to detecting a near singularity in a coefficient matrix. Numerically this should be done by determining singular values (cf. [7]). If a singular value is small we may replace it by zero for theoretical purposes and differentiate the relevant part of the constraint as in a higher index case. If an entire coefficient matrix is zero we thus have

- (i). index 2, if  $\mathbf{D} = \mathbf{O}$ ,
- (ii). index 3, if  $\mathbf{CB} = \mathbf{O}$ .

$h$	$e_{x_1}$	$e_{x_2}$	$e_{y_1}$	$e_{y_2}$	drift
$10^{-4}$	$0.55 \cdot 10^{-4}$	$0.99 \cdot 10^{-4}$	$0.39 \cdot 10^{-8}$	$0.12 \cdot 10^{-3}$	$0.55 \cdot 10^{-4}$
$10^{-5}$	$0.60 \cdot 10^{-4}$	$0.11 \cdot 10^{-3}$	$0.50 \cdot 10^{-7}$	$0.13 \cdot 10^{-3}$	$0.60 \cdot 10^{-4}$
$10^{-6}$	$0.49 \cdot 10^{-4}$	$0.90 \cdot 10^{-4}$	$0.54 \cdot 10^{-7}$	$0.12 \cdot 10^{-3}$	$0.49 \cdot 10^{-4}$
$10^{-7}$	$0.56 \cdot 10^{-4}$	$0.11 \cdot 10^{-3}$	$0.55 \cdot 10^{-7}$	$0.11 \cdot 10^{-3}$	$0.56 \cdot 10^{-4}$
$10^{-8}$	$0.58 \cdot 10^{-4}$	$0.11 \cdot 10^{-3}$	$0.55 \cdot 10^{-7}$	$0.11 \cdot 10^{-3}$	$0.58 \cdot 10^{-4}$
$10^{-9}$	$0.58 \cdot 10^{-4}$	$0.12 \cdot 10^{-3}$	$0.55 \cdot 10^{-7}$	$0.11 \cdot 10^{-3}$	$0.59 \cdot 10^{-4}$

Table 8.3: Results for the nearly index two index one DAE of Example 8.1 showing that the stabilization method can circumvent the  $h^{-1}$  effect in the algebraic variable  $y_2$ .

Differentiating once more we see that we should investigate, at least for constant coefficients, the singularity of

$$\mathbf{CAB}. \quad (9.1)$$

If the latter matrix is singular we clearly have index 4. In general it is the singularity of the matrix

$$\mathbf{CA}^i\mathbf{B}, \quad (9.2)$$

which may decide for at least index  $i + 3$ . We remark that there exists a program GELDA [11] determining the actual index numerically. The extreme situation where  $i \rightarrow \infty$  is now shown to be equivalent to the matrix pencil being singular. Defining

$$\hat{\mathbf{E}} := \begin{bmatrix} \mathbf{I}_n & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}; \quad \hat{\mathbf{A}} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}, \quad (9.3)$$

we recall the notion of the *matrix pencil*  $(\hat{\mathbf{E}}, \hat{\mathbf{A}})$  (cf. [4, 7]). The matrix pencil  $(\hat{\mathbf{E}}, \hat{\mathbf{A}})$  is called singular if  $\hat{\mathbf{A}} - \lambda\hat{\mathbf{E}}$  is *singular* for all values of  $\lambda$ ; otherwise it is called *regular*.

We first prove a simple lemma.

**Lemma 9.1** *If the matrix pencil  $(\hat{\mathbf{E}}, \hat{\mathbf{A}})$  is singular, then  $\mathbf{D}$  is singular.*

*Proof.* Apparently for any  $\lambda$  there exists a vector  $\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}$  such that

$$\begin{aligned} \mathbf{Ax} + \mathbf{By} &= \lambda\mathbf{x}, \\ \mathbf{Cx} + \mathbf{Dy} &= \mathbf{0}. \end{aligned}$$



If  $\mathbf{D}$  were nonsingular we could write

$$[\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}]\mathbf{x} = \lambda\mathbf{x}.$$

This is a "regular" eigenvalue problem, which makes sense for a finite number ( $\leq n$ ) of values  $\lambda$  only. For any other value of  $\lambda$  we must have  $\mathbf{x} = \mathbf{0}$ , whence  $\mathbf{y} = \mathbf{0}$ ; this is, however, excluded, thus leading to a contradiction. Hence  $\mathbf{D}$  must be singular.  $\square$

We now come to the main result.

**Theorem 9.2** *The matrix pencil  $(\hat{\mathbf{E}}, \hat{\mathbf{A}})$  is singular iff for some  $\mathbf{z} \in \mathbb{R}^m$ ,  $\mathbf{z}^T\mathbf{D} = \mathbf{0}$  and  $\mathbf{z}^T\mathbf{C}\mathbf{A}^i\mathbf{B} = \mathbf{0}$  for all  $i$ .*

*Proof.* According to Lemma 9.1 it is not restrictive to assume that  $\mathbf{C}$  and  $\mathbf{D}$  have the following form:

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_1 \\ \mathbf{C}_2 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} \mathbf{O} \\ \mathbf{D}_2 \end{bmatrix}, \quad \mathbf{C}_1 \in \mathbb{R}^{p \times n}, \quad \mathbf{C}_2 \in \mathbb{R}^{(m-p) \times n}, \quad \mathbf{D} \in \mathbb{R}^{(m-p) \times m},$$

where  $\mathbf{D}_2$  has full row rank, i.e.  $p (\geq 1)$ , is the geometric multiplicity of eigenvectors belonging to the eigenvalue 0. (Note that such a form can always be obtained by a similarity transformation  $\begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{T}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{T} \end{bmatrix}$  for some  $\mathbf{T}$ .) Actually we may even assume that  $\mathbf{D}_2 = [\mathbf{O}|\mathbf{D}_{22}]$ , where  $\mathbf{D}_{22} \in \mathbb{R}^{(m-p) \times (m-p)}$  is nonsingular. The special form of  $\mathbf{D}$  implies a projection,  $\mathbf{P}_1$  say, with  $\mathbf{P}_1\mathbf{D} = \mathbf{D}$ . Then  $\mathbf{Q}_1 := \mathbf{I} - \mathbf{P}_1 = \begin{bmatrix} \mathbf{I}_p & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \in \mathbb{R}^{m^2}$ . By assumption  $\mathbf{Q}_1$  is not trivial ( $p \geq 1$ ).

We now show the "only if". Given a  $\lambda$ , there exists some nontrivial vector  $\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \neq \mathbf{0}$  with  $[\hat{\mathbf{A}} - \lambda\hat{\mathbf{E}}]\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}$ . Hence with  $\hat{\mathbf{T}}_1 := \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{Q}_1\mathbf{C} & \mathbf{I} \end{bmatrix}$ , we obtain

$$\hat{\mathbf{T}}_1 [\hat{\mathbf{A}} - \lambda\hat{\mathbf{E}}] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \lambda\mathbf{I} & \mathbf{B} \\ \mathbf{Q}_1\mathbf{C}\mathbf{A} + \mathbf{C} & \mathbf{Q}_1\mathbf{C}\mathbf{B} + \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \lambda\mathbf{I} & \mathbf{B} \\ \mathbf{C}_1\mathbf{A} & \mathbf{C}_1\mathbf{B} \\ \mathbf{C}_2 & \mathbf{D}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}.$$

Applying Lemma 9.1 to the latter system we conclude that  $\begin{bmatrix} \mathbf{C}_1\mathbf{B} \\ \mathbf{D}_2 \end{bmatrix}$  must be singular.

Hence for some projection  $\mathbf{P}_2$  we have  $\mathbf{P}_2 \begin{bmatrix} \mathbf{C}_1\mathbf{B} \\ \mathbf{D}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{C}_1\mathbf{B} \\ \mathbf{D}_2 \end{bmatrix}$ . Due to the fact that  $\mathbf{D}_2$  has full rank, we realize that  $\mathbf{P}_2 := \begin{bmatrix} \tilde{\mathbf{P}}_2 & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix}$ , where  $\tilde{\mathbf{P}}_2 \in \mathbb{R}^{p^2}$  (and of rank  $\geq 1$ ).

Define  $\tilde{\mathbf{Q}}_2 := \mathbf{I} - \tilde{\mathbf{P}}_2$  and  $\hat{\mathbf{T}}_2 := \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \tilde{\mathbf{Q}}_2\mathbf{C}_1 & \mathbf{I} \end{bmatrix}$ . Then

$$\hat{\mathbf{T}}_2\hat{\mathbf{T}}_1 [\hat{\mathbf{A}} - \lambda\hat{\mathbf{E}}] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \lambda\hat{\mathbf{I}} & \mathbf{B} \\ \tilde{\mathbf{Q}}_2\mathbf{C}_1\mathbf{A}^2 & \tilde{\mathbf{Q}}_2\mathbf{C}_1\mathbf{A}\mathbf{B} \\ \mathbf{C}_2 & \mathbf{D}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}.$$

In general we can thus find a sequence of nontrivial projection matrices  $\tilde{\mathbf{Q}}_i$ , with  $\text{rank}(\tilde{\mathbf{Q}}_i) =: p_i \geq 1$ , such that  $(\tilde{\mathbf{Q}}_i \cdots \tilde{\mathbf{Q}}_2 \mathbf{C}_1 \mathbf{A}^i \mathbf{B})$  has  $\text{rank} \geq 1$ . In particular this implies that for some  $\mathbf{z} \in \mathbb{R}^m$  and for all  $i$   $\mathbf{z}^T \mathbf{C}_1 \mathbf{A}^i \mathbf{B} = \mathbf{0}$ . From the construction it trivially follows that  $\mathbf{z}^T \mathbf{D} = \mathbf{0}$ .

Next we prove the "if". So let  $\mathbf{D}$  be as above and  $\mathbf{z}^T \mathbf{C}_1 \mathbf{A}^i \mathbf{B} = \mathbf{0}$  for some  $\mathbf{z} \in \mathbb{R}^m$  and all  $i$ . First we consider the case that  $\lambda \notin S(\mathbf{A})$ ,  $\lambda \neq 0$  ( $S(\mathbf{A})$  is the spectrum of  $\mathbf{A}$ ). Let  $\tilde{\mathbf{A}} := \frac{1}{\lambda} \mathbf{A}$ . Then

$$\tilde{\mathbf{A}}^i = (\tilde{\mathbf{A}} - \mathbf{I})[\tilde{\mathbf{A}}^{i-1} + \cdots + \mathbf{I}] + \mathbf{I}.$$

Since  $\mathbf{z}^T \mathbf{C}_1 \mathbf{A}^j \mathbf{B} = \mathbf{0}$ , we obtain

$$\mathbf{z}^T \mathbf{C}_1 (\tilde{\mathbf{A}} - \mathbf{I})^{-1} \tilde{\mathbf{A}}^i \mathbf{B} = \mathbf{0} + \mathbf{z}^T \mathbf{C}_1 (\tilde{\mathbf{A}} - \mathbf{I})^{-1} \mathbf{B}.$$

Defining the characteristic polynomial of  $\tilde{\mathbf{A}}$  as  $\psi_{\tilde{\mathbf{A}}}(\mu) := \sum_{i=0}^n \alpha^i \mu^i$ , we thus deduce

$$\mathbf{0} = \mathbf{z}^T \mathbf{C}_1 (\tilde{\mathbf{A}} - \mathbf{I})^{-1} \psi_{\tilde{\mathbf{A}}}(\tilde{\mathbf{A}}) \mathbf{B} = \psi_{\tilde{\mathbf{A}}}(1) \mathbf{z}^T \mathbf{C}_1 (\tilde{\mathbf{A}} - \mathbf{I})^{-1} \mathbf{B}$$

Since  $1 \notin S(\tilde{\mathbf{A}})$ , we thus find

$$\mathbf{z}^T \mathbf{C}_1 (\tilde{\mathbf{A}} - \mathbf{I})^{-1} \mathbf{B} = \mathbf{0},$$

whence

$$\mathbf{z}^T \mathbf{C}_1 (\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{B} = \mathbf{0}$$

This then can be used to see that the matrix

$$[-\mathbf{C}(\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{B} + \mathbf{D}]$$

is singular (premultiply by  $(\mathbf{z}^T, 0, \dots, 0)$ ). So there exists a nontrivial vector  $\mathbf{y}$  such that

$$[-\mathbf{C}(\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{B} + \mathbf{D}] \mathbf{y} = \mathbf{0}.$$

Upon substituting  $\mathbf{x} := -(\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{B} \mathbf{y}$ , we see that

$$\begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}, \quad \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \neq \mathbf{0}.$$

If  $\lambda \notin S(\mathbf{A})$  but  $\lambda = 0$  we can obtain more directly that  $\mathbf{z}^T \mathbf{C}_1 \mathbf{A}^{-1} \mathbf{B} = \mathbf{0}$  and a similar proof follows.

Now let  $\lambda \in S(\mathbf{A})$ ,  $\lambda \neq 0$ . In this case is not restrictive to assume that  $\mathbf{A} = \begin{bmatrix} \lambda \mathbf{I}_s & \mathbf{O} \\ \mathbf{O} & \mathbf{K} \end{bmatrix}$ , where  $s$  is the geometric multiplicity of  $\lambda$ . Partition  $\mathbf{C}_1 = [\mathbf{C}_{11} | \mathbf{C}_{12}]$ , where  $\mathbf{C}_{11}$  contains  $s$  columns. Then for  $\tilde{\mathbf{A}} := \frac{1}{\lambda} \mathbf{A}$  and  $\tilde{\mathbf{K}} := \frac{1}{\lambda} \mathbf{K}$  we have

$$\begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & (\tilde{\mathbf{K}} - \mathbf{I})^{-1} \tilde{\mathbf{K}}^i \end{bmatrix} = \begin{bmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \left\{ \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \tilde{\mathbf{K}}^{i-1} \end{bmatrix} + \cdots + \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \right\} + \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & (\tilde{\mathbf{K}} - \mathbf{I})^{-1} \end{bmatrix}$$

Hence with  $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix}$  ( $\mathbf{B}_1$  having  $s$  rows)

$$\mathbf{z}^T (\mathbf{C}_{11} \mathbf{B}_1 + \mathbf{C}_{12} (\tilde{\mathbf{K}} - \mathbf{I})^{-1} \tilde{\mathbf{K}}^i \mathbf{B}_2) = \mathbf{0} + \mathbf{z}^T (\mathbf{C}_{11} \mathbf{B}_1 + \mathbf{C}_{12} (\tilde{\mathbf{K}} - \mathbf{I})^{-1} \mathbf{B}_2)$$

Using the characteristic polynomial  $\psi_{\tilde{\mathbf{K}}}$  with  $\psi_{\tilde{\mathbf{K}}}(1) \neq 0$  and  $\psi_{\tilde{\mathbf{K}}}(\tilde{\mathbf{K}}) = \mathbf{0}$  we obtain

$$\mathbf{z}^T \mathbf{C}_{12} (\tilde{\mathbf{K}} - \mathbf{I})^{-1} \mathbf{B}_2 = \mathbf{0} = \mathbf{z}^T \mathbf{C}_{12} (\tilde{\mathbf{K}} - \lambda \mathbf{I})^{-1} \mathbf{B}_2.$$

As a consequence

$$\mathbf{z}^T \mathbf{C}_{11} \mathbf{B}_1 = \mathbf{0}.$$

If  $\mathbf{C}_1$  denotes the first  $s$  columns of  $\mathbf{C}$ , then we conclude therefore that either  $\mathbf{C}_1$  or  $\mathbf{B}_1$  is a singular square matrix if  $s = m$ . We may rewrite the matrix  $(\hat{\mathbf{A}} - \lambda \hat{\mathbf{E}})$  as

$$\begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} & \mathbf{B} \\ \mathbf{C} & \mathbf{O} \end{bmatrix} = \begin{bmatrix} \mathbf{O} & \mathbf{O} & \mathbf{B}_1 \\ \mathbf{O} & \mathbf{K} - \lambda \mathbf{I} & \mathbf{B}_2 \\ \mathbf{C}_1 & \mathbf{C}_2 & \mathbf{D} \end{bmatrix} \quad \text{and partition } \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{y} \end{bmatrix}.$$

Note that  $\mathbf{B}_1 \in \mathbb{R}^{s \times m}$ ,  $\mathbf{C}_1 \in \mathbb{R}^{m \times s}$ . If  $s > m$ , we can always find a nontrivial vector  $\mathbf{x}_1 \in \mathbb{R}^s$  with  $\mathbf{C}_1 \mathbf{x}_1 = \mathbf{0}$ ,  $\mathbf{x}_1 \neq \mathbf{0}$ . Hence

$$[\hat{\mathbf{A}} - \lambda \hat{\mathbf{E}}] \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}.$$

If  $s = m$  and  $\mathbf{C}_1$  is singular we can proceed as before. If  $s = m$  and only  $\mathbf{B}_1$  is singular as well as for the case  $s < m$  we can consider row vectors and left multiplication to show singularity.

Finally if  $\lambda \in S(\mathbf{A})$ ,  $\lambda = 0$  we can shorten the beginning of the proof as we did for the case  $\lambda \notin S(\mathbf{A})$ .  $\square$

We remark that it follows from Theorem 9.2 that  $\mathbf{C} \mathbf{A}^i \mathbf{B}$  is singular for all  $i$ , if the pencil  $(\hat{\mathbf{E}}, \hat{\mathbf{A}})$  is singular. The result in this theorem is slightly stronger, however. To show the mechanism of  $\mathbf{A}$  producing a non singular lower right block in the updating process eventually, consider the next example.

**Example 9.3** Consider a DAE where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \mathbf{C} = [0 \quad 1] \quad \text{and } \mathbf{D} = \mathbf{O}. \quad (9.4)$$

Then  $\mathbf{CB} = \mathbf{O}$ , but  $\mathbf{CAB} = 1$ , whence the index is 3. The fact that  $\mathbf{CA}^i\mathbf{B} = \mathbf{O}$  for all even values of  $i$ , is of no importance in this example.  $\square$

Next consider a nearly singular problem.

**Example 9.4** Consider a DAE where

$$\mathbf{A} = \begin{bmatrix} -1 & -1 \\ 0 & 0 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \mathbf{C} = [0 \quad 1] \quad \text{and } D = \varepsilon. \quad (9.5)$$

The solution reads.

$$\begin{aligned} \mathbf{x}(t) &= \begin{bmatrix} e^{-t} & (1 + \varepsilon^{-1})(e^{-t} - 1) \\ 0 & 1 \end{bmatrix} \mathbf{x}(0) - \varepsilon^{-1} \int_0^t \begin{bmatrix} e^{-(t-s)} \\ 0 \end{bmatrix} q(s) ds \\ &\quad + \int_0^t \begin{bmatrix} e^{-(t-s)} & (1 + \varepsilon^{-1})(e^{-(t-s)} - 1) \\ 0 & 1 \end{bmatrix} \mathbf{p}(s) ds, \\ y(t) &= -\varepsilon^{-1} x_2(0) - \varepsilon^{-1} \int_0^t \mathbf{p}(s) ds - \varepsilon^{-1} q(t). \end{aligned}$$

Since there is no fast decaying integrating function under the integral sign, practical integration does not work to "remove" the stability constants which are of  $\mathcal{O}(\varepsilon^{-1})$ . It is simple to see that  $\det(\hat{\mathbf{A}} - \lambda \hat{\mathbf{E}}) = -\varepsilon(\lambda + 1)^2$ . Hence, we have a singular pencil for  $\varepsilon = 0$ . For  $\varepsilon = 0$  we clearly also have  $\mathbf{CA}^i\mathbf{B} = \mathbf{O} \forall i \in \mathbb{N}$ . It is obvious that the stability constants are of order  $\mathcal{O}(\varepsilon^{-1})$  for  $\mathbf{x}$  as well as for  $y$ , however often we apply partial integration.  $\square$

## References

- [1] U. ASCHER AND L. R. PETZOLD. Projected implicit Runge-Kutta for differential algebraic equations. *SIAM J. Numer. Anal.* 28 (1991), 1097–1120.
- [2] U. M. ASCHER, R. M. M. MATTHEIJ, AND R. D. RUSSELL. *Numerical Solution of Boundary Value Problems for Ordinary Differential Equations*. SIAM, Philadelphia, 1995.
- [3] K. E. BRENNAN, S. L. CAMPBELL, AND L. R. PETZOLD. *Numerical Solution of Initial-Value Problems in Differential Algebraic Equations*. North-Holland, New York, 1989.
- [4] F. R. GANTMACHER. *The Theory of Matrices*, vol. 2. Chelsea, 1964.
- [5] C. W. GEAR, G. K. GUPTA, AND B. J. LEIMKUEHLER. Automatic integration of Euler-Lagrange equations with constraints. *J. Comp. Appl. Math.* 12&13 (1985), 77–90.
- [6] C. W. GEAR AND L. R. PETZOLD. Ode methods for the solution of differential algebraic systems. *SIAM J. Numer. Anal.* 21 (1984), 367–384.
- [7] G. H. GOLUB AND C. F. V. LOAN. *Matrix Computations*. John Hopkins University Press, Baltimore, 1990.
- [8] E. GRIEPENTROG AND R. MÄRZ. *Differential Algebraic Equations and Their Numerical Treatment*, vol. 88. Teubner-Texte zur Math., Leipzig, 1986.
- [9] E. HAIRER AND G. WANNER. *Solving Ordinary Differential Equations II, Stiff and Differential Algebraic Problems*. Springer-Verlag, Berlin, 1991.
- [10] L. V. KALACHEV AND R. E. O’MALLEY. The regularization of linear differential-algebraic equations. *SIAM J. Math. Anal.* 27, 3 (1996), 258–273.
- [11] P. KUNKEL, V. MEHRMANN, W. RATH, AND J. WEICKERT. A new software package for linear differential-algebraic equations. dedicated to c. william gear on the occasion of his 60th birthday. *SIAM J. Sci. Comput.* 18, 1 (1997), 115–138.
- [12] M. LENTINI AND R. MÄRZ. The conditioning of boundary value problems in transferable differential-algebraic equations. *SIAM J. Numer. Anal.* 27, 4 (1990), 1001–1015.
- [13] G. SÖDERLIND. Remarks on the stability of high-index DAEs with respect to parametric perturbations. *Computing* 49 (1992), 303–314.

- [14] P. M. E. J. WIJCKMANS. *Conditioning of Differential Algebraic Equations and Numerical Solution of Multibody Dynamics*. PhD thesis, Eindhoven University of Technology, Eindhoven, 1996.

PREVIOUS PUBLICATIONS IN THIS SERIES:

Number	Author(s)	Title	Month
97-11	E.F. Kaasschieter J.D. van der Werff ten Bosch G.J. Mulder	A Numerical Fractional Flow Model for Air Sparging	September '97
97-12	He Yinnian	On the Convergence of Optimum Nonlinear Galerkin Method	September '97
97-13	He Yinnian R.M.M. Mattheij	Optimum Mixed Finite Element Nonlinear Galerkin Method for the Navier-Stokes Equations; I: Error Estimates for Spatial Discretization	September '97
97-14	He Yinnian	Optimum Mixed Finite Element Nonlinear Galerkin Method for the Navier-Stokes Equations; II: Stabil- ity Analysis for Time Discretization	September '97
97-15	He Yinnian	Optimum Mixed Finite Element Nonlinear Galerkin Method for the Navier-Stokes Equations; III: Convergence Analysis for Time Discretization	September '97
97-16	A.F.M. ter Elst C.M.P.A. Smulders	Reduced heat kernels on homoge- neous spaces	September '97
97-17	H.J.C. Huijberts	Minimal order linear model match- ing for nonlinear control systems	September '97
97-18	L.C.G.J.M. Habets	System equivalence for AR-systems over rings	November '97
97-19	R.M.M. Mattheij P.M.E.J. Wijckmans	Conditioning of Two Deck Differen- tial Algebraic Equations	November '97

