

Four Lectures on Differential-Algebraic Equations

Steffen Schulz

Humboldt Universität zu Berlin

June 13, 2003

Abstract

Differential-algebraic equations (DAEs) arise in a variety of applications. Therefore their analysis and numerical treatment plays an important role in modern mathematics. This paper gives an introduction to the topic of DAEs. Examples of DAEs are considered showing their importance for practical problems. Several well known index concepts are introduced. In the context of the tractability index existence and uniqueness of solutions for low index linear DAEs is proved. Numerical methods applied to these equations are studied.

Mathematics Subject Classification: 34A09, 65L80

Keywords: differential-algebraic equations, numerical integration methods

Introduction

In this report we consider implicit differential equations

$$f(x'(t), x(t), t) = 0 \tag{1}$$

on an interval $\mathcal{J} \subset \mathbb{R}$. If $\frac{\partial f}{\partial x'}$ is nonsingular, then it is possible to formally solve (1) for x' in order to obtain an ordinary differential equation. However, if $\frac{\partial f}{\partial x'}$ is singular, this is no longer possible and the solution x has to satisfy certain algebraic constraints. Thus equations (1) where $\frac{\partial f}{\partial x'}$ is singular are referred to as *differential-algebraic equations* or DAEs.

These notes aim at giving an introduction to differential-algebraic equations and are based on four lectures given by the author during his stay at the University of Auckland in 2003.

The first section deals with examples of DAEs. Here problems from different kinds of applications are considered in order to stress the importance of DAEs when modeling practical problems.

In the second section each DAE is assigned a number, the index, to measure its complexity concerning both theoretical and numerical treatment. Several index notions are introduced, each of them stressing different aspects of the DAE considered. Special emphasis is given to the tractability index for linear DAEs.

The definition of the tractability index in the second section gives rise to a detailed analysis concerning existence and uniqueness of solutions. The main tool is a procedure to decouple the DAE into its dynamical and algebraic part. In section three this analysis is carried out for linear DAEs with low index as it was established by März [25].

The results obtained, especially the decoupling procedure, are used in the fourth section to study the behaviour of numerical methods when applied to linear DAEs. The material presented in this section is mainly taken from [18].

1 Examples of differential-algebraic equations

Modelling with differential-algebraic equations plays a vital role, among others, for constrained mechanical systems, electrical circuits and chemical reaction kinetics. In this section we will give examples of how DAEs are obtained in these fields. We will point out important characteristics of differential-algebraic equations that distinguish them from ordinary differential equations.

More information about differential-algebraic equations can be found in [2, 15] but also in [32].

1.1 Constrained mechanical systems

Consider the mathematical pendulum in figure 1.1. Let m be the pendulum's mass which is attached to a rod of length l [15]. In order to describe the pendulum in Cartesian coordinates we write down the potential energy

$$U(x, y) = mgh = mgl - mgy \quad (1.1)$$

where $(x(t), y(t))$ is the position of the moving mass at time t . The earth's acceleration of gravity is given by g , the pendulum's height is h . If we denote derivatives of x and y by \dot{x} and \dot{y} respectively, the kinetic energy is given by

$$T(\dot{x}, \dot{y}) = \frac{1}{2}m(\dot{x}^2 + \dot{y}^2). \quad (1.2)$$

The term $\dot{x}^2 + \dot{y}^2$ describes the pendulum's velocity. The constraint is found to be

$$0 = g(x, y) = x^2 + y^2 - l^2. \quad (1.3)$$

(1.1)-(1.3) are used to form the Lagrange function

$$L(q, \dot{q}) = T(\dot{x}, \dot{y}) - U(x, y) - \lambda g(x, y).$$

Here q denotes the vector $q = (x, y, \lambda)$. Note that λ serves as a Lagrange multiplier. The equations of motion are now given by Euler's equations

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_k} \right) - \frac{\partial L}{\partial q_k} = 0, \quad k = 1, 2, 3.$$

We arrive at the system

$$\begin{aligned} m\ddot{x} + 2\lambda x &= 0, \\ m\ddot{y} - mg + 2\lambda y &= 0, \\ g(x, y) &= 0. \end{aligned} \quad (1.4)$$

By introducing additional variables $u = \dot{x}$ and $v = \dot{y}$ we see that (1.4) is indeed of the form (1).

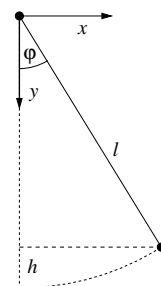


Figure 1.1: The mathematical pendulum

When solving (1.4) as an initial value problem, we observe that each initial value $(x(t_0), y(t_0)) = (x_0, y_0)$ has to satisfy the constraint (1.3) (consistent initialization). No initial condition can be posed for λ , as λ is determined implicitly by (1.4).

Of course the pendulum can be modeled by the second order ordinary differential equation

$$\ddot{\varphi} = -\frac{g}{l} \sin \varphi$$

when the angle φ is used as the dependent variable. However for practical problems a formulation in terms of a system of ordinary differential equations is often not that obvious, if not impossible.

1.2 Electrical circuits

Modern simulation of electrical networks is based on modelling techniques that allow an automatic generation of the model equations. One of the techniques most widely used is the modified nodal analysis (MNA) [7, 8].

1.2.1 A simple example

To see how the modified nodal analysis works, consider the simple circuit in figure 1.2 taken from [39]. It consists of a voltage source $v_V = v(t)$, a resistor with conductance G and a capacitor with capacitance $C > 0$. The layout of the circuit can be described by

$$A_a = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \\ 1 & 0 & -1 \end{pmatrix},$$

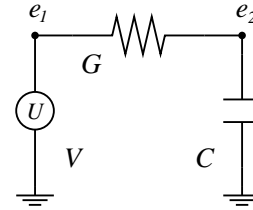


Figure 1.2: A simple circuit

where the columns of A_a correspond to the voltage, resistive and capacitive branches respectively. The rows represent the network's nodes, so that -1 and 1 indicate the nodes that are connected by each branch under consideration. Thus A_a assigns a polarity to each branch.

By construction the rows of A_a are linearly dependent. However, after deleting one row the remaining rows describe a set of linearly independent equations, The node corresponding to the deleted row will be denoted as the ground node. The matrix

$$A = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \end{pmatrix}$$

is called the incidence matrix. It is now possible to formulate basic physical laws in terms of the incidence matrix A [20]. Denote with i and v the vector of branch currents and voltage drops respectively and introduce the vector e of node potentials. For each node the node potential is it's voltage with respect to the ground node.

- Kirchhoff's Current Law (KCL): For each node the sum of all currents is zero. $\} \Rightarrow Ai = 0$
- Kirchhoff's Voltage Law (KVL): For each loop the sum of all voltages is zero. $\} \Rightarrow v = A^T e$

For the circuit in figure 1.2 KCL and KVL read

$$-i_V + i_G = 0, \quad -i_G + i_C = 0 \quad (1.5a)$$

and

$$v_V = -e_1, \quad v_G = e_1 - e_2, \quad v_C = e_2 \quad (1.5b)$$

respectively. If we assume ideal linear devices the equations modelling the resistor and the capacitor are

$$i_G = Gv_G, \quad i_C = C \frac{dv_C}{dt}. \quad (1.5c)$$

Finally we have

$$v_V = v(t) \quad (1.5d)$$

for the independent source which is thought of as the input signal driving the system. The system (1.5) is called the *sparse tableau*. The equations of the modified nodal analysis are obtained from the sparse tableau by expressing voltages in terms of node potential via (1.5b) and currents, where possible, by device equations (1.5c):

$$\left. \begin{array}{l} -i_V + G(e_1 - e_2) = 0 \\ -G(e_1 - e_2) + C \frac{de_2}{dt} = 0 \\ -e_1 = v \end{array} \right\} \\ \Leftrightarrow \begin{pmatrix} 0 \\ C \\ 0 \end{pmatrix} \left((0 \ 1 \ 0) \begin{pmatrix} e_1 \\ e_2 \\ i_V \end{pmatrix} \right)' + \begin{pmatrix} G & -G & -1 \\ -G & G & 0 \\ -1 & 0 & 0 \end{pmatrix} \begin{pmatrix} e_1 \\ e_2 \\ i_V \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ v \end{pmatrix} \quad (1.6)$$

The MNA equations reveal typical properties of DAEs:

- (i) Only certain parts of $x = (e_1, e_2, i_V)^T$ need to be differentiable. It is sufficient if e_1 and i_V are continuous.
- (ii) Any initial condition $x(t_0) = x_0$ needs to be consistent, i.e. there is a solution passing through x_0 . Here this means that we can pose an initial condition for e_2 or i_V only.

For (1.6) it is sufficient to solve the ordinary differential equation

$$e_2'(t) = -C^{-1}G(v(t) + e_2(t)).$$

$e_2(t)$ can be thought of as the output signal. The remaining components of the solution are uniquely determined as

$$e_1(t) = -v(t), \quad i_V(t) = G(e_1(t) - e_2(t)).$$

Another important feature that distinguishes DAEs from ordinary differential equations is that the solution process often involves differentiation rather than integration. This is illustrated in the next example.

1.2.2 Another simple example

If we replace the independent voltage in figure 1.2 source by a current source $i_I = i(t)$ and the capacitor by an inductor with inductance L , we arrive at the circuit in figure 1.3. The sparse tableau now reads

$$-i_I + i_G = 0, \quad -i_G + i_L = 0, \quad (1.7a)$$

$$v_I = -e_1, \quad v_G = e_1 - e_2, \quad v_L = e_2, \quad (1.7b)$$

$$i_G = Gv_G, \quad v_L = L \frac{di_L}{dt}, \quad (1.7c)$$

$$i_I = i(t). \quad (1.7d)$$

Thus the modified nodal analysis leads to

$$\begin{aligned} G(e_1 - e_2) &= i(t) \\ -G(e_1 - e_2) + i_L &= 0 \\ L \frac{di_L}{dt} - e_2 &= 0 \end{aligned} \quad (1.8)$$

The solution is given by

$$\begin{aligned} i_L &= i(t), \\ e_2 &= L \frac{di_L}{dt} = L \frac{di(t)}{dt}, \\ e_1 &= e_2 + G^{-1}i(t) = L \frac{di(t)}{dt} + G^{-1}i(t), \end{aligned}$$

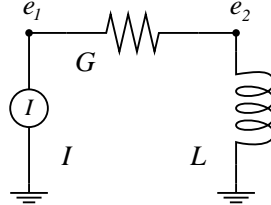


Figure 1.3: Another simple circuit

under the assumption that the current $i(t)$ is differentiable. Notice that all component values are fixed. To solve for e_2 we need to differentiate the current i .

1.3 A transistor amplifier

We will now present a more substantial example adapted from [6]. Consider the transistor amplifier circuit in figure 1.4. P. Rentrop has received this example from K. Glashoff and H.J. Oberle and documented it in [34].

The circuit consists of eight nodes, $U_e(t) = 0.1 \sin(200\pi t)$ is an arbitrary 100 Hz input signal and e_8 , the node potential of the 8th node, is the amplified output. The circuit contains two transistors. We model the behaviour of these semiconductor devices by voltage controlled current sources

$$\begin{aligned} I_{gate} &= (1 - \alpha) g(e_{gate} - e_{source}), \\ I_{drain} &= \alpha g(e_{gate} - e_{source}), \\ I_{source} &= g(e_{gate} - e_{source}) \end{aligned}$$

with a constant $\alpha = 0.99$, g is the nonlinear function

$$g : \mathbb{R} \rightarrow \mathbb{R}, \quad v \mapsto g(v) = \beta \left(\exp \left(\frac{v}{U_F} \right) \right), \quad \beta = 10^{-6}, \quad U_F = 0.026.$$

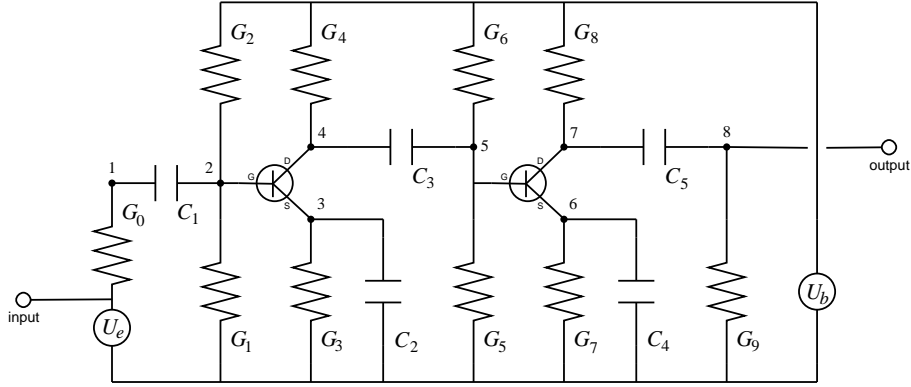


Figure 1.4: Circuit diagram for the transistor amplifier

It is also possible to use PDE models (partial differential equations) to model semiconductor devices. This approach leads to abstract differential-algebraic systems studied in [23, 35, 40].

The modified nodal analysis can now be carried out as in the previous examples. Consider for instance the second node. KCL implies that

$$\begin{aligned}
 0 &= -i_{C_1} - i_{R_1} - i_{R_2} - i_{gate,2} \\
 &= -C_1 v'_{C_1} - v_{G_1} G_1 - v_{G_2} G_2 - (1 - \alpha) g(e_2 - e_3) \\
 &= -C_1 (e_2 - e_1)' - e_2 G_1 - (e_2 - U_b) G_2 + (\alpha - 1) g(e_2 - e_3) \\
 &= C_1 (e_1 - e_2)' - e_2 (G_1 + G_2) + U_b G_2 + (\alpha - 1) g(e_2 - e_3).
 \end{aligned}$$

$U_b = 6$ is the working voltage of the circuit and the remaining constant parameters of the model are chosen to be

$$G_0 = 10^{-3}, \quad G_k = \frac{1}{9} \cdot 10^{-3}, \quad k = 1, \dots, 9, \quad C_k = 10^{-6}, \quad k = 1, \dots, 5.$$

A similar derivation for the other nodes leads to the quasi-linear system

$$A(Dx(t))' = b(x(t)) \tag{1.9}$$

with

$$A = \begin{pmatrix} C_1 & 0 & 0 & 0 & 0 \\ -C_1 & 0 & 0 & 0 & 0 \\ 0 & -C_2 & 0 & 0 & 0 \\ 0 & 0 & C_3 & 0 & 0 \\ 0 & 0 & -C_3 & 0 & 0 \\ 0 & 0 & 0 & -C_4 & 0 \\ 0 & 0 & 0 & 0 & C_5 \\ 0 & 0 & 0 & 0 & -C_5 \end{pmatrix}, \quad b(x) = \begin{pmatrix} -U_e G_0 + e_1 G_0 \\ -U_b G_2 + e_2 (G_1 + G_2) - (\alpha - 1) g(e_2 - e_3) \\ -g(e_2 - e_3) + e_3 G_3 \\ -U_b G_4 + e_4 G_4 + \alpha g(e_2 - e_3) \\ -U_b G_6 + e_5 (G_5 + G_6) - (\alpha - 1) g(e_5 - e_6) \\ -g(e_5 - e_6) + e_6 G_7 \\ -U_b G_8 + e_7 G_8 + \alpha g(e_5 - e_6) \\ e_8 G_9 \end{pmatrix}.$$

A numerical solution of (1.9) can be calculated using DASSL or RADAU5, see [6, 14].

A mathematically more general version of (1.9) is

$$A(x(t), t)(D(t)x(t))' = b(x(t), t) \quad (1.10)$$

with a solution dependent matrix A . We identified x_i with the node potential e_i . Let us assume that $N_0(t) = \ker A(x(t), t)D(t)$ does not depend on x . We will follow [16] and investigate (1.10) in more detail. With

$$f(y, x, t) = A(x(t), t)y - b(x(t), t),$$

(1.10) can be written as

$$f((D(t)x(t))', x(t), t) = 0. \quad (1.11)$$

Denote $B(y, x, t) = f'_x(y, x, t)$ and let $Q(t)$ be a continuous projector function onto $N_0(t)$. Calculate

$$G_1(y, x, t) = A(x, t)D(t) + B(y, x, t)Q(t).$$

For the transistor amplifier (1.11) in figure 1.4 this matrix is always nonsingular. We want to use this matrix in conjunction with the Implicit Function Theorem to derive an ordinary differential equation that determines the dynamical flow of (1.10).

Let $D(t)^-$ be defined by

$$\begin{aligned} DD^-D &= D, & DD^- &= I_5, \\ D^-DD^- &= D^-, & D^-D &= P := I_8 - Q. \end{aligned}$$

I_k denotes the identity in \mathbb{R}^k and $D(t)^-$ is a generalized reflexive inverse of $D(t)$. For more information on generalized matrix inverses see section 2.3.1 on page 18.

For a solution x of (1.11) define

$$u(t) = D(t)x(t), \quad w(t) = D(t)^-u'(t) + Q(t)x(t).$$

Observe that $A(Dx)' = ADw$ and $x = Px + Qx = D^-Dx + Qx = D^-u + Qw$. Thus it holds that

$$(1.11) \Leftrightarrow ADw + b(x, t) \Leftrightarrow F(w, u, t) := f(Dw, D^-u + Qw, t) = 0.$$

Note that

$$u' = R'u + Dw,$$

since $Dw = DD^-u + DQx = (Ru)' = u' - R'u$. The mapping F can be studied without requiring x to be a solution. Let $(y_0, x_0, t_0) \in \mathbb{R}^{5+8+1}$, such that

$$f(y_0, x_0, t_0) = 0.$$

For $w_0 = D(t_0)^-y_0 + Q(t_0)x_0$, $u_0 = D(t_0)x_0$ it follows that

- $F(w_0, u_0, t_0) = f(y_0, x_0, t_0) = 0$,
- $F'_w(w_0, u_0, t_0) = G_1(y_0, x_0, t_0)$ is nonsingular.

Due to the Implicit Function Theorem there is a $\varrho > 0$ and a smooth mapping

$$\omega : B_{\varrho}(u_0, t_0) \times \mathfrak{I} \rightarrow \mathbb{R}^m$$

satisfying

$$\omega(u_0, t_0) = w_0, \quad F(\omega(u, t), u, t) = 0 \quad \forall (u, t) \in B_{\varrho}(u_0, t_0).$$

We use ω to define

$$x(t) = D(t)^{-1}u(t) + Q(t)\omega(u(t), t), \quad t \in \mathfrak{I}.$$

where u is the solution of the ordinary differential equation

$$u'(t) = R'(t)u(t) + D(t)\omega(u(t), t), \quad u(t_0) = D(t_0)x_0. \quad (1.12)$$

x is indeed a solution of (1.10), since

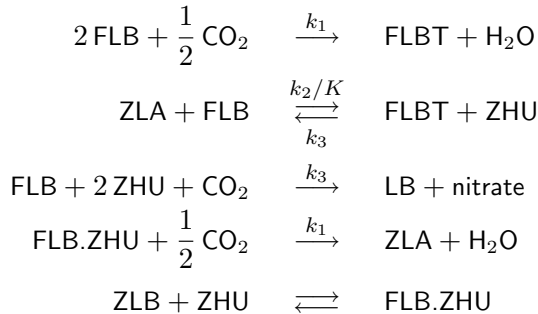
$$f((D(t)x(t))', x(t), t) = f(u', D^{-1}u + Q\omega(u, t), t) = F(\omega, u, t) = 0.$$

This example shows that there is a formulation of the problem in terms of an ordinary differential equation (1.12) as was the case for the mathematical pendulum in the first example. However, (1.12) is available only theoretically as it was obtained using the Implicit Function Theorem. Thus we have to deal directly with the DAE formulation (1.10) when solving the problem. Nevertheless, (1.12) will play a vital part in analyzing (1.10) and in analyzing numerical methods applied to (1.10).

In section 3 it will be shown how (1.12) can be obtained explicitly for linear DAEs. Section 4 is devoted to showing that there are numerical methods that, when applied directly to (1.10), behave as if they were integrating (1.12), given that (1.10) satisfies some additional properties. In this case results concerning convergence and order of numerical methods can be transferred directly from ODE theory to DAEs.

1.4 The Akzo Nobel Problem

The last example originates from the Akzo Nobel Central Research in Arnhem, the Netherlands, and is again taken from [6]. It describes a chemical process in which two species, FLB and ZLU, are mixed while carbon dioxide is continuously added. The resulting species of importance is ZLA. The reaction equations are given in [5].



The last equation describes an equilibrium where the constant

$$K_s = \frac{[\text{FLB} \cdot \text{ZHU}]}{[\text{FLB}] \cdot [\text{ZHU}]}$$

plays a role in parameter estimation. Square brackets denote concentrations.

The chemical process is appropriately described by the reaction velocities

$$\begin{aligned}
r_1 &= k_1 \cdot [\text{FLB}]^4 \cdot [\text{CO}_2]^{\frac{1}{2}}, \\
r_2 &= k_2 \cdot [\text{FLBT}] \cdot [\text{ZHU}], \\
r_3 &= \frac{k_2}{K} \cdot [\text{FLB}] \cdot [\text{ZLA}], \\
r_4 &= k_3 \cdot [\text{FLB}] \cdot [\text{ZHU}]^2, \\
r_5 &= k_4 \cdot [\text{FLB.ZHU}]^2 \cdot [\text{CO}_2]^{\frac{1}{2}},
\end{aligned}$$

see [6] for details. The inflow of carbon dioxide per volume unit is denoted by F_{in} and satisfies

$$F_{in} = \kappa A \cdot \left(\frac{p(\text{CO}_2)}{H} - [\text{CO}_2] \right).$$

κA is the mass transfer coefficient, H the Henry constant and $p(\text{CO}_2)$ is the partial carbon dioxide pressure [6]. It is assumed that $p(\text{CO}_2)$ is independent of $[\text{CO}_2]$. The various constants are given by

$$\begin{aligned}
k_1 &= 18.7, & k_4 &= 0.42, & K_s &= 115.83, \\
k_2 &= 0.58, & K &= 34.4, & p(\text{CO}_2) &= 0.9, \\
k_3 &= 0.09, & \kappa A &= 3.3, & H &= 737.
\end{aligned}$$

If we identify the concentrations $[\text{FLB}]$, $[\text{CO}_2]$, $[\text{FLBT}]$, $[\text{ZHU}]$, $[\text{ZLA}]$, $[\text{FLB.ZHU}]$ with x_1, \dots, x_6 respectively, we obtain the differential-algebraic equation

$$\begin{pmatrix} 1 & & & & & \\ & 1 & & & & \\ & & 1 & & & \\ & & & 1 & & \\ & & & & 1 & \\ & & & & & 0 \end{pmatrix} x'(t) = \begin{pmatrix} -2r_1 + r_2 - r_3 - r_4 \\ -\frac{1}{2}r_1 & & -r_4 & -\frac{1}{2}r_5 + F_{in} \\ r_1 - r_2 + r_3 \\ -r_2 + r_3 - 2r_4 \\ r_2 - r_3 & & +r_5 \\ K_s x_1 x_4 - x_6 \end{pmatrix}. \quad (1.13)$$

This DAE can be analyzed in a similar way as the previous example. The matrix

$$G_1 = AD + BQ = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0.42x_6\sqrt{x_2} \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -0.84x_6\sqrt{x_2} \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

is always nonsingular. Here, $A = D = \text{diag}(1, 1, 1, 1, 1, 0)$ was chosen.

2 Index concepts for DAEs

In the last section we saw that DAEs differ in many ways from ordinary differential equations. For instance the circuit in figure 1.3 lead to a DAE where a differentiation process is involved when solving the equations. This differentiation needs to be carried out numerically, which is an unstable operation. Thus there are some problems to be expected when solving these systems. In this section we try to measure the difficulties arising in the theoretical and numerical treatment of a given DAE.

2.1 The Kronecker index

Let's take linear differential-algebraic equations with constant coefficients as a starting point. These equations are given as

$$Ex'(t) + Fx(t) = q(t), \quad t \in \mathcal{J}, \quad (2.1)$$

with $E, F \in L(\mathbb{R}^m)$. Even for (2.1) existence and uniqueness of solutions is not apriori clear.

Example 2.1 For the DAE

$$\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} x'(t) + \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} x(t) = 0$$

a solution $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ is given by $x_2(t) = g(t)$ and $x_1(t) = -\int_{t_0}^t g(s) ds$, where the function $g \in C(\mathcal{J}, \mathbb{R})$ can be chosen arbitrarily. \square

In order to exclude examples like 2.1 we consider the matrix pencil $\lambda E + F$. The pair (E, F) is said to form a regular matrix pencil, if there is a λ such that $\det(\lambda E + F) \neq 0$. A simultaneous transformation of E and F into Kronecker normal form makes a solution of (2.1) possible.

Theorem 2.2 (Kronecker [19]) *Let (E, F) form a regular matrix pencil. Then there exist nonsingular matrices U and V such that*

$$UEV = \begin{pmatrix} I & 0 \\ 0 & N \end{pmatrix}, \quad U F V = \begin{pmatrix} C & 0 \\ 0 & I \end{pmatrix},$$

where $N = \text{diag}(N_1, \dots, N_k)$ is a block-diagonal matrix of Jordan-blocks N_i to the eigenvalue 0. \square

The proof can be found in [9] or [15]. Notice that due to the special structure of N there is $\mu \in \mathbb{N}$ such that $N^{\mu-1} \neq 0$ but $N^\mu = 0$. μ is known as N 's index of nilpotency. It does not depend on the special choice of U and V .

We solve (2.1) by introducing the transformation

$$x = V \begin{pmatrix} u \\ v \end{pmatrix}, \quad \begin{pmatrix} a(t) \\ b(t) \end{pmatrix} = Uq(t).$$

Thus (2.1) is equivalent to

$$\begin{aligned} UEV(V^{-1}x(t))' + UFVV^{-1}x(t) &= Uq(t) \\ \Leftrightarrow \begin{pmatrix} I & 0 \\ 0 & N \end{pmatrix} \begin{pmatrix} u(t) \\ v(t) \end{pmatrix}' + \begin{pmatrix} C & 0 \\ 0 & I \end{pmatrix} \begin{pmatrix} u(t) \\ v(t) \end{pmatrix} &= \begin{pmatrix} a(t) \\ b(t) \end{pmatrix}. \end{aligned} \quad (2.2)$$

The first equation is an ordinary differential equation

$$u'(t) + Cu(t) = a(t)$$

for the u component. The second equation reads

$$\begin{aligned} v(t) &= b(t) - Nv'(t) \\ &= b(t) - N(b'(t) - Nv''(t)) = b(t) - Nb'(t) + N^2v''(t) \\ &= \dots = \sum_{i=0}^{\mu-1} (-N)^i b^{(i)}(t) \end{aligned} \quad (2.3)$$

determining the v component completely by repeated differentiation of the right hand side b . Since numerical differentiation is an unstable process, the index μ is a measure of numerical difficulty when solving (2.1).

Definition 2.3 *Let (E, F) form a regular matrix pencil. The (Kronecker) index of (2.1) is 0 if E is nonsingular and μ , i.e. N 's index of nilpotency, otherwise.*

2.2 The differentiation index

How can definition 2.3 be generalized to the case of time dependent coefficients or even to nonlinear DAEs? If we consider (2.3) again, it turns out that

$$v'(t) = \sum_{i=0}^{\mu-1} (-N)^i b^{(i+1)}(t),$$

meaning that exactly μ differentiations transform (2.2) into a system of explicit ordinary differential equations. This idea was generalized by Gear, Petzold and Campbell [4, 10, 11]. The following definition is taken from [15].

Definition 2.4 *The nonlinear DAE*

$$f(x'(t), x(t), t) = 0 \quad (2.4)$$

has (differentiation) index μ if μ is the minimal number of differentiations

$$f(x'(t), x(t), t) = 0, \quad \frac{df(x'(t), x(t), t)}{dt} = 0, \dots, \quad \frac{d^\mu f(x'(t), x(t), t)}{dt^\mu} = 0 \quad (2.5)$$

such that the equations (2.5) allow to extract an explicit ordinary differential system $x'(t) = \varphi(x(t), t)$ using only algebraic manipulations.

We now want to look at four examples to get a feeling of how to calculate the differentiation index. We always assume that the functions involved are smooth enough to apply definition 2.4.

Example 2.5 For linear DAEs with constant coefficients forming a regular matrix pencil we have differentiation index μ if and only if the Kronecker index is μ . \square

Example 2.6 Consider the system

$$x' = f(x, y) \tag{2.6a}$$

$$0 = g(x, y). \tag{2.6b}$$

The second equation yields

$$0 = \frac{dg(x, y)}{dt} = g_x(x, y)x' + g_y(x, y)y'.$$

If $g_y(x, y)$ is nonsingular in a neighbourhood of the solution, (2.6) is transformed to

$$x' = f(x, y) \tag{2.6a}$$

$$y' = -g_y(x, y)^{-1}g_x(x, y)x' = -g_y(x, y)^{-1}g_x(x, y)f(x, y) \tag{2.6b'}$$

and the differentiation index is $\mu = 1$.

The DAE (1.6) modelling the circuit in figure 1.2 is of the form (2.6) with

$$x = e_2, \quad y = \begin{pmatrix} e_1 \\ i_V \end{pmatrix}, \quad f(x, y) = \frac{G}{C}(e_1 - e_2) \quad \text{and} \quad g(x, y) = \begin{pmatrix} G(e_1 - e_2) - i_V \\ e_1 + v_V \end{pmatrix}.$$

Note that $g_y(x, y) = \begin{pmatrix} G & -1 \\ 1 & 0 \end{pmatrix}$ is nonsingular so that (1.6) is an index 1 equation. \square

Example 2.7 The system

$$x' = f(x, y) \tag{2.8a}$$

$$0 = g(x) \tag{2.8b}$$

can be studied in a similar way. (2.8b) gives

$$0 = \frac{dg(x)}{dt} = g_x(x)x' = g_x(x)f(x, y) = h(x, y). \tag{2.8b'}$$

Comparing with example 2.6 we know that (2.8a), (2.8b') is an index 1 system if $h_y(x, y)$ remains nonsingular in a neighbourhood of the solution. If this condition holds, (2.8) is of index 2, as two differentiations produce

$$x' = f(x, y) \tag{2.8a}$$

$$y' = -h_y(x, y)^{-1}h_x(x, y)f(x, y) \\ = -(g_x(x)f_y(x, y))^{-1}\left(g_{xx}(x)(f(x, y), f(x, y)) + g_x(x)f_x(x, y)f(x, y)\right). \tag{2.8b''}$$

(2.8b') defines the “hidden constraint” of the index 2 equation (2.8).

The DAE (1.8) modelling the circuit in figure 1.3 can be written as

$$i'_L = \frac{1}{L}e_2 = f(i_L, e_2) \quad (2.10a)$$

$$0 = i_L - i_I = g(i_L). \quad (2.10b)$$

The remaining variable e_1 is determined by $e_1 = e_2 + G^{-1}i_I$, where i_I is the input current. (2.10) is of the form (2.8) with $x = i_L$ and $y = e_2$. $h_y(x, y) = g_x f_y = 1 \cdot \frac{1}{L}$ is nonsingular and the index is 2. \square

Example 2.8 Finally take a look at the system

$$x' = f(x, y) \quad (2.11a)$$

$$y' = g(x, y, z) \quad (2.11b)$$

$$0 = h(x). \quad (2.11c)$$

Differentiation of (2.11c) yields

$$0 = \frac{dh(x)}{dt} = h_x(x)x' = h_x(x)f(x, y) = \hat{h}(x, y) \quad (2.11c')$$

and

$$\mathbf{r}' = \begin{pmatrix} x \\ y \end{pmatrix}' = \begin{pmatrix} f(x, y) \\ g(x, y, z) \end{pmatrix} = \mathbf{f}(\mathbf{r}, \boldsymbol{\eta}) \quad (2.11a) \text{ and } (2.11b) \quad (2.12a)$$

$$0 = \hat{h}(x, y) = \mathbf{g}(\mathbf{r}) \quad (2.11c') \quad (2.12b)$$

is of the form (2.8) with $\mathbf{r} = \begin{pmatrix} x \\ y \end{pmatrix}$ and $\boldsymbol{\eta} = z$. Define

$$\mathbf{h}(\mathbf{r}, \boldsymbol{\eta}) = \mathbf{g}_{\mathbf{r}}(\mathbf{r})\mathbf{f}(\mathbf{r}, \boldsymbol{\eta})$$

and compare with (2.8b') to find that (2.12) is of the index 2 if

$$\begin{aligned} \mathbf{h}_{\boldsymbol{\eta}}(\mathbf{r}, \boldsymbol{\eta}) &= \mathbf{g}_{\mathbf{r}}(\mathbf{r})\mathbf{f}_{\boldsymbol{\eta}}(\mathbf{r}, \boldsymbol{\eta}) = \begin{pmatrix} \hat{h}_x & \hat{h}_y \end{pmatrix} \begin{pmatrix} f_z \\ g_z \end{pmatrix} \\ &= (h_{xx}(f, \cdot) + h_x f_x \quad h_{xy}(f, \cdot) + h_x f_y) \begin{pmatrix} 0 \\ g_z \end{pmatrix} = h_x f_y g_z \end{aligned}$$

remains nonsingular. This shows that (2.11) is an index 3 system if the matrix $h_x(x)f_y(x, y)g_z(x, y, z)$ is invertible in a neighbourhood of the solution (x, y, z) .

Hidden constraints are given by (2.11c') but also by

$$\mathbf{h}(\mathbf{r}, \boldsymbol{\eta}) = \mathbf{g}_{\mathbf{r}}(\mathbf{r})\mathbf{f}(\mathbf{r}, \boldsymbol{\eta}) = h_{xx}(f, f) + h_x f_x f + h_x f_y g = 0,$$

which is condition (2.8b') in terms of the index 2 system (2.12).

Consider again the mathematical pendulum from section 1.1 in the formulation

$$\begin{aligned} x' = u &= f_1(x, y, u, v) \\ y' = v &= f_2(x, y, u, v) \\ u' = -\frac{2}{m}\lambda x &= g_1(x, y, u, v, \lambda) \\ v' = +g - \frac{2}{m}\lambda y &= g_2(x, y, u, v, \lambda) \\ 0 = x^2 + y^2 - l^2 &= h(x, y). \end{aligned} \quad (2.13)$$

For $l > 0$ the value $h_{(x,y)}f_{(u,v)}g_{\lambda} = -\frac{4}{m}(x^2 + y^2)$ is always nonsingular so that (2.13) is an index 3 problem. \square

2.3 The tractability index

In definition 2.4 the function f is assumed to be smooth enough to calculate the derivatives (2.5). In applications this smoothness is often not given. For instance in circuit simulation input signals are continuous but often not differentiable.

In this section we want to study the tractability index introduced by Griepentrog, März [13]. In fact we consider the generalization of the tractability index proposed by März [25]. The idea is to replace the smoothness requirements for the coefficients by the requirement on certain subspaces to be smooth.

To define the tractability index we introduce linear DAEs with properly stated leading terms. A second matrix $D(t)$ is used when formulating the DAE as

$$A(t)(D(t)x(t))' + B(t)x(t) = q(t). \quad (2.14)$$

In contrast to the standard formulation

$$E(t)x(t)' + F(t)x(t) = q(t) \quad (2.15)$$

the leading term in (2.14) precisely figures out which derivatives are actually involved.

The formulation (2.14) was first used in [1] to study linear DAEs and their adjoint equations. For (2.15) the adjoint equation

$$(E^*y)' - F^*y = p$$

is of a different type. For the more general formulation (2.14) the adjoint equation fits nicely into this general form:

$$D^*(A^*y)' - B^*y = p.$$

In this section we consider linear DAEs (2.14) with matrix coefficients

$$A \in C(\mathfrak{J}, L(\mathbb{R}^n, \mathbb{R}^m)), \quad D \in C(\mathfrak{J}, L(\mathbb{R}^m, \mathbb{R}^n)), \quad B \in C(\mathfrak{J}, L(\mathbb{R}^m)).$$

Neither A nor D needs to be a projector function. Note that $A(t)$ and $D(t)$ are rectangular matrices in general. However, A and D are assumed to be well matched in the following sense.

Definition 2.9 *The leading term of (2.14) is properly stated if*

$$\ker A(t) \oplus \operatorname{im} D(t) = \mathbb{R}^n, \quad t \in \mathfrak{J},$$

and there is a continuously differentiable projector function $R \in C^1(\mathfrak{J}, L(\mathbb{R}^n))$ with

$$\operatorname{im} R(t) = \operatorname{im} D(t), \quad \ker R(t) = \ker A(t) \quad t \in \mathfrak{J}.$$

By definition $A(t)$ and $D(t)$ have a common constant rank if the leading term is properly stated [25].

Definition 2.10 *A function $x : \mathfrak{J} \rightarrow \mathbb{R}^m$ is said to be a solution of (2.14) if*

$$x \in C_D^1(\mathfrak{J}, \mathbb{R}^m) = \{x \in C(\mathfrak{J}, \mathbb{R}^m) \mid Dx \in C^1(\mathfrak{J}, \mathbb{R}^n)\}$$

satisfies (2.14) pointwise.

Let us point out that a solution x is a continuous function, but the part $Dx : \mathfrak{J} \rightarrow \mathbb{R}^n$ is differentiable.

We now define a sequence of matrix functions and possibly time-varying subspaces. All relations are meant pointwise for $t \in \mathfrak{J}$. Let $G_0 = AD$, $B_0 = B$ and for $i \geq 0$

$$\left. \begin{aligned} N_i &= \ker G_i, \\ S_i &= \{z \in \mathbb{R}^m \mid B_i z \in \operatorname{im} G_i\} = \{z \in \mathbb{R}^m \mid Bz \in \operatorname{im} G_i\}, \\ Q_i &= Q_i^2, \quad \operatorname{im} Q_i = N_i, \quad P_i = I - Q_i, \\ G_{i+1} &= G_i + B_i Q_i, \\ B_{i+1} &= B_i P_i - G_{i+1} D^- C'_{i+1} D P_0 \cdots P_i, \\ C_{i+1} &= D P_0 \cdots P_{i+1} D^-. \end{aligned} \right\} \quad (2.16)$$

Here, $D^- : \mathfrak{J} \rightarrow L(\mathbb{R}^n, \mathbb{R}^m)$ denotes the reflexive generalized inverse of D such that

$$DD^-D = D, \quad D^-DD^- = D^-, \quad DD^- = R, \quad D^-D = P_0. \quad (2.17)$$

Note that D^- is uniquely determined by (2.17) and depends only on the choice of Q_0 . Section 2.3.1 contains more details about generalized matrix inverses.

Definition 2.11 *The DAE (2.14) with properly stated leading term is said to be a regular DAE with tractability index μ on the interval \mathfrak{J} if there is a sequence (2.16) such that*

- G_i has constant rank r_i on \mathfrak{J} ,
- $Q_i \in C(\mathfrak{J}, L(\mathbb{R}^m))$, $DP_0 \cdots P_i D^- \in C^1(\mathfrak{J}, L(\mathbb{R}^n))$, $i \geq 0$,
- $Q_{i+1} Q_j = 0$, $j = 0, \dots, i$, $i \geq 0$,
- $0 \leq r_0 \leq \dots \leq r_{\mu-1} < m$ and $r_\mu = m$.

(2.14) is said to be a regular DAE if it is regular with some index μ .

This index criterion does not depend on the special choice of the projector functions Q_i [28]. As proposed in [24] the sequence (2.16) can be calculated automatically. Thus the index can be calculated without the use of derivative arrays [27].

Example 2.12 Consider the DAE

$$\begin{pmatrix} t \\ 1 \end{pmatrix} \left(\begin{pmatrix} -1 & t \end{pmatrix} \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} \right)' + \begin{pmatrix} 1 & -t \\ 0 & 0 \end{pmatrix} \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} = 0$$

taken from [25]. With $\ker A(t) = \{0\}$, $\operatorname{im} D(t) = \mathbb{R}$ the leading term is properly stated. Calculate

$$G_0(t) = A(t)D(t) = \begin{pmatrix} -t & t^2 \\ -1 & t \end{pmatrix} \quad \text{and} \quad N_0(t) = \{z \in \mathbb{R}^2 \mid \exists \alpha \in \mathbb{R}, z = \alpha \begin{pmatrix} t \\ 1 \end{pmatrix}\}$$

to find that $N_0(t) \subset \ker B(t)$. Independently of the choice of Q_0 in (2.16) we have

$$G_1(t) = G_0(t) + B(t)Q(t) = G_0(t).$$

Similarly it follows that $G_i(t) = G_0(t)$ for every $i \geq 0$. This is *not* a regular DAE in the sense of definition 2.11. Note that for every $\gamma \in C(\mathfrak{J}, \mathbb{R})$ a solution is given by $x(t) = \gamma(t) \begin{pmatrix} t \\ 1 \end{pmatrix}$. Solutions are therefore not uniquely determined. This is the case in spite of the fact that for every t the local matrix pencil $\lambda AD + (B + AD')$ of the reformulated DAE

$$0 = A(t)D(t)x'(t) + (B(t) + A(t)D'(t))x(t) = \begin{pmatrix} -t & t^2 \\ -1 & t \end{pmatrix} x'(t) + x(t)$$

is regular. \square

The following lemma shows that definition 2.11 is indeed a generalization of the Kronecker index, i.e. in the case of constant coefficients, the Kronecker index and the tractability index for regular DAEs coincide. To show this, define the subspaces

$$S_{EF} = \{z \in \mathbb{R}^m \mid Fz \in \text{im } E\}, \quad N_E = \ker E.$$

for given matrices $E, F \in L(\mathbb{R}^m)$. Obviously for fixed $t \in \mathfrak{J}$ we have $N_i(t) = N_{G_i(t)}$ and $S_i(t) = S_{G_i(t)B_i(t)}$ in sequence (2.16).

Lemma 2.13 *For matrices $E, F \in L(\mathbb{R}^m)$ the following statements are equivalent:*

- 1° $N_E \cap S_{EF} = \{0\}$
- 2° For every projector Q_E onto N_E the matrix $E + FQ_E$ is nonsingular.
- 3° $N_E \oplus S_{EF} = \mathbb{R}^m$
- 4° (E, F) form a regular matrix pencil with Kronecker index 1.

Proof: (1° \Rightarrow 2°) $(E + FQ_E)z = 0$ implies $Q_E z \in S_{EF}$. Since $Q_E z \in N_E$ too, we have $Q_E z \in N_E \cap S_{EF} = \{0\}$ and $Q_E z = 0$. Thus $0 = Ez + FQ_E z = Ez$ and $z \in N_E = \text{im } Q_E$. Therefore $z = Q_E z = 0$.

(2° \Rightarrow 3°) $G_{EF} = E + FQ_E$ is nonsingular. Show that $Q_* = Q_E G_{EF}^{-1} F$ is the projector onto N_E along S_{EF} .

(3° \Rightarrow 4°) There is exactly one projector Q_* onto N_E along S_{EF} . Since 3° \Rightarrow 1° \Rightarrow 2°, we find $Q_* = Q_* G_{EF}^{-1} F$ with $G_{EF} = E + FQ_*$. Let $P_* = I - Q_*$.

Show that $\lambda E + F$ is nonsingular for $\lambda \notin \text{spec}(P_* G_{EF}^{-1} F)$ so that (E, F) form a regular matrix pencil. Due to theorem 2.2 there are nonsingular matrices $U, V \in GL_{\mathbb{R}}(m)$ such that

$$V E U = \begin{pmatrix} I & \\ & N \end{pmatrix} = \bar{E}, \quad V F U = \begin{pmatrix} C & \\ & I \end{pmatrix} = \bar{F}.$$

It follows that $N_{\bar{E}} = \ker \bar{E} = U^{-1} N_E$ and $S_{\bar{E}\bar{F}} = \{z \in \mathbb{R}^m \mid \bar{F}z \in \text{im } \bar{E}\} = U^{-1} S_{EF}$ so that

$$N_{\bar{E}} \cap S_{\bar{E}\bar{F}} = U^{-1}(N_E \cap S_{EF}) = \{0\}. \quad (2.18)$$

On the other hand

$$\begin{aligned} N_{\bar{E}} &= \left\{ \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \in \mathbb{R}^m \mid z_1 = 0, z_2 \in \ker N \right\} \quad \text{and} \\ S_{\bar{E}\bar{F}} &= \left\{ \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \in \mathbb{R}^m \mid \begin{pmatrix} C z_1 \\ z_2 \end{pmatrix} \in \text{im } \bar{E} \right\} = \left\{ \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \in \mathbb{R}^m \mid z_2 \in \text{im } N \right\}, \end{aligned}$$

meaning that $\text{im } N \cap \ker N = \{0\}$ and $N = 0$. Thus the Kronecker index is 1.

($4^\circ \Rightarrow 1^\circ$) Kronecker index 1 gives $N = 0$ and $S_{\bar{E}\bar{F}} = \{0\}$, $N_{\bar{E}} \cap S_{\bar{E}\bar{F}} = \{0\}$. Use (2.18) to see $N_E \cap S_{EF} = U(N_{\bar{E}} \cap S_{\bar{E}\bar{F}}) = \{0\}$. \square

As in the previous section we now want to calculate the index of the DAEs modelling the electrical circuits in figure 1.2 and 1.3.

Example 2.14 For (1.6) we calculate $G_0 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & 0 \end{pmatrix}$ and $N_0 = \mathbb{R} \times \{0\} \times \mathbb{R}$. Choose $Q_0 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ to find that $G_1 = \begin{pmatrix} G & 0 & -1 \\ -G & C & 0 \\ -1 & 0 & 0 \end{pmatrix}$ is nonsingular. For the circuit in figure 1.2 we therefore have index 1. \square

Example 2.15 Equation (1.8) can be written as

$$\begin{pmatrix} 0 \\ 0 \\ L \end{pmatrix} \left((0 \ 0 \ 1) \begin{pmatrix} e_1 \\ e_2 \\ i_L \end{pmatrix} \right)' + \begin{pmatrix} G & -G & 0 \\ -G & G & 1 \\ 0 & -1 & 0 \end{pmatrix} \begin{pmatrix} e_1 \\ e_2 \\ i_L \end{pmatrix} = \begin{pmatrix} i(t) \\ 0 \\ 0 \end{pmatrix} \quad (1.8')$$

leading to $G_0 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & L \end{pmatrix}$, $N_0 = \mathbb{R} \times \mathbb{R} \times \{0\}$. With $Q_0 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$ it turns out that $G_1 = \begin{pmatrix} G & 0 & 0 \\ -G & G & 0 \\ 0 & -1 & L \end{pmatrix}$ is singular and $N_1 = \{z \in \mathbb{R}^3 \mid \exists \alpha \in \mathbb{R}, z_1 = z_2 = \alpha L, z_3 = \alpha\}$. $Q_1 = \begin{pmatrix} 0 & 0 & L \\ 0 & 0 & L \\ 0 & 0 & 1 \end{pmatrix}$ is a projector onto N_1 satisfying $Q_1 Q_0 = 0$. Finally $G_2 = \begin{pmatrix} G & -G & 0 \\ -G & G & 1 \\ 0 & -1 & L \end{pmatrix}$ is nonsingular. Thus the index is 2. Note that the terms C'_{i+1} disappear in (2.16) as Q_0 does not depend on t . \square

Nevertheless, in general the derivatives of C_{i+1} appearing in the definition of B_{i+1} in sequence (2.16) are necessary in order to determine the index correctly. We will illustrate this in the next example which can be found in [25] as well.

Example 2.16 The DAE

$$x'_2 = q_1 - x_1 = f(x_1) \quad (2.19a)$$

$$x'_3 = q_2 - (1 + \eta)x_2 - \eta t(q_1 - x_1) = g(x_1, x_2, x_3) \quad (2.19b)$$

$$0 = q_3 - \eta t x_2 - x_3 = h(x_2, x_3) \quad (2.19c)$$

is easily checked to have (differentiation) index 3 as repeated differentiation of (2.19c) yields

$$0 = q'_3 - q_2 + x_2$$

$$0 = q''_3 - q'_2 + q_1 - x_1$$

$$x'_1 = q'''_3 - q''_2 + q'_1.$$

The index does not depend on the value of η . We now write (2.19) as

$$\begin{pmatrix} 1 & 0 \\ \eta t & 1 \\ 0 & 0 \end{pmatrix} \left(\begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \right)' + \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 + \eta & 0 \\ 0 & \eta t & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} q_1 \\ q_2 \\ q_3 \end{pmatrix} \quad (2.19')$$

with a properly stated leading term and calculate the sequence (2.16)

$$\begin{aligned} G_0 &= \begin{pmatrix} 0 & 1 & 0 \\ 0 & \eta t & 1 \\ 0 & 0 & 0 \end{pmatrix}, & Q_0 &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, & G_1 &= \begin{pmatrix} 1 & 1 & 0 \\ 0 & \eta t & 1 \\ 0 & 0 & 0 \end{pmatrix}, & Q_1 &= \begin{pmatrix} 0 & -1 & 0 \\ 0 & 1 & 0 \\ 0 & -\eta t & 0 \end{pmatrix}, \\ G_2 &= \begin{pmatrix} 1 & 1 & 0 \\ 0 & \eta t + 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}, & Q_2 &= \begin{pmatrix} 0 & \eta t & 1 \\ 0 & -\eta t & -1 \\ 0 & \eta t(\eta t + 1) & \eta t + 1 \end{pmatrix}, & G_3 &= \begin{pmatrix} 1 & 1 & 0 \\ 0 & \eta t + 1 & 1 \\ 0 & \eta t & 1 \end{pmatrix}. \end{aligned}$$

Since $\det G_3 = 1$, (2.19') is a regular DAE with index 3 independently of η . However, if we dropped the terms C_{i+1} in (2.16) and defined $\mathcal{G}_{i+1} = \mathcal{G}_i + \mathcal{B}_i Q_i$, $\mathcal{B}_{i+1} = \mathcal{B}_i P_i$ with $\mathcal{G}_0 = AD$ and $\mathcal{B}_0 = B$ we would obtain

$$\mathcal{G}_2 = \begin{pmatrix} 1 & 1 & 0 \\ 0 & \eta t + 1 + \eta & 1 \\ 0 & 0 & 0 \end{pmatrix}, \quad \mathcal{Q}_2 = \frac{1}{1+\eta} \begin{pmatrix} 0 & \eta t & 1 \\ 0 & -\eta t & -1 \\ 0 & (\eta t + 1 + \eta)\eta t & \eta t + 1 + \eta \end{pmatrix}, \quad \mathcal{G}_3 = \begin{pmatrix} 1 & 1 & 0 \\ 0 & \eta t + 1 + \eta & 1 \\ 0 & \eta t & 1 \end{pmatrix}.$$

$\det \mathcal{G}_3 = 1 + \eta$ shows that \mathcal{G}_3 is singular for $\eta = -1$. Thus the use of the simpler version of B_i would lead to an index criterion not recognizing the index properly. \square

The previous example gives rise to investigating the relationship between G_i and \mathcal{G}_i further. Due to $G_i P_i = G_i$ the matrix G_{i+1} may be written as

$$G_{i+1} = (G_i + B_{i-1} P_{i-1} Q_i)(I - P_i D^- C_i' D P_0 \cdots P_{i-1} Q_i).$$

For low indices we thus find

$$G_0 = \mathcal{G}_0, \quad G_1 = \mathcal{G}_1, \quad G_2 = \mathcal{G}_2(I - P_1 D^- C_1' D P_0 Q_1)$$

with the nonsingular factor $I - P_1 D^- C_1' D P_0 Q_1$. The matrices G_2 and \mathcal{G}_2 have therefore common rank and we had to choose an index 3 example in 2.16 to show the necessity of the second term in the definition of B_{i+1} .

We don't have to restrict ourselves to linear DAEs (2.14). Nonlinear DAEs

$$A(x(t), t)(d(x, t), t)' + b(x(t), t) = 0 \tag{2.20}$$

can also be considered. For (2.20) the index μ is defined in such a way that all linearizations along solutions have the same index μ in the sense of definition 2.11. The index 1 case is studied extensively in [16]. We have already made use of this approach when investigating the transistor amplifier example in section 1.3. More information on nonlinear DAEs can be found in [27, 29].

2.3.1 Some technical details

In order to define the sequence (2.16) we introduced the generalized reflexive inverse D^- of D . Here we want to provide a short summary of the properties of generalized matrix inverses [41].

For a rectangular matrix $M \in L(\mathbb{R}^m, \mathbb{R}^n)$, a matrix $\tilde{M} \in L(\mathbb{R}^n, \mathbb{R}^m)$ is called a generalized inverse of M if

$$\tilde{M} M \tilde{M} = \tilde{M}.$$

If the condition

$$M\tilde{M}M = M$$

holds as well, then \tilde{M} is called a reflexive generalized inverse of M . Observe that for any reflexive generalized inverse \tilde{M} of M the matrices

$$(M\tilde{M})^2 = M\tilde{M}M\tilde{M} = M\tilde{M}, \quad (\tilde{M}M)^2 = \tilde{M}M\tilde{M}M = \tilde{M}M$$

are projectors. Reflexive generalized inverses are not uniquely determined. Uniqueness is obtained if we require $M\tilde{M}$ and $\tilde{M}M$ to be *special* projectors. We could, for instance, require them to be ortho-projectors

$$(M\tilde{M})^T = M\tilde{M}, \quad (\tilde{M}M)^T = \tilde{M}M.$$

In this case \tilde{M} is called the Moore-Penrose inverse of M , often denoted by M^+ .

In the case of DAEs with properly stated leading terms we appropriated the projectors $P_0(t) \in L(\mathbb{R}^m)$ and $R(t) \in L(\mathbb{R}^n)$ to determine $D^-(t) \in \mathbb{R}^n, \mathbb{R}^m$ uniquely. $D^-(t)$ is the reflexive generalized inverse of $D(t)$ defined by

$$DD^-D = D, \quad D^-DD^- = D^-, \quad DD^- = R, \quad D^-D = P_0. \quad (2.21)$$

If there was another generalized inverse \tilde{D}^- satisfying (2.21), then

$$\tilde{D}^- = \tilde{D}^-D\tilde{D}^- = \tilde{D}^-R = \tilde{D}^-DD^- = P_0DD^- = D^-DD^- = D^-.$$

In definition 2.11 the condition

$$Q_{i+1}Q_j = 0, \quad j = 0, \dots, i, \quad i \geq 0 \quad (2.22)$$

is required. We will show briefly that the projectors Q_i in sequence (2.16) can always be chosen to satisfy (2.22).

If for a given DAE (2.14) there was an index i_* such that $N_{i_*+1} \cap N_{i_*} \neq \{0\}$, then (2.14) would not be a regular DAE as all G_i would be singular. Thus $N_0 \cap N_1 = \{0\}$ is a necessary condition for a regular DAE and the projector Q_1 onto N_1 can be chosen such that $N_0 \subset \ker Q_1$.

For an index $i \geq 1$ let the projectors Q_j for $j = 1, \dots, i$ satisfy $Q_jQ_k = 0$, $k = 0, \dots, j-1$. Then $N_{i+1} \cap N_i = \{0\}$ implies $N_{i+1} \cap N_j = \{0\}$ for $j = 1, \dots, i$ and Q_{i+1} can be chosen such that $N_0 \oplus N_1 \oplus \dots \oplus N_i \subset \ker Q_{i+1}$.

2.4 Other index concepts

As seen in the previous sections a DAE can be assigned an index in several ways. In the case of linear equations with constant coefficients all index notions coincide with the Kronecker index. Apart from that, each index definition stresses different aspects of the DAE under consideration. While the differentiation index aims at finding possible reformulations in terms of ordinary differential equations, the tractability index is used to study DAEs without the use of derivative arrays.

There are several other index concepts available. Here we want to introduce some of them briefly.

2.4.1 The perturbation index

The perturbation index was introduced for nonlinear DAEs

$$f(x'(t), x(t)) = 0 \tag{2.23}$$

by Hairer, Lubich and Roche in [14]. (2.23) has perturbation index μ along a solution x on $\mathfrak{J} = [0, T]$ if μ is the smallest integer such that, for all functions \hat{x} having a defect

$$f(\hat{x}'(t), \hat{x}(t)) = \delta(t),$$

there exists on \mathfrak{J} an estimate

$$\|\hat{x}(t) - x(t)\| \leq C(\|\hat{x}(0) - x(0)\| + \max_{0 \leq \xi \leq t} \|\delta(\xi)\| + \dots + \max_{0 \leq \xi \leq t} \|\delta^{\mu-1}(\xi)\|)$$

whenever the expression on the right-hand side is sufficiently small. Here C denotes a constant which depends only on f and the length of \mathfrak{J} .

The perturbation index measures the sensitivity of solutions with respect to perturbations of the given problem [15].

2.4.2 The geometric index

Here we present the geometric index as it is introduced in [38]. Consider the autonomous DAE

$$f(x', x) = 0 \tag{2.24}$$

and assume that $M_0 = f^{-1}(0)$ is a smooth submanifold of $\mathbb{R}^m \times \mathbb{R}^m$. Then the DAE (2.24) can be written as

$$(x', x) \in M_0.$$

Each solution has to satisfy $x \in W_0 = \pi(M_0)$, where $\pi : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ is the canonical projection onto the second component. If W_0 is a submanifold of \mathbb{R}^m , then (x', x) belongs to the tangent bundle TW_0 of W_0 . In other words

$$(x', x) \in M_1 = M_0 \cap TW_0.$$

M_1 is called the first reduction of M_0 . Iterate this process to obtain a sequence M_0, M_1, M_2, \dots of manifolds where M_{i+1} is the first reduction of M_i and

$$(x', x) \in \bigcap_{i \geq 0} M_i.$$

The geometric index is defined as the smallest integer μ such that $M_\mu = M_{\mu+1}$. This index notion was introduced in [33] and studied extensively in [31] by Rabier and Rheinboldt.

2.4.3 The strangeness index

This index notion is a generalization of the Kronecker index to DAEs

$$E(t)x'(t) + F(t)x(t) = q(t), \quad t \in \mathfrak{J} \subset \mathbb{R}, \quad (2.25)$$

with time-dependent coefficients. It is due to Kunkel and Mehrmann [22]. The matrices U and V in theorem 2.2 now depend on t , i.e. (2.25) is transformed to

$$UEVy' + (UFV - UEV')y = Uq \quad \Leftrightarrow \quad \hat{E}y' + \hat{F}y = \hat{q}.$$

The pairs of matrix functions (E, F) and (\hat{E}, \hat{F}) are said to be globally equivalent. For fixed $t \in \mathfrak{J}$ define matrices $T(t)$, $\hat{T}(t)$, $Z(t)$ and $V(t)$ such that the column vectors of $T(t)$, $\hat{T}(t)$, $Z(t)$ and $V(t)$ span the subspaces $\ker E(t)$, $\text{im } E^T$, $\ker E^T$ and $\text{im}(Z(t)^T N(t)T(t))^\perp$, respectively. Use these matrices to define

$$\begin{aligned} r(t) &= \text{rank } E(t), & d(t) &= r(t) - s(t), \\ a(t) &= \text{rank } (Z(t)^T N(t)T(t)), & u(t) &= m - r(t) - a(t) - s(t), \\ s(t) &= \text{rank } (V(t)^T Z(t)^T N(t)\hat{T}(t)). \end{aligned}$$

We assume that the functions r , s and a are constant on \mathfrak{J} . Then (E, F) is globally equivalent to the pair

$$\left(\left(\begin{array}{cccc} 0 & 0 & 0 & 0 \\ 0 & I_d & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right), \left(\begin{array}{ccccc} 0 & F_{12} & 0 & F_{14} & F_{15} \\ 0 & 0 & 0 & F_{24} & F_{25} \\ 0 & 0 & I_a & 0 & 0 \\ I_s & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array} \right) \right) = (E_1, F_1).$$

The proof can be found in [21]. The value s is called the strangeness of the pair (E, F) . Denote (E, F) by (E_0, F_0) and $s_0 = s$. Similarly we define the strangeness s_1 of the pair (E_1, F_1) . If we repeat the procedure described above, we arrive at a sequence of globally equivalent pairs (E_i, F_i) , $i \geq 0$, each having strangeness s_i . The strangeness index or s-index is then defined by

$$\mu = \min\{i = 0, 1, 2, \dots \mid s_i = 0\}.$$

Relations between the tractability index and the strangeness index are given in [36].

3 Solvability of linear DAEs with properly stated leading term

In this section we consider linear differential-algebraic equations

$$A(t)(D(t)x(t))' + B(t)x(t) = q(t), \quad t \in \mathfrak{J} \quad (3.1)$$

with properly stated leading terms as in definition 2.9. A , B and D are continuous matrix functions with

$$D(t) \in L(\mathbb{R}^m, \mathbb{R}^n), \quad A(t) \in L(\mathbb{R}^n, \mathbb{R}^m)$$

$$B(t) \in L(\mathbb{R}^m, \mathbb{R}^m), \quad q(t) \in \mathbb{R}^m$$

A function $x : \mathfrak{J} \rightarrow \mathbb{R}^m$ is said to be a solution of (3.1) if

$$x \in C_D^1(\mathfrak{J}, \mathbb{R}^m) = \{x \in C(\mathfrak{J}, \mathbb{R}^m) \mid Dx \in C^1(\mathfrak{J}, \mathbb{R}^n)\}$$

satisfies (3.1) pointwise.

As in the previous section we define for $t \in \mathfrak{J}$ pointwise $G_0 = AD$, $B_0 = B$ and for $i \geq 0$

$$\left. \begin{aligned} N_i &= \ker G_i, \\ S_i &= \{z \in \mathbb{R}^m \mid B_i z \in \text{im } G_i\} = \{z \in \mathbb{R}^m \mid Bz \in \text{im } G_i\}, \\ Q_i &= Q_i^2, \quad \text{im } Q_i = N_i, \quad P_i = I - Q_i, \\ G_{i+1} &= G_i + B_i Q_i, \\ B_{i+1} &= B_i P_i - G_{i+1} D^- C'_{i+1} D P_0 \cdots P_i, \\ C_{i+1} &= D P_0 \cdots P_{i+1} D^-. \end{aligned} \right\} \quad (3.2)$$

D^- is again the reflexive generalized inverse of D from section 2.3.

For completeness we repeat the definition of index μ from the previous section.

Definition 3.1 *The DAE (3.1) with properly stated leading term is said to be a regular DAE with tractability index μ on the interval \mathfrak{J} if there is a sequence (3.2) such that*

- G_i has constant rank r_i on \mathfrak{J} ,
- $Q_i \in C(\mathfrak{J}, L(\mathbb{R}^m))$, $D P_0 \cdots P_i D^- \in C^1(\mathfrak{J}, L(\mathbb{R}^n))$, $i \geq 0$,
- $Q_{i+1} Q_j = 0$, $j = 0, \dots, i$, $i \geq 0$,
- $0 \leq r_0 \leq \dots \leq r_{\mu-1} < m$ and $r_\mu = m$.

(3.1) is said to be a regular DAE if it is regular with some index μ .

The material presented here is mainly taken from [25], [1] and [26].

3.1 Decoupling of linear index-1 DAEs

Let (3.1) be a regular index 1 DAE with properly stated leading term. Due to definition 3.1 the Matrix G_1 is nonsingular.

Lemma 3.2 *The matrices of sequence (3.2) satisfy*

$$(a) P_0 = D^-D, P_0D^- = D^-, DP_0 = D, DP_0D^- = DD^- = R,$$

$$(b) RD = DD^-D = D,$$

$$(c) A = AR = ADD^-,$$

$$(d) Q_0 = G_1^{-1}BQ_0,$$

$$(e) P_0 = G_1^{-1}AD,$$

$$(f) P_0x = P_0y \Leftrightarrow DP_0x = DP_0y \underset{(a)}{\Leftrightarrow} Dx = Dy.$$

Proof: (a) and (b) are just the properties of the generalized reflexive inverse D^- . Remember that $R \in C^1(\mathcal{J}, L(\mathbb{R}^n))$ is the smooth projector function realizing the decomposition $\ker A(t) \oplus \text{im } D(t) = \mathbb{R}^n$ provided by the properly stated leading term. $\ker A = \ker R$ implies (c). $G_1Q_0 = ADQ_0 + BQ_0^2 = BQ_0$ proves (d). Similarly $G_1P_0 = ADP_0 + BQ_0P_0 = AD$ shows (e). For (f) we only have to show “ \Leftarrow ”. If $DP_0z = 0$ then $P_0z \in \ker D = \ker AD = \ker P_0$ and thus $P_0z = 0$. $\ker D = \ker AD$ holds due to the properly stated leading term. \square

Let's assume that x is a solution of the DAE (3.1). Scaling with G_1^{-1} yields

$$A(Dx)' + Bx = q \Leftrightarrow G_1^{-1}A(Dx)' + G_1^{-1}Bx = G_1^{-1}q. \quad (3.3)$$

Note that

- $G_1^{-1}A(Dx)' \underset{(c)}{=} G_1^{-1}AR(Dx)' = G_1^{-1}ADD^-(Dx)' \underset{(e)}{=} P_0D^-(Dx)',$
- $G_1^{-1}Bx = G_1^{-1}BP_0x + G_1^{-1}BQ_0x \underset{(d)}{=} G_1^{-1}BP_0x + Q_0x.$

Thus multiplication of (3.3) by P_0 and Q_0 from the left shows that

$$\begin{aligned} A(Dx)' + Bx = q &\Leftrightarrow G_1^{-1}A(Dx)' + G_1^{-1}Bx = G_1^{-1}q \\ &\Leftrightarrow \left\{ \begin{array}{l} P_0D^-(Dx)' + P_0G_1^{-1}BP_0x = P_0G_1^{-1}q \\ Q_0G_1^{-1}BP_0x + Q_0x = Q_0G_1^{-1}q \end{array} \right\} \\ &\underset{(f)}{\Leftrightarrow} \left\{ \begin{array}{l} DP_0D^-(Dx)' + DP_0G_1^{-1}BP_0x = DP_0G_1^{-1}q \\ Q_0G_1^{-1}BP_0x + Q_0x = Q_0G_1^{-1}q \end{array} \right\} \\ &\underset{(a)}{\Leftrightarrow} \left\{ \begin{array}{l} R(Dx)' + DG_1^{-1}BP_0x = DG_1^{-1}q \\ Q_0G_1^{-1}BP_0x + Q_0x = Q_0G_1^{-1}q \end{array} \right\} \\ &\underset{(b)}{\Leftrightarrow} \left\{ \begin{array}{l} (Dx)' - R'Dx + DG_1^{-1}BP_0x = DG_1^{-1}q \\ Q_0G_1^{-1}BP_0x + Q_0x = Q_0G_1^{-1}q \end{array} \right\} \\ &\underset{(a)}{\Leftrightarrow} \left\{ \begin{array}{l} (Dx)' = R'(Dx) - DG_1^{-1}BD^-(Dx) + DG_1^{-1}q \\ Q_0x = -Q_0G_1^{-1}BD^-(Dx) + Q_0G_1^{-1}q \end{array} \right\} \end{aligned}$$

Every solution x can therefore be written as

$$\begin{aligned} x &= P_0x + Q_0x = D^-(Dx) + Q_0x = D^-(Dx) - Q_0G_1^{-1}BD^-(Dx) + Q_0G_1^{-1}q \\ &= (I - Q_0G_1^{-1}B)D^-u + Q_0G_1^{-1}q \end{aligned} \quad (3.4)$$

where $u = Dx$ is a solution of the ODE

$$u' = R'u - DG_1^{-1}BD^-u + DG_1^{-1}q. \quad (3.5)$$

Definition 3.3 *The explicit ordinary differential equation (3.5) is called the inherent regular ODE of the index-1 equation (3.1).*

Lemma 3.4 (i) *im D is a (time varying) invariant subspace of (3.5).*

(ii) *(3.5) is independent of the choice of Q_0 .*

Proof:

(i) Because of $\text{im } D = \text{im } R = \ker(I - R)$ multiplication of (3.5) by $I - R$ gives

$$(I - R)u' = (I - R)R'u = -(I - R)'Ru$$

and $v = (I - R)u$ satisfies the ODE $v' = (I - R)'v$.

If there is $t_* \in \mathcal{J}$ such that $u(t_*) = R(t_*)u(t_*) \in \text{im } D(t_*)$, then $v(t_*) = 0$. This means $v(t) = 0$ and thus $u(t) = R(t)u(t)$ for every t .

(ii) Let \hat{Q}_0 be another projector with $\text{im } \hat{Q}_0 = N_0$ and let \hat{P}_0, \hat{D}^- be defined as in (2.16). Then $\hat{G}_1 = G_1(I + Q_0\hat{Q}_0P_0)$ implies $\hat{G}_1^{-1} = (I - Q_0\hat{Q}_0P_0)G_1^{-1}$ and $D\hat{G}_1^{-1} = DG_1^{-1}$. Finally note that

$$\begin{aligned} D\hat{G}_1^{-1}B\hat{D}^- &\stackrel{(a)}{=} D\hat{G}_1^{-1}B\hat{P}_0D^- = D\hat{G}_1^{-1}BD^- - D\hat{G}_1^{-1}B\hat{Q}_0D^- \\ &\stackrel{(d)}{=} DG_1^{-1}BD^- - D\hat{Q}_0D^- = DG_1^{-1}BD^-. \end{aligned} \quad \square$$

The decoupling procedure above and lemma 3.4 enable us to prove existence and uniqueness of solutions for the index 1 DAE (3.1).

Theorem 3.5 *Let (3.1) be a regular index 1 DAE. For each $d \in \text{im } D(t_0)$, $t_0 \in \mathcal{J}$, the initial value problem*

$$A(t)(D(t)x(t))' + B(t)x(t) = q(t), \quad D(t_0)x(t_0) = d \quad (3.6)$$

is uniquely solvable in $C_D^1(\mathcal{J}, \mathbb{R}^m)$.

Proof: There is exactly one solution $u \in C^1(\mathcal{J}, \mathbb{R}^m)$ of the inherent ODE

$$u' = R'u - DG_1^{-1}BD^-u + DG_1^{-1}q$$

satisfying the initial condition $u(t_0) = d$. Lemma 3.4 shows that $u(t) = R(t)u(t)$ for every t . Therefore

$$x = (I - Q_0G_1^{-1}B)D^-u + Q_0G_1^{-1}q \in C_D^1(\mathcal{J}, \mathbb{R}^m)$$

is a solution of (3.6) satisfying $Dx = u$. The decoupling process shows the uniqueness. \square

Note that the initial condition $D(t_0)x(t_0) = d$ for $d \in \text{im } D(t_0)$ can be replaced by $D(t_0)x(t_0) = D(t_0)x^0$, $x^0 \in \mathbb{R}^m$.

3.2 Decoupling of linear index-2 DAEs

We now want to repeat the same argument for linear index 2 differential-algebraic equations. We assume that (3.1) is an index 2 DAE with properly stated leading term. Due to definition 3.1 and lemma 2.13 we have $N_1(t) \oplus S_1(t) = \mathbb{R}^m$. In this section we choose Q_1 to be the canonical projector onto N_1 along S_1 . Lemma 2.13 also implies $Q_1 Q_0 = Q_1 G_2^{-1} B_1 Q_0 = 0$ as required in definition 3.1.

For the sequence (3.2) to make sense we have to assume $DP_1 D^- \in C^1(\mathcal{J}, L(\mathbb{R}^n))$. Then $DQ_1 D^- = -DP_1 D^- + DD^- = -DP_1 D^- + R$ is also smooth. Note that $DQ_1 D^-$ and $DP_1 D^-$ are projector functions.

In addition to (a), (b), (c) and (f) from lemma 3.2 we now have

Lemma 3.6

- (g) $Q_1 = Q_1 G_2^{-1} B_1$,
- (h) $G_2^{-1} AD = P_1 P_0$,
- (i) $G_2^{-1} B = G_2^{-1} B P_0 P_1 + P_1 D^- (DP_1 D^-)' DQ_1 + Q_1 + Q_0$,
- (j) $Q_1 x = Q_1 y \Leftrightarrow DQ_1 x = DQ_1 y$,
- (k) $\Omega \Omega' \Omega = 0$ for every projector function $\Omega \in C^1(\mathcal{J}, L(\mathbb{R}^n))$.

Proof: (g) follows from lemma 2.13, (h) can be proved similar to (e) in lemma 3.2, but (i) is a consequence of $B = B P_0 + B Q_0 = B P_0 P_1 + B P_0 Q_1 + B Q_0$ and

$$G_2 Q_0 = B Q_0, \quad G_2 Q_1 = B_1 Q_1, \quad B P_0 Q_1 = B_1 Q_1 + G_2 P_1 D^- (DP_1 D^-)' DQ_1.$$

To show (j) assume that $DQ_1 z = 0$. Then $Q_1 z \in \ker D = \ker P_0$ and $Q_1 z = Q_1^2 z = Q_1 P_0 Q_1 z = 0$.

Finally $0 = (I - \Omega)\Omega$ implies $0 = (I - \Omega)'\Omega + (I - \Omega)\Omega' = -\Omega'\Omega + (I - \Omega)\Omega'$. \square

In order to decouple (3.1) in the index 2 case we again assume that x is a solution of the DAE. Since G_2 is nonsingular, we find

$$\begin{aligned} A(Dx)' + Bx = q &\Leftrightarrow G_2^{-1} A(Dx)' + G_2^{-1} Bx = G_2^{-1} q & (3.7) \\ &\Leftrightarrow P_1 D^- (Dx)' + G_2^{-1} B P_0 P_1 x + P_1 D^- (DP_1 D^-)' DQ_1 x + Q_1 x + Q_0 x = G_2^{-1} q \end{aligned}$$

using (a), (c), (h) and (i). Due to $I = P_1 + Q_1 = P_0 P_1 + Q_0 P_1 + Q_1$ we can decouple (3.7) by multiplying with $P_0 P_1$, $Q_0 P_1$ and Q_1 respectively. (3.1) is therefore equivalent to the system

$$P_0 P_1 D^- (Dx)' + P_0 P_1 G_2^{-1} B P_0 P_1 x + P_0 P_1 D^- (DP_1 D^-)' DQ_1 x = P_0 P_1 G_2^{-1} q, \quad (3.8a)$$

$$\begin{aligned} Q_0 P_1 D^- (Dx)' + Q_0 P_1 G_2^{-1} B P_0 P_1 x \\ + Q_0 P_1 D^- (DP_1 D^-)' DQ_1 x + Q_0 x = Q_0 P_1 G_2^{-1} q, \end{aligned} \quad (3.8b)$$

$$Q_1 G_2^{-1} B P_0 P_1 x + Q_1 x = Q_1 G_2^{-1} q. \quad (3.8c)$$

With (a) and (f) equation (3.8a) takes the form

$$DP_1D^-(Dx)' + DP_1G_2^{-1}BP_0P_1x + DP_1D^-(DP_1D^-)'DQ_1x = DP_1G_2^{-1}q.$$

Use the product rule of differentiation to find

$$DP_1D^-(Dx)' = (DP_1x)' - (DP_1D^-)'(Dx).$$

On the other hand

$$DP_1D^-(DP_1D^-)'DQ_1 = (DP_1D^-)'DQ_1$$

as $(DP_1D^-)'DP_1P_0Q_1 = 0$, so that (3.8a) is equivalent to

$$(DP_1x)' - (DP_1D^-)'(DP_1x) + DP_1G_2^{-1}BD^-(DP_1x) = DP_1G_2^{-1}q. \quad (3.8a')$$

A similar analysis involving (g), (j) and (k) from lemma 3.6 yields

$$\begin{aligned} -Q_0Q_1D^-(DQ_1x)' + Q_0Q_1D^-(DQ_1D^-)'(DP_1x) \\ + Q_0P_1G_2^{-1}BD^-(DP_1x) + Q_0x = Q_0P_1G_2^{-1}q, \end{aligned} \quad (3.8b')$$

$$DQ_1x = DQ_1G_2^{-1}q. \quad (3.8c')$$

Each solution x of (3.1) can thus be written as

$$\begin{aligned} x &= P_0x + Q_0x = D^-Dx + Q_0x \\ &= D^-(DP_1x + DQ_1x) + Q_0x \\ &= KD^-u - Q_0Q_1D^-(DQ_1D^-)'u + (Q_0P_1 + P_0Q_1)G_2^{-1}q + Q_0Q_1D^-(DQ_1G_2^{-1}q)' \end{aligned} \quad (3.9)$$

where

$$K = I - Q_0P_1G_2^{-1}B$$

and $u = DP_1x$ satisfies the ordinary differential equation

$$u' - (DP_1D^-)'u + DP_1G_2^{-1}BD^-u = DP_1G_2^{-1}q.$$

As in the index 1 case this ODE will be referred to as the inherent regular ODE.

Definition 3.7 *The explicit ordinary differential equation*

$$u' = (DP_1D^-)'u - DP_1G_2^{-1}BD^-u + DP_1G_2^{-1}q \quad (3.10)$$

is called the inherent regular ODE of the index 2 equation (3.1).

For the index 2 case we now prove the lemma corresponding to lemma 3.4 from the previous section.

Lemma 3.8 (i) *im DP_1 is a (time varying) invariant subspace of (3.10).*

(ii) *(3.10) is independent of the choice of Q_0 and thus uniquely determined by the problem data.*

Proof: To prove (i), carry out a similar analysis as in the proof of lemma 3.4 but with R replaced by DP_1D^- . To see (ii) consider another projector \hat{Q}_0 with $\text{im } \hat{Q}_0 = N_0$ and the relation $\hat{G}_1 = G_1(I + Q_0\hat{Q}_0P_0)$. The subspaces $\hat{N}_1 = (I - Q_0\hat{Q}_0P_0)N_1$ and $\hat{S}_1 = S_1$ are given in terms of N_1 and S_1 so that $\hat{Q}_1 = (I + Q_0\hat{Q}_0P_0)Q_1$ is the canonical projector onto \hat{N}_1 along \hat{S}_1 . This implies $D\hat{P}_1\hat{D}^- = DP_1D^-$. Use the representation

$$\hat{G}_2^{-1} = (I + Q_0\hat{P}_0P_1P_0)G_2^{-1}$$

to see that $DP_1G_2^{-1}$ and $DP_1G_2^{-1}BD^-$ are independent of the choice of Q_0 . \square

As in the previous section we are now able to prove existence and uniqueness of solutions for regular index 2 DAEs with properly stated leading terms. We make use of the function space

$$C_{DQ_1G_2^{-1}}^1(\mathfrak{J}, \mathbb{R}^m) = \{ z \in C(\mathfrak{J}, \mathbb{R}^m) \mid DQ_1G_2^{-1}z \in C^1(\mathfrak{J}, \mathbb{R}^n) \}.$$

Theorem 3.9 *Let (3.1) be a regular index 2 DAE with $q \in C_{DQ_1G_2^{-1}}^1(\mathfrak{J}, \mathbb{R}^m)$. For each $d \in \text{im } D(t_0)P_1(t_0)$, $t_0 \in \mathfrak{J}$, the initial value problem*

$$A(t)(D(t)x(t))' + B(t)x(t) = q(t), \quad D(t_0)P_1(t_0)x(t_0) = d \quad (3.11)$$

is uniquely solvable in $C_D^1(\mathfrak{J}, \mathbb{R}^m)$.

Proof: Solve the inherent regular ODE (3.10) with initial value $u(t_0) = d$. Lemma 3.8 yields $u(t) \in \text{im } D(t)P_1(t)$ for every t and

$$x = KD^-u - Q_0Q_1D^-(DQ_1D^-)'u + (Q_0P_1 + P_0Q_1)G_2^{-1}q + Q_0Q_1D^-(DQ_1G_2^{-1}q)'$$

is the desired solution of (3.11). \square

The initial condition $D(t_0)P_1(t_0)x(t_0) = d$ can be replaced by $D(t_0)P_1(t_0)x(t_0) = D(t_0)P_1(t_0)x^0$ for $x^0 \in \mathbb{R}^m$.

3.3 Remarks

In sections 1.3 and 1.4 we presented examples of nonlinear differential-algebraic equations $f((Dx)', x, t) = 0$, where the solution could be expressed as

$$x(t) = D(t)^-u(t) + Q(t)\omega(u(t), t), \quad t \in \mathfrak{J}.$$

u was the solution of

$$u'(t) = R'(t)u(t) + D(t)\omega(u(t), t), \quad u(t_0) = D(t_0)x_0 \quad (3.12)$$

and ω was implicitly defined by

$$F(\omega, u, t) = f(D\omega, D^-u + Q\omega, t) = 0.$$

The ordinary differential equation (3.12) is thus only available theoretically.

In this section we made use of the sequence (3.2) established in the context of the tractability index in order to perform a refined analysis of linear DAEs with properly stated leading terms. We were able to find explicit expressions of (3.12) for these equations with index 1 and 2.

This detailed analysis lead us to results about existence and uniqueness of solutions for DAEs with low index. We were able to figure out precisely what initial conditions are to be posed, namely $D(t_0)x(t_0) = D(t_0)x^0$ and $D(t_0)P_1(t_0)x(t_0) = D(t_0)P_1(t_0)x^0$ in the index 1 and index 2 case respectively.

These initial conditions guarantee that solutions u of the inherent regular ODE (3.5) and (3.10) lie in the corresponding invariant subspace. Let us stress that only those solutions of the regular inherent ODE that lie in the invariant subspace are relevant for the DAE. Even if this subspace varies with t we know the dynamical degree of freedom to be $\text{rank } G_0$ and $\text{rank } G_0 + \text{rank } G_1 - m$ for index 1 and 2 respectively [25].

The results presented can be generalized for arbitrary index μ . The inherent regular ODE for an index μ DAE with properly stated leading term is given in [25]. There it is also proved that the index μ is invariant under linear transformations and refactorizations of the original DAE and the inherent regular ODE remains unchanged.

Finally let us point out that we assumed A , D and B to be continuous only. The required smoothness of the coefficients in the standard formulation

$$Ex' + Fx = q \tag{3.13}$$

was replaced by the requirement on certain subspaces to be spanned by smooth functions. Namely, the projectors R , DP_1D^- and DQ_1D^- are differentiable if DN_1 and DS_1 are spanned by continuously differentiable functions [1].

However, if the DAE

$$A(Dx)' + Bx = q \tag{3.14}$$

is given with smooth coefficients and we orient on C^1 -solutions, then comparisons with concepts for (3.13) can be made via

$$ADx' + (B - AD')x = q.$$

On the other hand, if E has constant rank on \mathcal{J} and $P_E \in C^1(\mathcal{J}, L(\mathbb{R}^m))$ is a projector function onto $\ker E$, we can reformulate (3.13) as

$$E(P_E x)' + (F - EP_E')x = q$$

with a properly stated leading term.

4 Numerical methods for linear DAEs with properly stated leading term

The last part is devoted to studying the application of numerical methods to linear DAEs of index $\mu = 1$ and $\mu = 2$. From the previous section we know that (3.4) and (3.9) are representations of the exact the solution, respectively. In fact, it turns out that (3.4) is just a special cases of (3.9). To see this, observe that for $\mu = 1$ the matrix G_1 is nonsingular so that $Q_1 = 0$, $P_1 = I$ and $G_2 = G_1$. We therefore treat index 1 and index 2 equations simultaneously in this section. We will show how to apply Runge-Kutta methods to DAEs

$$A(t)(D(t)x(t))' + B(t)x(t) = q(t) \quad (4.1)$$

with properly stated leading terms. Results presented here follow the lines of [17, 18, 26]. Runge-Kutta methods for DAEs are also studied in [14].

When using the s -stage Runge-Kutta method

$$\frac{c}{\beta^T} \Big| \frac{\mathcal{A}}{\beta^T}, \quad \mathcal{A} = (\alpha_{ij}) \in L(\mathbb{R}^s), \quad c = \mathcal{A}e, \quad \beta \in \mathbb{R}^s, \quad e = (1, \dots, 1)^T \in \mathbb{R}^s,$$

to solve an ordinary differential equation

$$x'(t) = F(x(t), t) \quad (4.2)$$

numerically with stepsize h , an approximation x_{l-1} to the exact solution $x(t_{l-1})$ is used to calculate the approximation x_l to $x(t_l) = x(t_{l-1} + h)$ via

$$x_l = x_{l-1} + h \sum_{i=1}^s \beta_i X'_{li} \quad (4.3a)$$

where X'_{li} is defined by

$$X'_{li} = F(X_{li}, t_{li}), \quad i = 1, \dots, s, \quad (4.3b)$$

and $t_{li} = t_{l-1} + c_i h$ are intermediate timesteps. The internal stages X_i are given by

$$X_{li} = x_{l-1} + h \sum_{j=1}^s \alpha_{ij} X'_{lj}. \quad (4.3c)$$

Observe that (4.3a) and (4.3c) depend on the method and only (4.3b) depends on the equation (4.2). If the ODE (4.2) is replaced by the DAE

$$f(x'(t), x(t), t) = 0$$

we also replace (4.3b) by

$$f(X'_{li}, X_{li}, t_{li}) = 0, \quad i = 1, \dots, s \quad (4.3b')$$

in the Runge-Kutta scheme.

The matrix $\frac{\partial f}{\partial x'}$ is singular. Therefore some components of the increments X'_{li} need to be calculated from (4.3c) as seen in the following trivial example.

Example 4.1 If $f(x', x, t) = x - q(t)$, then $x(t) = q(t)$. The numerical method (4.3a), (4.3b'), (4.3c) now reads

$$x_l = x_{l-1} + h \sum_{i=1}^s \beta_i X'_{li}, \quad q(t_{li}) = X_{li} = x_{l-1} + h \sum_{j=1}^s \alpha_{ij} X'_{lj}.$$

This system can be solved if and only if \mathcal{A} is nonsingular. \square

In the following, we always assume \mathcal{A} to be nonsingular. This leads to an expression of X'_{li} in terms of X_{lj} .

Lemma 4.2 Let $\mathcal{A} = (\alpha_{ij})$ be nonsingular and $\mathcal{A}^{-1} = (\tilde{\alpha}_{ij})$. Then

$$X_{li} = x_{l-1} + h \sum_{j=1}^s \alpha_{ij} X'_{lj}, \quad i=1, \dots, s \quad \Leftrightarrow \quad X'_{li} = \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (X_{lj} - x_{l-1}), \quad i=1, \dots, s.$$

Proof: If \otimes denotes the Kronecker product and $e_m = (1, \dots, 1)^T \in \mathbb{R}^m$ then

$$\begin{pmatrix} X_{l1} \\ \vdots \\ X_{ls} \end{pmatrix} = e_m \otimes x_{l-1} + h(\mathcal{A} \otimes I_m) \begin{pmatrix} X'_{l1} \\ \vdots \\ X'_{ls} \end{pmatrix} \Leftrightarrow \begin{pmatrix} X'_{l1} \\ \vdots \\ X'_{ls} \end{pmatrix} = \frac{1}{h} (\mathcal{A}^{-1} \otimes I_m) \left[\begin{pmatrix} X_{l1} \\ \vdots \\ X_{ls} \end{pmatrix} - e_m \otimes x_{l-1} \right] \quad \square$$

Now consider the linear DAE (4.1) with continuous matrix functions

$$A(t) \in L(\mathbb{R}^n, \mathbb{R}^m), \quad D(t) \in L(\mathbb{R}^m, \mathbb{R}^n), \quad B(t) \in L(\mathbb{R}^m, \mathbb{R}^m).$$

and a properly stated leading term.

When applying the numerical scheme (4.3a),(4.3b'),(4.3c) we don't want to lose the additional information provided by the properly stated leading term. According to lemma 4.2 we therefore replace (4.3c) by

$$[DX]_{li}' = \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (D_{lj} X_{lj} - D_{l-1} x_{l-1}) \quad (4.3c')$$

and solve the system

$$A_{li} [DX]_{li}' + B_{li} X_{li} = q_{li}, \quad i = 1, \dots, s \quad (4.3b'')$$

for X_{li} . Here we write $D_{l-1} = D(t_{l-1})$, $D_{li} = D(t_{li})$, $A_{li} = A(t_{li})$ and so on. Using this ansatz the output value

$$x_l = x_{l-1} + h \sum_{i=1}^s \beta_i \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (X_{lj} - x_{l-1}) = (1 - \beta^T \mathcal{A}^{-1} e) x_{l-1} + \sum_{i=1}^s \sum_{j=1}^s \beta_i \tilde{\alpha}_{ij} X_{lj}$$

is computed. For RadauIIA methods this expression simplifies considerably.

Definition 4.3 The s -stage RadauIIA method is uniquely determined by requiring $C(s)$, $D(s)$, $c_s = 1$ and choosing c_1, \dots, c_{s-1} to be the zeros of the Gauss-Legendre polynomial p_s .

For the conditions $C(s)$, $D(s)$ see [3]. The Gauss-Legendre polynomial p_s is orthogonal to every polynomial of degree less than s . RadauIIA methods are A - and L -stable and have order $p = 2s - 1$. The last row of \mathcal{A} coincides with β^T [3, 15].

Lemma 4.4 For the s -stage RadauIIA method $1 - \beta^T \mathcal{A}^{-1} e = 0$ holds and the output value computed by (4.3a), (4.3b''), (4.3c), (4.3c') is given by the last stage X_{l_s} .

Proof: $1 - \beta^T \mathcal{A}^{-1} e = 1 - Z_s(\mathcal{A}) \mathcal{A}^{-1} e = 1 - (0, \dots, 0, 1) e = 0$ and

$$x_l = (1 - \beta^T \mathcal{A}^{-1} e) x_{l-1} + \sum_{i=1}^s \sum_{j=1}^s \beta_i \tilde{\alpha}_{ij} X_{lj} = ((0, \dots, 0, 1) \otimes I_m) \begin{pmatrix} X_{l1} \\ \vdots \\ X_{ls} \end{pmatrix} = X_{l_s}. \quad \square$$

To summarize these results we present the following algorithm for solving the DAE (4.1) using RadauIIA methods.

Algorithm 4.5 Given an approximation x_{l-1} to the exact solution $x(t_{l-1})$ and a stepsize h , solve

$$A_{li} [DX]_{li}' + B_{li} X_{li} = q_{li}, \quad i = 1, \dots, s \quad (4.3b'')$$

for X_{li} where $[DX]_{li}'$ is given by

$$[DX]_{li}' = \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (D_{lj} X_{lj} - D_{l-1} x_{l-1}). \quad (4.3c')$$

Return the output value $x_l = X_{l_s}$ as an approximation to $x(t_l) = x(t_{l-1} + h)$.

The exact solution x of (4.1) satisfies

$$x(t) \in \mathcal{M}_0(t) = \{z \in \mathbb{R}^m \mid B(t)z - q(t) \in \text{im}(A(t)D(t))\} \quad \forall t.$$

Since $X_{li} \in \mathcal{M}_0(t_{li})$ for every i and $c_s = 1$ we have

$$x_l = X_{l_s} \in \mathcal{M}_0(t_{l_s}) = \mathcal{M}_0(t_l)$$

for every RadauIIA method. Thus the RadauIIA approximation satisfies the algebraic constraint and RadauIIA methods are especially suited for solving DAEs [14, 15].

4.1 Decoupling of the discretized equation

Algorithm 4.5 replaces the DAE

$$A(Dx)' + Bx = q \quad (4.1)$$

by the discretized problem

$$A_{li} [DX]_{li}' + B_{li} X_{li} = q_{li}, \quad i = 1, \dots, s. \quad (4.4)$$

As seen in section 3.2, the analytic solution x of index 1 and index 2 equations (4.1) can be represented as

$$x = KD^-u - Q_0 Q_1 D^- (DQ_1 D^-)' u + (Q_0 P_1 + P_0 Q_1) G_2^{-1} q + Q_0 Q_1 D^- (DQ_1 G_2^{-1} q)' \quad (4.5)$$

where $K = I - Q_0 P_1 G_2^{-1} B$ and the component $u = DP_1 x$ satisfies the inherent regular ordinary differential equation

$$u' - (DP_1 D^-)' u + DP_1 G_2^{-1} B D^- u = DP_1 G_2^{-1} q. \quad (4.6)$$

If we applied the Runge-Kutta method directly to the inherent regular ODE, due to lemma 4.2 we would obtain

$$\frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (U_{lj} - u_{l-1}) - (DP_1 D^-)'_{li} U_{li} + (DP_1 G_2^{-1} B D^-) U_{li} = (DP_1 G_2^{-1} q)_{li} \quad (4.7)$$

for $i = 1, \dots, s$. Our aim is to show that the Runge-Kutta method, when applied to (4.1), behaves as if it was integrating the inherent regular ODE (4.6).

We start by repeating the decoupling procedure from section 3.2 for the discretized equation (4.4). Doing so (4.4) is found to be equivalent to the system

$$\left. \begin{aligned} (DP_1 D^-)_{li} [DX]_{li}' + (DP_1 G_2^{-1} B P_0 P_1)_{li} X_{li} \\ + (DP_1 D^-)'_{li} D_{li} Q_{1,li} X_{li} &= (DP_1 G_2^{-1} q)_{li} \\ -(Q_0 Q_1 D^-)_{li} [DX]_{li}' + (Q_0 P_1 G_2^{-1} B P_0 P_1)_{li} X_{li} \\ + (Q_0 P_1 D^-)'_{li} (DP_1 D^-)'_{li} D_{li} Q_{1,li} X_{li} + Q_{0,li} X_{li} &= (Q_0 P_1 G_2^{-1} q)_{li} \\ D_{li} Q_{1,li} X_{li} &= (DQ_1 G_2^{-1} q)_{li} \end{aligned} \right\} \quad (4.8)$$

for $i = 1, \dots, s$. The decoupled system (4.8) immediately implies the convergence of RadauIIA methods applied to (4.1) on compact intervals \mathfrak{J} if the stepsize h tends to zero [18]. A careful analysis of (4.8) also leads to the main result in this section.

Theorem 4.6 *Let (4.1) be an index μ equation, $\mu \in \{1, 2\}$. Let the subspaces $D(\cdot)S_1(\cdot)$ and $D(\cdot)N_1(\cdot)$ be constant. Then the difference between the exact solution and the solution obtained by using a RadauIIA method can be written as*

$$\begin{aligned} x(t_l) - x_l &= K_l D_l^- (u(t_l) - u_l) \\ &+ (Q_0 Q_1 D^-)_{li} \left\{ (DQ_1 G_2^{-1} q)'_l - \frac{1}{h} \sum_{j=0}^k \tilde{\alpha}_{sj} \left((DQ_1 G_2^{-1} q)_{lj} - (DQ_1 G_2^{-1} q)_{l-1} \right) \right\}. \end{aligned}$$

Here u_l is exactly the RadauIIA approximation to the solution $u(t_l)$ of the inherent regular ODE (4.6).

Note that $\frac{1}{h} \sum_{j=0}^k \tilde{\alpha}_{sj} \left((DQ_1 G_2^{-1} q)_{lj} - (DQ_1 G_2^{-1} q)_{l-1} \right)$ is exactly the Runge-Kutta approximation to $(DQ_1 G_2^{-1} q)'_l$. The proof of theorem 4.6 will use the following lemma.

Lemma 4.7 *$DP_1 D^-$ and $DQ_1 D^-$ are projector functions satisfying*

$$(i) \quad DS_1 = \text{im } DP_1 = \text{im } DP_1 D^-, \quad DN_1 = \text{im } DQ_1 = \text{im } DQ_1 D^-.$$

If the subspaces DS_1 and DN_1 are constant, so that there are constant projectors V, W onto DS_1 and DN_1 respectively, then the following relations hold:

$$(ii) \quad DP_1 D^- V = V, \quad DP_1 D^- W = 0, \quad DQ_1 D^- W = W, \quad DQ_1 D^- V = 0,$$

$$(iii) \quad (DP_1 D^-)' V = 0, \quad (DP_1 D^-)' W = 0, \quad (DQ_1 D^-)' W = 0, \quad (DQ_1 D^-)' V = 0.$$

Proof: DP_1D^- and DQ_1D^- are projector functions due to lemma 3.2 and 3.6. The same lemmas imply (i), so that $DP_1D^-V = V$ and $DQ_1D^-W = W$ hold as well. These relations together with (i) show (ii). Finally use (ii) to prove (iii) by noting that V and W are constant projectors and therefore do not depend on t . \square

Proof of theorem 4.6: The proof will be divided into four parts. In ① we analyze $(DP_1D^-)_{li}[DX]_{li}'$ and $(Q_0Q_1D^-)_{li}[DX]_{li}'$, so that we can find a representation of the numerical solution in part ②. This representation will depend on $U_{ls} = D_{ls}P_{1,ls}X_{ls}$. In ③ we show that $u_l = U_{ls}$ is exactly the RadauIIA solution of the inherent regular ODE. The poof will be completed by comparing the analytic and the numeric solution in part ④.

① Analyze $(DP_1D^-)_{li}[Dx]_{li}'$ and $(Q_0Q_1D^-)_{li}[Dx]_{li}'$

Write $U_{li} = D_{li}P_{1,li}X_{li}$ and $u_{l-1} = D_{l-1}P_{1,l-1}x_{l-1}$. Then

$$\begin{aligned} (DP_1D^-)_{li}[Dx]_{li}' &= \frac{1}{h}(DP_1D^-)_{li} \sum_{j=1}^s \tilde{\alpha}_{ij}(D_{lj}X_{lj} - D_{l-1}x_{l-1}) \\ &= \frac{1}{h}(DP_1D^-)_{li} \sum_{j=1}^s \tilde{\alpha}_{ij} \left(U_{lj} + D_{lj}Q_{1,li}X_{lj} - u_{l-1} - (DQ_1x)_{l-1} \right). \end{aligned}$$

Use lemma 4.7 to see that

$$(DP_1D^-)_{li}(U_{lj} - u_{l-1}) = (DP_1D^-)_{li}V(U_{lj} - u_{l-1}) = V(U_{lj} - u_{l-1}) = U_{lj} - u_{l-1}$$

and

$$(DP_1D^-)_{li}(D_{lj}Q_{1,li}X_{lj} - (DQ_1x)_{l-1}) = (DP_1D^-)_{li}W(D_{lj}Q_{1,li}X_{lj} - (DQ_1x)_{l-1}) = 0.$$

We arrive at

$$(DP_1D^-)_{li}[Dx]_{li}' = \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (U_{lj} - u_{l-1}).$$

Similarly, lemma 4.7 implies

$$(DQ_1D^-)_{li}[Dx]_{li}' = \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} \left(D_{lj}Q_{1,li}X_{lj} - (DQ_1x)_{l-1} \right).$$

Because of

$$(Q_0Q_1D^-)_{li} = (Q_0(Q_1P_0Q_1)D^-)_{li} = ((Q_0Q_1D^-)(DQ_1D^-))_{li},$$

it follows that

$$(Q_0Q_1D^-)_{li}[DX]_{li}' = \frac{1}{h}(Q_0Q_1D^-)_{li} \sum_{j=1}^s \tilde{\alpha}_{ij} \left(D_{lj}Q_{1,li}X_{lj} - (DQ_1x)_{l-1} \right).$$

② Obtain a representation of the numerical solution x_l

The discretized system (4.8) now reads

$$\left. \begin{aligned} \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (U_{lj} - u_{l-1}) + (DP_1 G_2^{-1} B D^-)_{li} U_{li} \\ + (DP_1 D^-)'_{li} D_{li} Q_{1,li} X_{li} = (DP_1 G_2^{-1} q)_{li} \\ - \frac{1}{h} (Q_0 Q_1 D^-)_{li} \sum_{j=1}^s \tilde{\alpha}_{ij} (D_{lj} Q_{1,li} X_{lj} - (DQ_1 x)_{l-1}) \\ + (Q_0 P_1 G_2^{-1} B D^-)_{li} U_{li} \\ + (Q_0 P_1 D^-)_{li} (DP_1 D^-)'_{li} D_{li} Q_{1,li} X_{li} + Q_{0,li} X_{li} = (Q_0 P_1 G_2^{-1} q)_{li} \\ D_{li} Q_{1,li} X_{li} = (DQ_1 G_2^{-1} q)_{li} \end{aligned} \right\}$$

but due to lemma 4.7 this reduces to

$$\left. \begin{aligned} \frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (U_{lj} - u_{l-1}) + (DP_1 G_2^{-1} B D^-)_{li} U_{li} = (DP_1 G_2^{-1} q)_{li} \\ - \frac{1}{h} (Q_0 Q_1 D^-)_{li} \sum_{j=1}^s \tilde{\alpha}_{ij} (D_{lj} Q_{1,li} X_{lj} - (DQ_1 x)_{l-1}) \\ + (Q_0 P_1 G_2^{-1} B D^-)_{li} U_{li} + Q_{0,li} X_{li} = (Q_0 P_1 G_2^{-1} q)_{li} \\ D_{li} Q_{1,li} X_{li} = (DQ_1 G_2^{-1} q)_{li} \end{aligned} \right\}$$

The numerical solution can thus be written as

$$\begin{aligned} x_l = X_{ls} &= P_{0,ls} X_{ls} + Q_{0,ls} X_{ls} = D_{ls}^- (D_{ls} P_{1,ls} X_{ls} + D_{ls} Q_{1,ls} X_{ls}) + Q_{0,ls} X_{ls} \\ &= (I - (Q_0 P_1 G_2^{-1} B)_{ls}) D_{ls}^- U_{ls} + (P_0 Q_1 + Q_0 P_1)_l (G_2^{-1} q)_l \quad (4.9) \\ &\quad + \frac{1}{h} (Q_0 Q_1 D^-)_l \sum_{j=1}^s \tilde{\alpha}_{sj} \left((DQ_1 G_2^{-1} q)_{lj} - (DQ_1 G_2^{-1} q)_{l-1} \right). \end{aligned}$$

The stage approximations U_{lj} satisfy the recursion

$$\frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{ij} (U_{lj} - u_{l-1}) + (DP_1 G_2^{-1} B D^-)_{li} U_{li} = (DP_1 G_2^{-1} q)_{li}. \quad (4.10)$$

③ (4.10) is the RadauIIA method applied to the inherent regular ODE

Again, lemma 4.7 implies

$$(DP_1 D^-)'_{li} U_{li} = (DP_1 D^-)'_{li} V U_{li} = 0$$

in (4.7). This shows that (4.10) and (4.7) coincide. Therefore, and due to $c_s = 1$, $u_l = U_{ls}$ is exactly the Runge-Kutta solution of the inherent regular ODE (4.7).

④ Compare the analytic and the numeric solution

Use Lemma 4.7 to see, that in (4.5)

$$(DQ_1 D^-)'_l u(t_l) = (DQ_1 D^-)'_l V u(t_l) = 0.$$

Now the assertion follows by comparing (4.5) and (4.9). \square

Theorem 4.6 is the central tool in analyzing the behaviour of RadauIIA methods when applied to DAEs (4.1). In the case of index $\mu = 1$ theorem 4.6 shows that discretization and the decoupling procedure commute.

Corollary 4.8 *Let the DAE (4.1) be of index 1. Assume that $\text{im } D(t)$ is constant. Then we have for any RadauIIA method*

$$x(t_l) - x_l = K_l D_l^- (u(t_l) - u_l), \quad K = I - Q_0 G_1^{-1} B.$$

Proof: If the index is 1, we have $Q_1 = 0$ and $P_1 = I$. Thus $N_1 = \{0\}$ and $S_1 = \mathbb{R}^n$. Since $\text{im } D(t)$ is constant, the subspaces DS_1 and DN_1 are constant as well. We can therefore apply theorem 4.6. \square

Due to corollary 4.8 the following diagram commutes for index 1 equations with constant $\text{im } D$.

$$\begin{array}{ccc}
\boxed{\begin{array}{c} A(Dx)' + Bx = q \\ (4.1) \end{array}} & \xrightarrow[\text{discretization}]{\text{RadauIIA}} & \boxed{\begin{array}{c} A_{li}[DX]'_{li} + B_{li}X_{li} = q_{li} \\ (4.4) \end{array}} \\
\downarrow \text{decoupling} & & \downarrow \text{decoupling} \\
\boxed{\begin{array}{c} x = KD^-u + Q_0G_1^{-1}q \\ u' + DG_1^{-1}BD^-u = DG_1^{-1}q \end{array}} & \xrightarrow[\text{discretization}]{\text{RadauIIA}} & \boxed{\begin{array}{c} x_l = K_l D_l^- u_l + Q_{0l} G_{1l}^{-1} q_l \\ \frac{1}{h} \sum_{j=0}^k \tilde{\alpha}_{ij} (U_{lj} - u_{l-1}) + D_l G_{1l}^{-1} B_l D_l^- u_l = D_l G_{1l}^{-1} q_l \end{array}}
\end{array}$$

If the index is 2, we cannot expect the corresponding diagram to commute. However, the term

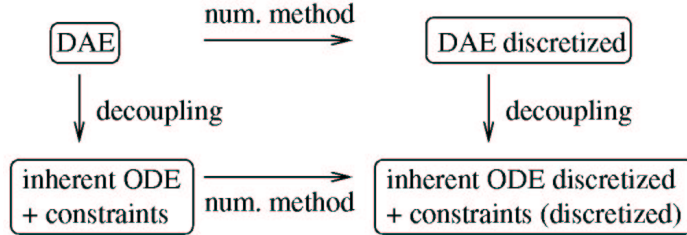
$$\frac{1}{h} \sum_{j=1}^s \tilde{\alpha}_{sj} \left((DQ_1 G_2^{-1} q)_{lj} - (DQ_1 G_2^{-1} q)_{l-1} \right) = [DQ_1 G_2^{-1} q]'_{t_l}$$

appearing in theorem 4.6 is exactly the RadauIIA approximation to $(DQ_1 G_2^{-1} q)'_l$ (lemma 4.2) so that

$$x(t_l) - x_l = K_l D_l^- (u(t_l) - u_l) + Q_{0l} Q_{1l} D_l^- \left\{ (DQ_1 G_2^{-1} q)'_l - [DQ_1 G_2^{-1} q]'_{t_l} \right\}.$$

In this sense we have the following statement:

When applying a RadauIIA method to problems of index $\mu \in \{1, 2\}$ with constant subspaces DS_1 and DN_1 , then discretization and decoupling commute.



Definition 4.9 *The DAE (4.1) of index $\mu \in \{1, 2\}$ is said to be numerically qualified, if*

- $\mu = 1$ and $\text{im } D$ is constant,
- $\mu = 2$ and DS_1, DN_1 are constant.

The commutativity of discretization and the decoupling process is the desired property for DAEs since it guarantees a good behaviour of the numerical method. Even though the numerical method is applied to the DAE directly, it behaves as if it was integrating the regular inherent ODE (4.6). In this case results concerning convergence on compact intervals \mathfrak{J} hold automatically. The RadauIIA method applied to a numerically qualified DAEs is convergent with the same order as for ODEs. Results obtained for ODEs concerning the reflexion of qualitative behaviour by the numerical solution can be transferred directly to DAEs using theorem 4.6. More information about stability preserving integration of index 1 and 2 DAEs can be found in [17, 18].

However, the representation (4.9) shows that the Runge-Kutta scheme is weakly unstable when applied to index 2 DAEs. This is due to the inherent differentiation and becomes apparent for small stepsizes h .

We focused on the application of RadauIIA methods. This restriction is not necessary. All results presented here can be proved in a similar way for BDF methods. The application of general linear methods to DAEs is currently being studied.

4.2 A numerical example

Consider the index 2 example due to Gear and Petzold [12].

$$\begin{aligned} \begin{pmatrix} 0 & 0 \\ 1 & \eta t \end{pmatrix} x'(t) + \begin{pmatrix} 1 & \eta t \\ 0 & 1 + \eta \end{pmatrix} x(t) &= \begin{pmatrix} e^{-t} \\ 0 \end{pmatrix}, \quad \eta \in \mathbb{R} \text{ constant} \\ \Leftrightarrow \begin{cases} x_1(t) + \eta t x_2(t) = e^{-t} \\ x_1'(t) + \eta t x_2'(t) + (1 + \eta) x_2(t) = 0 \end{cases} \end{aligned} \quad (4.11)$$

In [12] it is shown that the BDF method fails completely for $\eta = -1$ and is exponentially unstable for all other parameter values $-1 < \eta < -0.5$. In [14] (4.11) is said to pose difficulties to every numerical method.

Numerical results are given in figure 4.1. (4.11) was solved on the interval $[0, 3]$ using the implicit Euler method, the BDF₂-formula and the RadauIIA method with two stages. The step-size used was $h = 10^{-1.5}$. The exact solution is given by $x_1(t) = (1 - \eta t)e^{-t}$ and $x_2(t) = e^{-t}$, so that $x^0 = (1, 1)^T$ is a consistent initial value. All numerical methods used fail even for moderate values of η due to the exponential instability.

Consider the following reformulation

$$\begin{pmatrix} 0 \\ 1 \end{pmatrix} ((1 \quad \eta t) x(t))' + \begin{pmatrix} 1 & \eta t \\ 0 & 1 \end{pmatrix} x(t) = \begin{pmatrix} e^{-t} \\ 0 \end{pmatrix}. \quad (4.12)$$

(4.12) now has a properly stated leading term and $DN_1 = \mathbb{R}$, $DS_1 = \{0\}$ show that the reformulated problem is numerically qualified. We therefore know discretization and the decoupling process commute. This means that solving the reformulated problem yields the correct numerical results as figure 4.2 shows.

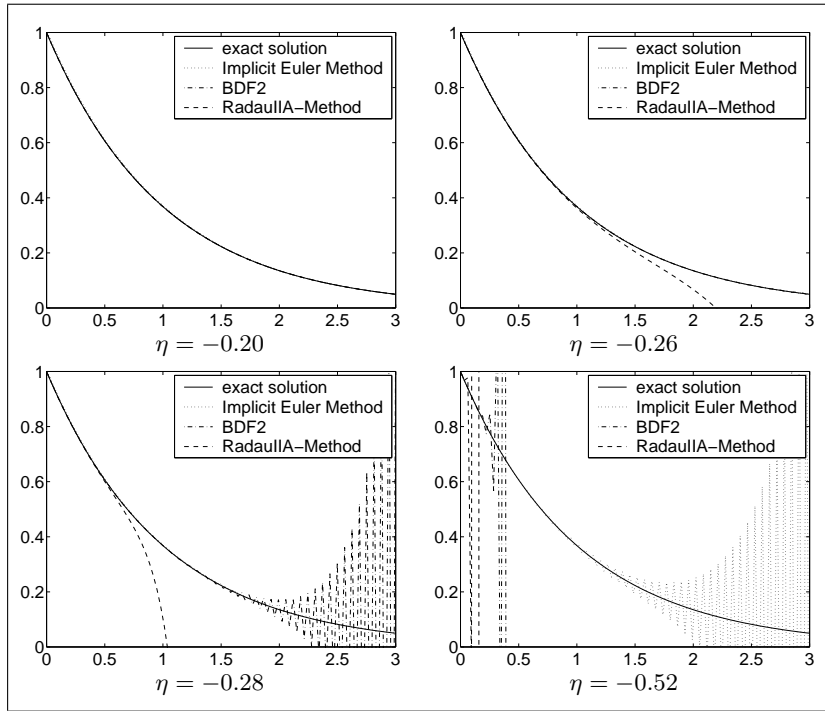


Figure 4.1: Numerical solutions (2. component) of (4.11).

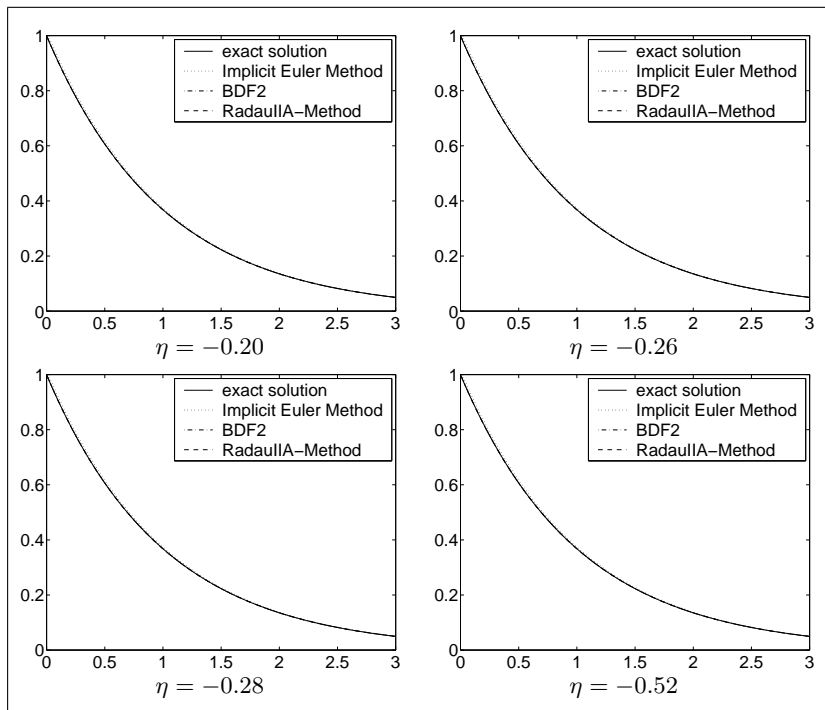


Figure 4.2: Numerical solutions (2. component) of (4.12).

Acknowledgements

First of all I thank Prof John Butcher for inviting me to the University of Auckland. It was an honour for me to be part of his group. I learned a lot about mathematics (but not only about mathematics).

This report would not have been possible without the support of Prof Roswitha März. She and her group at the Humboldt Universität zu Berlin taught me all I know about differential-algebraic equations. I look forward to continuing my work with them.

I greatly appreciate the help of Nicolette Moir and John Rugis who read this report and made very useful comments. The whole Numerical Analysis group at the University of Auckland was a very kind audience and I thank them for listening to me when I was presenting this material.

In terms of financial assistance I'd like to express my gratitude to the Studienstiftung des deutschen Volkes for supporting my stay here in Auckland.

Finally I thank all my friends at the Carlton House and in Grafton Hall for turning my stay here into a wonderful time ... but finally I thank you, Sabine, for letting me go for nine months but being always close to me at the same time.

References

- [1] Balla, K., März, R.: *A unified approach to linear differential algebraic equations and their adjoint equations*, Journal for Analysis and its Applications, vol. 21, no.3, pp. 783-802 (2002)
- [2] Brenan, K.E., Campbell, S.L., Petzold, L.: *Numerical solution of initial-value problems in differential-algebraic equations*, North-Holland, New York (1989)
- [3] Butcher, J.C.: *Numerical Methods for Ordinary Differential Equations*, John Wiley & Sons (2003)
- [4] Campbell, S.L., Gear, C.W.: *The index of general nonlinear DAEs*, Numer. Math., vol. 72, pp. 173-196 (1995)
- [5] CBS-reaction-meeting Köln. Handouts, May 1993. Br/ARLO-CRC.
- [6] Centrum voor Wiskunde en Informatica (CWI), Test Set for IVP Solvers, <http://ftp.cwi.nl/IVPtestset/descrip.htm>
- [7] Estévez Schwarz, D., Tischendorf, C.: *Structural analysis of electric circuits and consequences for MNA*, International Journal of Circuit Theory and Applications, 28, 131-162 (2000)
- [8] Feldmann, U., Günther, M.: *CAD-based electric-circuit modeling in industry – I. Mathematical structure and index of network equations*, Survey on Mathematics for Industry, Springer (1999)
- [9] Gantmacher, F.R.: *The theory of matrices*, New York, Chelsea Pub. Co. (1959)
- [10] Gear, C.W.: *Differential-algebraic equations, indices and integral algebraic equations*, SIAM Journal on Numerical Analysis, vol. 27 (1990)
- [11] Gear, C.W., Petzold, L.: *Differential/algebraic systems and matrix pencils*, Matrix Pencils, B. Kagstrom & A. Ruhe (eds.), Lecture Notes in Mathematics 973, Springer Verlag, pp.75-89 (1983)
- [12] Gear, C.W., Petzold, L.: *ODE methods for the solution of differential/algebraic systems*, SIAM Journal on Numerical Analysis, vol. 21, pp. 716-728 (1984)
- [13] Griepentrog, E., März, R.: *Differential-algebraic equations and their numerical treatment*, Teubner, Leipzig (1986)
- [14] Hairer, E., Lubich, C., Roche, M.: *The Numerical Solution of Differential-Algebraic Systems by Runge-Kutta Methods*, Lecture Notes in Mathematics 1409, Springer, Berlin Heidelberg (1989)
- [15] Hairer, E., Wanner, G.: *Solving ordinary differential equations II: stiff and differential algebraic problems.*, Springer, Berlin Heidelberg New York Tokyo (1991)
- [16] Higuera, I., März, R.: *Differential algebraic equations with properly stated leading terms*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 00-20 (2000), to appear in Computers and Mathematics with Applications

- [17] Higuera, I., März, R., Tischendorf, C.: *Stability preserving integration of index-1 DAEs*, Applied Numerical Mathematics 45, pp. 175-200 (2003)
- [18] Higuera, I., März, R., Tischendorf, C.: *Stability preserving integration of index-2 DAEs*, Applied Numerical Mathematics 45, pp. 201-229 (2003)
- [19] Kronecker, L.: *Algebraische Reduktion der Scharen bilinearer Formen*, Akademie der Wissenschaften Berlin, Werke vol. III, pp. 141-155 (1890)
- [20] Kuh, E.S., Desoer, C.A.: *Basic Circuit Theory*, McGraw-Hill Book Company (1969)
- [21] Kunkel, P., Mehrmann, V.: *Analysis und Numerik linearer differential-algebraischer Gleichungen*, Technische Universität Chemnitz, Fakultät für Mathematik, Preprint SPG 94-27 (1994)
- [22] Kunkel, P., Mehrmann, V.: *Canonical forms for linear differential algebraic equations with variable coefficients*, Journal of Computational and Applied Mathematics, vol. 56, pp. 225-251 (1995)
- [23] Lamour, R., März, R., Tischendorf, C.: *PDAEs and Furter Mixed Systems as Abstract Differential Algebraic Systems*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 01-11 (2001)
- [24] Lamour, R.: *Index determination for DAEs*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 01-19 (2001)
- [25] März, R.: *The index of linear differential algebraic equations with properly stated leading terms*, Results in Mathematics 42, 308-338 (2002)
- [26] März, R.: *Differential algebraic systems anew*, Applied Numerical Mathematics 42, 315-335 (2002)
- [27] März, R.: *Characterizing differential algebraic equations without the use of derivative arrays*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 02-08 (2002)
- [28] März, R.: *Solvability of linear differential algebraic equations with properly stated leading terms*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 02-12 (2002)
- [29] März, R.: *Differential Algebraic Systems with Properly Stated Leading Term and MNA Equations*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 02-13 (2002)
- [30] Petzold, L.: *Differential algebraic equations are not ODEs*, SIAM Journal on Scientific Computing 3, pp. 367-384 (1982)
- [31] Rabier, P., Rheinboldt, W.: *A general existence and uniqueness theory for implicit differential-algebraic equations*, Differential and Integral Equations, 4(3), pp. 563-582 (1991)
- [32] Rabier, P., Rheinboldt, W.: *Theoretical and Numerical Analysis of Differential-Algebraic Equations*, Handboob of Numerical Analysis, vol. VIII, edited by P.G.Ciarlet and J.L.Lions, Elsevier Science B.V. (2002)

- [33] Reich, S.: *Beitrag zur Theorie der Algebrodifferentialgleichungen*, PhD Thesis, Technische Universität Dresden (1990)
- [34] Rentrop, P., Roche, M., Steinebach, G.: *The application of Rosenbrock-Wanner type methods with stepsize control in differential-algebraic equations*, Numerische Mathematik, 55:545-563 (1989)
- [35] Schulz, S.: *Ein PDAE-Netzwerkmodell als abstraktes differential-algebraisches System*, Humboldt-Universität zu Berlin, Institut für Mathematik, Diplomarbeit (2002)
- [36] Schumilina, I.: *Index 3 Algebro-Differentialgleichungen mit proper formuliertem Hauptterm*, Humboldt-Universität Berlin, Institut für Mathematik, Preprint 01-20 (2001)
- [37] Tischendorf, C.: *Topological index calculation of differential-algebraic equations in circuit simulation*, Surveys on Mathematics for Industry 8, Springer (1999)
- [38] Tischendorf, C.: *Numerical Analysis of Differential-Algebraic Equations*, Humboldt-Universität zu Berlin, <http://www.math.hu-berlin.de/~caren/dae-vorl.pdf>, (1999)
- [39] Tischendorf, C.: *Numerische Simulation elektrischer Netzwerke*, Humboldt-Universität zu Berlin, Institut für Mathematik, <http://www.math.hu-berlin.de/~caren/schaltungen.pdf>, (2000)
- [40] Tischendorf, C.: *Modeling Circuit Systems Coupled with Distributed Semiconductor Equations*, Humboldt-Universität zu Berlin, Institut für Mathematik, Preprint 03-2 (2003)
- [41] Zielke, G.: *Motivation und Darstellung von verallgemeinerten Matrixinversen*, Beiträge zur Numerischen Mathematik 7, pp. 177-218 (1979)