

HIGHER-INDEX DIFFERENTIAL-ALGEBRAIC EQUATIONS: ANALYSIS AND NUMERICAL TREATMENT

ROSWITHA MÄRZ

*Section of Mathematics, Humboldt University,
Berlin, German Democratic Republic*

1. Introduction

We are interested in singular ordinary differential equations (ODE's)

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

where the partial Jacobian $f'_y(y, x, t)$ is everywhere singular but has constant rank. Moreover, we suppose the nullspace of $f'_y(y, x, t)$ to be independent of (y, x) . Those uniformly singular ODE's are called *differential-algebraic equations* (DAE's). They originate from actual applications in different areas. In particular, dynamical systems subjected to constraints are described by DAE's (cf. [12]).

In the present paper, so-called index- k -tractable DAE's, $k \in \{2, 3\}$ ("higher index" DAE's), are considered. In Section 2, solvability statements for linear index- k -tractable DAE's are given. The linear map representing a linear initial value problem with appropriately formulated initial conditions is injective but has no bounded inverse in its natural setting.

Those problems are called *ill-posed*. In particular, even the constant-stepsize BDF applied to a linear constant-coefficient DAE of this kind becomes unstable (cf. [8]).

On the other hand, the BDF is reported to work well also for some higher index DAE's when using a careful special error control ([2], [6], [12], [19]). Up to now, for linear constant coefficient DAE's and for certain semi-explicit nonlinear index-2 or index-3 DAE's, the BDF is also proved to be convergent ([2], [5], [6], [12]). Clearly, this does not contradict the instability since the Theorem of Lax does not apply to ill-posed problems. Nevertheless, some confusion arose out of these facts.

In Section 3 and 5 the variable-order variable-step BDF applied to the quasilinear index-2-tractable DAE

$$A(t)x'(t) + q(x(t), t) = 0$$

is proved to have a weak instability only. Instead of a uniform stability bound S for stability, we now obtain a bound $\underline{h}^{-1}S$.

Thereby, $A(t)$ is assumed to have a constant nullspace, and a certain subspace is supposed to be smooth. Further, using detailed error estimations, the BDF is shown to be convergent for (1.2) and, moreover, to have the same order as in case of regular ODE's. Firstly, these facts are proved for linear DAE's (Section 3). Then, by commuting linearization and discretization, the mentioned results are obtained also for nonlinear DAE's.

The weak instability causes the BDF to be very sensitive with respect to the error control in the index-2 case. But, fortunately, implementations using an appropriate control of the defects in the nonlinear equations to be solved per step are reported to work reasonably ([6], [12]). However, the situation for index-3 DAE's becomes worse. At present, the best way of the practical treatment of an index-3 DAE seems to turn to an equivalent index-2 DAE by a so-called reduction step (which is an analytical technique) and then apply a BDF.

An alternative way to treat higher-index DAE's numerically appears in outlines by regularization methods and further special methods for solving ill-posed problems. Investigations in this concern are started recently. Section 4 informs on some interesting results.

2. Solvability of linear initial value problems

Let \mathcal{N} denote the set of all ordered pairs $\{A, B\}$ of continuous matrix functions $A, B: [t_0, T] \rightarrow L(\mathbf{R}^m)$ the first of which has a smooth nullspace, i.e.

$$(2.1) \quad N(t) := \ker(A(t)) = \text{span}\{n_1(t), \dots, n_{m-r}(t)\}, \quad t \in [t_0, T],$$

$n_1, \dots, n_{m-r} \in C^1([t_0, T], \mathbf{R}^m)$, where r denotes the constant rank of $A(t)$.

Write shortly $C := C([t_0, T], \mathbf{R}^m)$, $C^1 := C^1([t_0, T], \mathbf{R}^m)$.

Assuming $\{A, B\} \in \mathcal{N}$, we consider the linear equation

$$(2.2) \quad Ax' + Bx = q.$$

Using continuously differentiable projection functions $Q, P: [t_0, T] \rightarrow L(\mathbf{R}^m)$ so that

$$(2.3) \quad Q(t)^2 \equiv Q(t), \quad \text{im}(Q(t)) \equiv N(t), \quad P(t) \equiv I - Q(t),$$

we may reformulate (2.2) as

$$(2.4) \quad A((Px)' - P'x) + Bx = q$$

and further relate the map $\mathfrak{A}: C_N^1 \rightarrow C$,

$$(2.5) \quad C_N^1 := \{x \in C: Px \in C_N^1\},$$

$$(2.6) \quad \mathfrak{A}x := A(Px)' + (B - AP')x, \quad x \in C_N^1,$$

to equation (2.2). The natural norm on C_N^1 is

$$\|x\| := \|x\|_\infty + \|(Px)'\|_\infty, \quad x \in C_N^1.$$

Note that $\{C_N^1, \|\cdot\|\}$ is a Banach space, and \mathfrak{U} becomes bounded.

DEFINITION. \mathfrak{U} is called *tractable* if

$$(2.7) \quad \dim(\ker(\mathfrak{U})) < \infty.$$

Clearly, tractability means that the solutions of the homogeneous equation $Ax' + Bx = 0$ form a finite-dimensional function space so that it becomes possible to select a unique solution by a finite number of initial conditions.

Up to now, the question whether $\ker(\mathfrak{U})$ has finite dimension is answered only partially. For constant coefficients A, B , we have (2.7) if and only if the matrix pencil $\lambda A + B$ is regular, i.e. $\det(\lambda A + B) \neq 0$ ([1], cf. also [8]). Further, in the case of sufficiently smooth coefficients A, B , some results are obtained via reduction methods (cf. [7], [5]). (2.7) is also true for all DAE's (2.2) having a global index in the sense of Gear and Petzold ([5]). Finally, index- k -tractability, $k \in \{1, 2, 3\}$ implies (2.7) provided $\{A, B\} \in \mathcal{N}$ and some canonical subspace varies smoothly ([13], [14]).

It should be mentioned that pre-tractability ([16]) is necessary but not sufficient for tractability. For a conjecture how to characterize tractability in terms of the coefficients $\{A, B\}$ we refer to [15].

We follow the idea to rely on certain canonical subspaces (cf. [8], [13], [14]). For a given pair $\{A, B\} \in \mathcal{N}$, we are going to use the matrix functions

$$(2.8) \quad G_1 := A + BQ,$$

$$(2.9) \quad A_1 := G_1 - AP'Q,$$

$$(2.10) \quad G_2 := A_1 + BPQ_1,$$

$$(2.11) \quad A_2 := G_2 - A_1(PP_1)'PQ_1,$$

$$(2.12) \quad G_3 := A_2 + BPP_1Q_2,$$

as well as the subspaces

$$(2.13) \quad N(t) := \ker(A(t)),$$

$$(2.14) \quad S(t) := \{z \in \mathbf{R}^m : B(t)z \in \text{im}(A(t))\},$$

$$(2.15) \quad N_1(t) := \ker(A_1(t)),$$

$$(2.16) \quad S_1(t) := \{z \in \mathbf{R}^m : B(t)P(t)z \in \text{im}(A_1(t))\},$$

$$(2.17) \quad N_2(t) := \ker(A_2(t)),$$

$$(2.18) \quad S_2(t) := \{z \in \mathbf{R}^m : B(t)P(t)P_1(t)z \in \text{im}(A_2(t))\};$$

where $t \in [t_0, T]$. Thereby, $Q_1(t)$, $Q_2(t)$ denote projections onto $N_1(t)$ and $N_2(t)$, respectively, further

$$P_1(t) := I - Q_1(t), \quad P_2(t) := I - Q_2(t).$$

Using A_2 , we always suppose PP_1 to be continuously differentiable. The subspaces (2.13) – (2.18) are called *canonical subspaces* of the DAE (2.2).

DEFINITION ([8], [13], [14]). Assume $\{A, B\} \in \mathcal{N}$. The DAE (2.2) is called

1. *transferable* (or *index-1-tractable*) if $G_1(t)$ remains nonsingular for all t ,
2. *index-2-tractable* if $G_1(t)$ is singular but $G_2(t)$ nonsingular for all t ,
3. *index-3-tractable* if both $G_1(t)$, $G_2(t)$ are singular but $G_3(t)$ remains nonsingular on $[t_0, T]$.

The point of index- k -tractability is to characterize special classes of tractable maps \mathfrak{A} respectively tractable DAE's (2.2), and to generalize the class of DAE's having the global index k (cf. [8], [13], [14]). In the case of $k > 1$ we speak of *higher-index* equations.

Recall ([8], [13], [14]) that $G_i(t)$ and $A_i(t)$ are always singular or nonsingular simultaneously. This fact is useful for testing index- k -tractability in practice. Moreover, recall also that transferability means

$$(2.19) \quad N(t) \oplus S(t) = \mathbf{R}^m, \quad t \in [t_0, T],$$

while index-2-tractability is equivalent to

$$(2.20) \quad N(t) \cap S(t) \neq \{0\}, \quad N_1(t) \oplus S_1(t) = \mathbf{R}^m, \quad t \in [t_0, T].$$

Finally, index-3-tractability is equivalent to the subspace-properties

$$(2.21) \quad N(t) \cap S(t) \neq \{0\}, \quad N_1(t) \cap S_1(t) \neq \{0\}, \quad N_2(t) \oplus S_2(t) = \mathbf{R}^m, \\ t \in [t_0, T].$$

It is well known that the formulation of appropriate initial conditions for index- k -tractable DAE's depends on the DAE itself. Here, we use the conditions

$$(2.22) \quad \mathcal{L}x := \begin{cases} P(t_0)x(t_0) & \text{for } k = 1, \\ P(t_0)P_1(t_0)x(t_0) & \text{for } k = 2, \\ P(t_0)P_1(t_0)P_2(t_0)x(t_0) & \text{for } k = 3 \end{cases} = b.$$

Denote, for index- k -tractable DAE's,

$$(2.23) \quad M := \begin{cases} \text{im}(P(t_0)) & \text{if } k = 1, \\ \text{im}(P(t_0)P_1(t_0)) & \text{if } k = 2, \\ \text{im}(P(t_0)P_1(t_0)P_2(t_0)) & \text{if } k = 3. \end{cases}$$

Clearly, the linear map

$$(2.24) \quad \mathcal{L} := (\mathfrak{A}, \mathcal{L}): C_N^1 \rightarrow C \times M$$

is bounded. The equation $\mathcal{L}x = (q, b)$ represents the initial value problem (2.2), (2.22). For the numerical treatment of such an IVP, it should be known whether this problem is well-posed in Hadamard's sense, i.e. whether \mathcal{L} has a bounded inverse. Unfortunately, \mathcal{L} is a homeomorphism only in the case of transferable (= index-1-tractable) DAE's. This is why we should differentiate the transferable DAE's from the nontransferable ones basically.

THEOREM 2.1. *Suppose that (2.2) is transferable, the IVP (2.2), (2.22) is uniquely solvable on C_N^1 for all $q \in C, b \in M$.*

Then, it holds that $\text{im}(\mathfrak{A}) = C, \dim(\ker(\mathfrak{A})) = r$, and \mathcal{L} becomes a homeomorphism.

([2], Theorem 1.2.20).

THEOREM 2.2. *Let (2.2) be index-2-tractable and $Q_1(t)$ be chosen to project onto $N_1(t)$ along $S_1(t)$. Assume $Q_1 \in C^1([t_0, T], \mathbb{R}^m)$, additionally.*

Then, the IVP (2.2), (2.22) is uniquely solvable on C_N^1 for all $q \in \mathfrak{M}, b \in M$, where

$$(2.25) \quad \mathfrak{M} := \{p \in C : Q_1 G_2^{-1} p \in C^1\}.$$

It holds that $\text{im}(\mathfrak{A}) = \mathfrak{M}$,

$$\dim(\ker(\mathfrak{A})) = \text{rank}(P(t_0)P_1(t_0)) < r.$$

\mathfrak{M} is a proper nonclosed subset within C , and \mathcal{L} has no bounded inverse.

([13], Theorem 2.4).

EXAMPLE. For the semi-explicit linear index-2 DAE

$$\begin{bmatrix} I \\ 0 \end{bmatrix} x' + \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & 0 \end{bmatrix} x = q$$

we compute

$$Q = \begin{bmatrix} 0 & \\ & I \end{bmatrix}, \quad Q_1 = \begin{bmatrix} I & B_{12} \\ 0 & 0 \end{bmatrix},$$

$$Q_1 = \begin{bmatrix} B_{12}(B_{21}B_{12})^{-1}B_{21} & 0 \\ -(B_{21}B_{12})^{-1} & B_{21} & 0 \end{bmatrix}.$$

Note that the projector function $B_{12}(B_{21}B_{12})^{-1}B_{21}$ is assumed to be smooth in [6], [12] also. ■

Recall from [13] that, under the assumptions of Theorem 2.2, the IVP (2.2), (2.22) is equivalent to the form

$$(2.26) \quad \begin{aligned} x &= (I - QP_1 G_2^{-1} B - (QQ_1)')y + PQ_1 G_2^{-1} q + QP_1 G_2^{-1} q \\ &\quad - (QQ_1)' PQ_1 G_2^{-1} q + (QQ_1 G_2^{-1} q)', \\ y' &= (PP_1)' y - PP_1 G_2^{-1} B y + (PP_1)' PQ_1 G_2^{-1} q + PP_1 G_2^{-1} q, \\ y(t_0) &= b. \end{aligned}$$

THEOREM 2.3. Suppose that (2.2) is index-3-tractable, and moreover, that $Q_1(t)$, $Q_2(t)$ are chosen in such a way that

$$(2.27) \quad Q_1(t) Q_2(t) = 0$$

holds and $Q_2(t)$ projects onto $N_2(t)$ along $S_2(t)$. Additionally, assume

$$Q_1, Q_2, \tilde{Q}_1 \in C^1([t_0, T], L(\mathbf{R}^m)),$$

$$Q_1 Q_2 \in C^2([t_0, T], L(\mathbf{R}^m)),$$

where $Q_1(t)$ denotes the projection onto $N_1(t)$ along

$$\{z \in \mathbf{R}^m: (B(t) - A_1(t)(PP_1)'(t))P(t)z \in \text{im}(G_3(t)P_1(t))\}.$$

Then, the IVP (2.2), (2.22) is uniquely solvable on C_N^1 for all $q \in \mathfrak{M}$, $b \in M$, where

$$(2.28) \quad \mathfrak{M} := \{p \in C: Q_2 G_3^{-1} p \in C^1, \\ Q_1 P_2 G_3^{-1} p + Q_1 Q_2 (PP_1 Q_2 G_3^{-1} q)' \in C^1\}.$$

It holds that $\text{im}(\mathfrak{L}) = \mathfrak{M}$

$$\dim(\ker(\mathfrak{L})) = \text{rank}(P(t_0)P_1(t_0)P_2(t_0)) < r.$$

\mathfrak{M} is a proper nonclosed subset within C , and \mathcal{L} has no bounded inverse.

([14], Theorem 2.4).

In both cases of Theorems 2.2 and 2.3, the range $\text{im}(\mathcal{L}) = \mathfrak{M} \times M$ is nonclosed within $C \times M$. Further, \mathcal{L} is injective, but its inverse \mathcal{L}^{-1} acting from $C \times M$ onto C_N^1 is no more continuous. Those problems in which the continuous dependence of the solutions on the right-hand sides is missing are called *ill-posed* or *essentially ill-posed* in Tikhonov's sense.

Surely, we could turn to maps $\mathcal{L}: C_N^1 \rightarrow \mathfrak{M} \times M$ and use an appropriate stronger norm $\|\cdot\|_{\mathfrak{M}}$ on \mathfrak{M} to have a Banach space $\{\mathfrak{M}, \|\cdot\|_{\mathfrak{M}}\}$ and a homeomorphism \mathcal{L} as well (cf. [17]). However, this approach seems to make no sense in view of the numerical treatment we are interested in, and also in view of possible linearizations of nonlinear problems.

It should also be mentioned, that \mathcal{L} itself becomes unbounded in the setting $\mathcal{L}: C \rightarrow C \times M$, while $\mathcal{L}^{-1}: C \times M \rightarrow C$ also remains unbounded. This is easily proved considering (2.26) and the related decoupling for the index 3 case in [14].

3. On the BDF applied to linear index-2-tractable DAE's

The first question we should deal with is about feasibility. A variable-order variable-stepsize BDF applied to (2.2) is simply

$$(3.1) \quad A(t_j) \frac{1}{h_j} \sum_{i=0}^k a_{ij} x_{j-i} + B(t_j) x_j = q(t_j),$$

thus, per integration step, a linear system with the coefficient matrix

$$(3.2) \quad F_j := A(t_j) \frac{a_{j_0}}{h_j} + B(t_j)$$

has to be solved.

LEMMA 3.1. *If (2.2) is index-2-tractable and*

$$(3.3) \quad I - Q(t) Q_1(t) P'(t)$$

remains nonsingular for all t , where $Q_1(t)$ denotes the canonical projection onto $N_1(t)$ along $S_1(t)$, then F_j is nonsingular for sufficiently small h_j .

Proof. Write shortly $F = A/\tau + B$, $\tau := h_j/a_{j_0}$, and drop the arguments t_j as well as the index j . The equation $Fz = y$ is equivalent to

$$(A + BQ - AP'Q) \left(\frac{1}{\tau} Pz + Qz \right) + BPz = y - AP'Qz,$$

further to

$$(A_1 + BPQ_1) \left(\frac{1}{\tau} P_1 Pz + P_1 Qz + Q_1 z \right) + BPP_1 z = y - AP'Qz$$

or

$$(3.4) \quad \frac{1}{\tau} P_1 Pz + P_1 Qz + Q_1 z + G_2^{-1} BPP_1 z = G_2^{-1} y - G_2^{-1} AP'Qz.$$

Now we may use the identity $Q_1 = Q_1 G_2^{-1} BP$ (cf. [13]) and, moreover, the relations (cf. [5])

$$G_2^{-1} AP'Q = G_2^{-1} A_1 PP'Q = P_1 PP'Q, \quad PP_1 P = PP_1, \\ PP_1 Q = 0, \quad QP_1 P = -QQ_1, \quad QP_1 Q = Q, \quad QP_1 PP'Q = -QQ_1 P'Q.$$

Multiplying (3.4) by τPP_1 , QP_1 and Q_1 , we obtain the system

$$PP_1 z + \tau PP_1 G_2^{-1} BPP_1 z + \tau PP_1 P'Qz = \tau PP_1 G_2^{-1} y, \\ -\frac{1}{\tau} QQ_1 z + Qz + QP_1 G_2^{-1} BPP_1 z - QQ_1 P'Qz = QP_1 G_2^{-1} y, \\ Q_1 z = Q_1 G_2^{-1} y.$$

Hence, we have

$$(3.5) \quad (I + \tau PP_1 G_2^{-1} B) PP_1 z + \tau PP_1 P'Qz = \tau PP_1 G_2^{-1} y, \\ QP_1 G_2^{-1} BPP_1 z + (I - QQ_1 P')Qz = QP_1 G_2^{-1} y + \frac{1}{\tau} QQ_1 G_2^{-1} y$$

to determine the components $PP_1 z$, Qz . For small $\tau > 0$, the matrix $H := I + \tau PP_1 G_2^{-1} B$ is nonsingular, and also the Kronecker complement in (3.5)

$$K := I - QQ_1 P' - \tau QP_1 G_2^{-1} BH^{-1} PP_1 P'$$

is so. Finally, solving (3.5) with respect to $PP_1 z$, Qz , we obtain

$$(3.6) \quad z = \{(I - \tau H^{-1} PP_1 P')K^{-1} \left(QP_1 + \frac{1}{\tau} QQ_1 - \tau QP_1 G_2^{-1} BH^{-1} PP_1 \right) + PQ_1 + \tau H^{-1} PP_1\} G_2^{-1} y. \blacksquare$$

EXAMPLE ([5]).

$$A(t) = \begin{bmatrix} 0 & 0 \\ 1 & \eta t \end{bmatrix}, \quad B(t) = \begin{bmatrix} 1 & \eta t \\ 0 & 1 + \eta \end{bmatrix}, \quad \eta \in \mathbf{R}.$$

The related DAE (2.2) is index-2-tractable. We compute

$$Q(t) = \begin{bmatrix} 0 & -\eta t \\ 0 & 1 \end{bmatrix}, \quad A_1(t) = \begin{bmatrix} 0 & 0 \\ 1 & 1 + \eta t \end{bmatrix},$$

$$Q_1(t) = \begin{bmatrix} 1 + \eta t & \eta t(1 + \eta t) \\ -1 & -\eta t \end{bmatrix}, \quad G_2(t) = \begin{bmatrix} 1 & \eta t \\ 1 & 1 + \eta t \end{bmatrix},$$

$$P(t)P_1(t) = 0, \quad I - Q(t)Q_1(t)P'(t) = \begin{bmatrix} 1 & -\eta^2 t \\ 0 & 1 + \eta \end{bmatrix}.$$

Thus the matrix (3.3) is nonsingular for $\eta \neq -1$ but singular otherwise. Moreover, it is easy to check that F_j becomes singular for $\eta = -1$. ■

Note that (3.6) gives an explicit expression for F_j^{-1} . In particular, in the constant-nullspace case we have $P' = 0$, $K = I$,

$$(3.8) \quad F_j^{-1} = \left\{ QP_1 + PQ_1 + \frac{a_{j_0}}{h_j} QQ_1 + \frac{h_j}{a_{j_0}} (I - QP_1 G_2^{-1} B)H^{-1} PP_1 \right\} G_2^{-1} \Big|_{t_j}.$$

The BDF's are well known to become unstable when they are applied to higher-index DAE's. This fact is true even in the case of constant-coefficient DAE's and constant-stepsize BDF's (cf. [8], [16]).

In the following we use grids

$$(3.9) \quad \pi: t_0 < t_1 < \dots < t^n = T$$

belonging to a given gridclass Π . The smallest gridclass we are interested in is the set Π_{equ} of all equidistant grids (3.9). Assume always

$$\Pi_{\text{equ}} \subseteq \Pi.$$

Denote by h and \underline{h} the maximal and minimal stepsize of (3.9), respectively.

Introduce the linear map $\mathcal{L}_\Pi: \mathbf{R}^{m(n+1)} \rightarrow \mathbf{R}^{m(n+1)}$,

Multiplying (3.14) by P and Q , respectively, we derive

$$(3.16) \quad \frac{1}{h_j} PP_1(t_j) \sum_{i=0}^k a_{ij} z_{j-i} + PP_1(t_j) G_2(t_j)^{-1} B_j PP_1(t_j) z_j \\ = PP_1(t_j) G_2(t_j)^{-1} w_j,$$

$$(3.17) \quad -\frac{1}{h_j} QQ_1(t_j) \sum_{i=0}^k a_{ij} z_{j-i} + Qz_j + QP_1(t_j) G_2(t_j)^{-1} B_j PP_1(t_j) z_j \\ = QP_1(t_j) G_2(t_j)^{-1} w_j,$$

Inserting the relation

$$\frac{1}{h_j} PP_1(t_j) \sum_{i=0}^k a_{ji} z_{j-i} \\ = \frac{1}{h_j} \sum_{i=0}^k a_{ji} PP_1(t_{j-i}) z_{j-i} \\ + \sum_{i=0}^k a_{ji} P \frac{1}{h_j} (P_1(t_j) - P_1(t_{j-i})) (PP_1(t_{j-i}) z_{j-i} + PQ_1(t_{j-i}) z_{j-i})$$

into (3.16) and using (3.15) we verify the inequalities

$$\max_{0 \leq j \leq n} |PP_1(t_j) z_j| \leq S_1 \max_{0 \leq j \leq n} |w_j|, \\ \max_{0 \leq j \leq n} |Q_1(t_j) z_j| \leq S_2 \max_{0 \leq j \leq n} |w_j|.$$

Hence, we have also

$$(3.18) \quad \max_{0 \leq j \leq n} |Pz_j| \leq S_3 \max_{0 \leq j \leq n} |w_j|.$$

On the other hand, (3.17) yields, for $j = k, \dots, n$,

$$(3.19) \quad Qz_j = \frac{1}{h_j} QQ_1(t_j) \sum_{i=0}^k a_{ji} \{PP_1(t_{j-i}) + PQ_1(t_{j-i})\} z_{j-i} \\ - QP_1(t_j) G_2(t_j)^{-1} B_j PP_1(t_j) z_j + QP_1(t_j) G_2(t_j)^{-1} w_j \\ = \frac{1}{h_j} QQ_1(t_j) \sum_{i=0}^k a_{ji} Q_1(t_{j-i}) \tilde{w}_{j-i} \\ + \sum_{i=0}^k a_{ji} Q \frac{1}{h_j} (Q_1(t_j) - Q_1(t_{j-i})) PP_1(t_{j-i}) z_{j-i} \\ - QP_1(t_j) G_2(t_j)^{-1} B_j PP_1(t_j) z_j + QP_1(t_j) \tilde{w}_j,$$

where

$$\tilde{w}_j := \begin{cases} w_j & \text{for } j = 0, \dots, k-1, \\ G_2(t_j)^{-1} w_j & \text{for } j = k, \dots, n. \end{cases}$$

Now, we have

$$(3.20) \quad |Qz_j| \leq \frac{1}{h_j} |QQ_1(t_j) \sum_{i=0}^k a_{ji} Q_1(t_{j-i}) \tilde{w}_{j-i}| + S_4 \max_{0 \leq l \leq n} |w_l|.$$

(3.18), (3.20) imply (3.11). Moreover, for $z_j = x(t_j) - x_j$, $j = 0, \dots, n$, where $x(t_j)$ denotes the value of the exact solution at t_j , and x_j is the exact numerical solution generated by the BDF using given starting values x_0, \dots, x_{k-1} , the corresponding values w_j are

$$(3.21) \quad \begin{aligned} w_j &= x(t_j) - x_j \quad \text{for } j = 0, \dots, k-1, \\ w_j = \tau_j &:= \frac{1}{h_j} A_j \sum_{i=0}^k a_{ji} x(t_{j-i}) + B_j x(t_j) - q(t_j) \\ &= \frac{1}{h_j} A_j \sum_{i=0}^k a_{ji} x(t_{j-i}) - A_j (Px)'(t_j), \end{aligned}$$

for $j = k, \dots, n$. Since $\tau_j \in \text{im}(A_j)$ we may write

$$w_j = \tau_j = A_j A_j^+ \tau_j, \quad j = k, \dots, n,$$

further

$$(3.22) \quad \begin{aligned} \tilde{w}_j &= G_2(t_j)^{-1} A_j A_j^+ \tau_j \\ &= G_2(t_j)^{-1} A_1(t_j) P A_j^+ \tau_j = P_1(t_j) P A_j^+ \tau_j, \quad j = k, \dots, n \end{aligned}$$

Surely, for exact starting values $x_j = x(t_j)$, $j = 0, \dots, k-1$, (3.20) yields simply

$$|Q(x(t_j) - x_j)| \leq S_4 \max_{k \leq l \leq n} |\tau_l|, \quad j \geq k,$$

further

$$(3.23) \quad \max_{0 \leq j \leq n} |x(t_j) - x_j| \leq S_5 \max_{k \leq l \leq n} |\tau_l|.$$

Finally, we collect our results in

THEOREM 3.2. *Let the assumptions of Theorem 2.2 be satisfied and, further, $P' = 0$. Let the variable-order variable-stepsizes BDF (3.1) be stable on $\Pi(\eta_1, \eta_2, h_{\max})$ for explicit ODE's. Then this BDF applied to the DAE (2.2) is weakly unstable and satisfies (3.11). Moreover, this BDF is convergent and has the same order as for explicit ODE's.*

Theorem 3.2 generalizes the related convergence results in [5], [6], [12] obtained for semi-explicit DAE's. For the values $\tilde{x}_k, \dots, \tilde{x}_n \in R^m$ generated by the perturbed equation

$$(3.24) \quad F_j \tilde{x}_j = q(t_j) - \frac{1}{h_j} A_j \sum_{i=1}^k a_{ji} \tilde{x}_{j-i} + \delta_j, \quad j = k, \dots, n,$$

using starting values $\tilde{x}_0, \dots, \tilde{x}_{k-1} \in \mathbf{R}^m$, we obtain the error estimation

$$(3.25) \quad \max_{0 \leq j \leq n} |x(t_j) - \tilde{x}_j| \leq S_6 \max \left\{ \max_{0 \leq j \leq k-1} |x(t_j) - \tilde{x}_j|, \max_{k \leq j \leq n} |\tau_j - \delta_j| \right\} \\ + \max_{k \leq j \leq n} \frac{1}{h_j} \left| Q Q_1(t_j) \sum_{i=0}^k a_{ji} Q_1(t_{j-i}) \tilde{\delta}_{j-i} \right|,$$

where

$$\tilde{\delta}_j = \begin{cases} -G_2(t_j)^{-1} \delta_j & \text{for } j = k, \dots, n, \\ x(t_j) - \tilde{x}_j & \text{for } j = 0, \dots, k-1, \end{cases}$$

and δ_j represents defects in the linear equations due to roundoff errors. Thus, the components $Q_1(t_j) \tilde{\delta}_j / h_j$ must be kept small by an appropriate error control.

From (3.10) we know immediately that

$$(3.26) \quad \|\mathcal{L}_H\|_\infty \leq \underline{h}^{-1} K.$$

This yields

$$(3.27) \quad \|\mathcal{L}_H\|_\infty \|\mathcal{L}_H^{-1}\|_\infty \leq \underline{h}^{-2} K S.$$

Similarly, the matrices F_j of the linear systems to be solved per integration step have the condition numbers

$$\|F_j\| \|F_j^{-1}\| \sim h_j^{-2}$$

since $Q Q_1$ is always nontrivial for index-2-tractable DAE's.

In comparison with this, we recall that

$$\|\mathcal{L}_H\|_\infty \leq \underline{h}^{-1} K_0, \quad \|\mathcal{L}_H^{-1}\|_\infty \leq S_0, \\ \|F_j\| \|F_j^{-1}\| \leq \text{const}$$

is true for explicit ODE's. In the case of transferable DAE's we have (cf. [8]) also

$$\|\mathcal{L}_H\|_\infty \leq \underline{h}^{-1} K_1, \quad \|\mathcal{L}_H^{-1}\|_\infty \leq S_1,$$

but $\|F_j\| \|F_j^{-1}\| \sim h_j^{-1}$.

Thus, it is not surprising that integration methods applied to index-2-DAE's depend more sensitively on the error control. However, using a careful error control, integration methods may work well (cf. [6], [12]) also for a class of index-2-DAE's.

However, for index-3-DAE's the situation becomes worse. At the best we may expect

$$\|\mathcal{L}_H^{-1}\|_\infty \leq \underline{h}^{-2} K \quad \text{and} \quad \|F_j\| \|F_j^{-1}\| \sim h_j^{-3}.$$

A detailed analysis of BDF's applied to index-3-DAE's is in preparation.

4. Some remarks on regularization methods

Since the IVP (2.2), (2.22) is essentially ill-posed in Tikhonov's sense, if (2.2) is index-2-tractable, numerical methods for ill-posed problems are of interest also for those nontransferable DAE's. We are going to deal with two regularizations which are very closely connected with singular perturbations.

The linear equation (2.2) is approximated by

$$(4.1) \quad (A + \varepsilon B)x'_\varepsilon + Bx_\varepsilon = q$$

using the so-called *pencil-regularization* (e.g. [1], [3]). An alternative way is the use of the regularization

$$(4.2) \quad (A + \varepsilon BP)x'_\varepsilon + (B + \varepsilon BPP')x_\varepsilon = q$$

proposed in [18].

With (4.1) one aims for a regular ODE. This goal is not reached in general. For instance, for the special DAE given in (3.7) with $\eta = -1$, the regularized equation (4.1) is again an index-2-tractable DAE. On the other hand, by (4.1), transferable DAE's are approximated by regular stiff ODE's, i.e. (4.1) yields "parasitic" boundary layers (cf. [11]). This is why (4.2) is proposed in [18].

If $P' = 0$, both (4.1) and (4.2) may be applied. The initial condition for index-2-tractable DAE's is (cf. (2.22))

$$(4.3) \quad P(t_0)P_1(t_0)x(t_0) = b.$$

For (4.2), we may use the additional condition (cf. (2.26))

$$(4.4) \quad P(t_0)Q_1(t_0)x(t_0) = P(t_0)Q_1(t_0)G_2(t_0)^{-1}q(t_0).$$

Additionally, when using (4.1), we have to determine $Q(t_0)x(t_0)$ besides (4.3), (4.4). However, the formula given by (2.26), i.e.

$$(4.5) \quad Q(t_0)x(t_0) = Q(t_0)P_1(t_0)G_2(t_0)^{-1}\{q(t_0) - B(t_0)b\} + (QQ_1G_2^{-1}q)'(t_0) \\ - (QQ_1)'(t_0)\{b + P(t_0)Q_1(t_0)G_2(t_0)^{-1}q(t_0)\},$$

is not reasonable for the practical use. Thus, we must compute an approximation of the consistent initial value for (4.1).

THEOREM 4.1. *Let the assumptions of Theorem 2.2 be satisfied, $q \in \mathfrak{M}$, $b \in M$. Then (4.2) is transferable for sufficiently small $\varepsilon > 0$.*

Let $x, x_\varepsilon \in C_N^1$ denote the solutions of (2.2), (4.3) and (4.2), (4.3), (4.4), respectively.

Then, $\|x_\varepsilon - x\|_{H_N^1} \rightarrow 0$ ($\varepsilon \rightarrow 0$), further, $\|Px_\varepsilon - Px\|_{L_2} = O(\varepsilon)$.

Moreover, if $PQ_1, PP_1 \in C^2$, $PQ_1G_2^{-1}q \in C^2$, then

$$\|x_\varepsilon - x\|_{H_N^1} = O(\varepsilon^{1/2}).$$

The proof is given in [10] for a certain special case. The general case is considered in [9].

EXAMPLE. The differentiation problem $x'_2 - x_1 = 0$, $x_2 = q_2$ is approximated by (cf. (4.1))

$$(4.6) \quad \begin{aligned} -\varepsilon x'_1 + x'_2 - x_1 &= 0, \\ \varepsilon x'_2 + x_2 &= q_2, \\ x_2(t_0) &= q_2(t_0), \quad x_1(t_0) \approx q'_2(t_0), \end{aligned}$$

and by (cf. (4.2))

$$(4.7) \quad \begin{aligned} x'_2 - x_1 &= 0, \\ \varepsilon x'_2 + x_2 &= q_2, \\ x_2(t_0) &= q_2(t_0), \end{aligned}$$

respectively.

It should be mentioned that further assertions concerning the convergence of (4.1) and (4.2), respectively, are given in [1], [3], [9], [11].

5. Nonlinear DAE's

Let us start this section formulating some smoothness conditions for the general nonlinear DAE

$$(5.1) \quad f(x'(t), x(t), t) = 0.$$

ASSUMPTION (A): $f \in C(\mathcal{G}, \mathbf{R}^m)$, $\mathcal{G} := \mathbf{R}^m \times \mathbf{R}^m \times [t_0, T]$. Let $f'_y(y, x, t)$, $f'_x(y, x, t)$ exist for all $(y, x, t) \in \mathcal{G}$, and let $f'_y, f'_x \in C(\mathcal{G}, L(\mathbf{R}^m))$. Let the nullspace of $f'_y(y, x, t)$ be independent of (y, x) . Denote

$$(5.2) \quad N(t) := \ker(f'_y(y, x, t)).$$

Let $N(\cdot)$ be smooth, and denote by Q, P the corresponding projection functions onto $N(\cdot)$ and its complement. ■

The nonlinear map $\mathfrak{A}: C_N^1 \rightarrow C$,

$$(5.3) \quad (\mathfrak{A}x)(t) := f((Px)'(t) - P'(t)x(t), x(t), t),$$

is defined on the whole space C_N^1 , and is Fréchet-differentiable there. For each given $x_* \in C_N^1$, the linearized map represents a linear DAE with the coefficients

$$(5.4) \quad \{f'_y(M_*(\cdot)), f'_x(M_*(\cdot))\} \in \mathcal{A},$$

where

$$M_*(t) := ((Px_*)'(t) - P'(t)x_*(t), x_*(t), t).$$

DEFINITION. Let (A) be satisfied, $k \in \{2, 3\}$, $x_* \in C_N^1$. Then (5.1) is called *transferable around x_** and *index- k -tractable at x_** , if the linear DAE having the coefficients (5.4) is transferable and index- k -tractable, respectively.

To test whether index- k -tractability is given say on the ball $\{x_* \in C_N^1: \|x_* - x_0\| \leq \varrho\}$, $x_0 \in C_N^1$ fixed, we may use the following matrices (cf. [13], [14]):

$$\begin{aligned}
 (5.5) \quad G_1(y, x, t) &:= f'_y(y, x, t) + f'_x(y, x, t) Q(t), \\
 (5.6) \quad A_1(y, x, t) &:= G_1(y, x, t) - f'_y(y, x, t) P'(t) Q(t), \\
 (5.7) \quad G_2(y, x, t) &:= A_1(y, x, t) + f'_x(y, x, t) P(t) Q_1(x, t), \\
 (5.8) \quad A_2(y, x, t) &:= G_2(y, x, t) - A_1(y, x, t) K(y, x, t) P(t) Q_1(x, t); \\
 (5.9) \quad G_3(y, x, t) &:= A_2(y, x, t) + f'_x(y, x, t) P(t) P_1(x, t) Q_2(y, x, t).
 \end{aligned}$$

Moreover, the nullspace of $A_1(y, x, t)$ is assumed to be independent of y and $Q(t)x$. Denote

$$(5.10) \quad \ker(A_1(y, x, t)) \equiv: N_1(x, t) \equiv N_1(P(t)x, t).$$

$Q_1(x, t)$ is defined to be a projection onto $N_1(x, t)$, further

$$P_1(x, t) := I - Q_1(x, t).$$

Using $A_2(y, x, t)$, we assume additionally that

$$\Pi(x, t) := P(t) P_1(x, t) = P(t) P_1(P(t)x, t)$$

depends continuously differentiable on (x, t) . Finally,

$$(5.11) \quad K(y, x, t) := \Pi'_x(x, t)(y + P'(t)x) + \Pi'_t(x, t).$$

For index-2-tractability, we have to check the singularity of $G_1(y, x, t)$ as well as the nonsingularity of $G_2(y, x, t)$ on the related neighbourhood of the trajectory of x_0 within \mathcal{G} . For index-3-tractability, both $G_1(y, x, t)$, $G_2(y, x, t)$ have to be singular but $G_3(y, x, t)$ has to become nonsingular (cf. [13] Theorem 3.1, [14], Theorem 4.1).

Note that semi-explicit DAE's having the global index $k \in \{2, 3\}$ (cf. [5], [12]) are shown to be index- k -tractable ([13], [14]).

Now, let us turn to the semi-explicit system

$$(5.12) \quad u'(t) - v(t) = 0,$$

$$(5.13) \quad v'(t) + g(u(t), v(t), t) + h'_u(u(t), t)^T w(t) = 0,$$

$$(5.14) \quad h(u(t), t) = 0,$$

which describes, e.g., mechanical motions subjected to constraints (cf. [6], [12]). Assume that $h'_u(u, t) =: H(u, t)$ has full rank, i.e. the constraints (5.14) are linearly independent, $\ker(H(u, t)^T) = \{0\}$. Putting (5.12)–(5.14) in the general form (5.1), we obtain, with $x = (u, v, w)$.

$$f'_y(y, x, t) = \begin{bmatrix} I & & \\ & I & \\ & & 0 \end{bmatrix}, \quad f'_x(y, x, t) = \begin{bmatrix} 0 & -I & 0 \\ g'_u + h''_{uu}{}^T w & g'_v & H^T \\ H & 0 & 0 \end{bmatrix},$$

$$A_1(y, x, t) = \begin{bmatrix} I & & \\ & I & H(u, t)^T \\ & & 0 \end{bmatrix}, \quad Q = \begin{bmatrix} 0 & & \\ & 0 & \\ & & I \end{bmatrix},$$

Clearly, the nullspace of $A_1(y, x, t)$ is independent of y and $Qx = w$. Compute

$$Q_1 = \begin{bmatrix} 0 & & \\ & H^T(HH^T)^{-1}H & \\ & -(HH^T)^{-1}H & 0 \end{bmatrix}, \quad PP_1 = \begin{bmatrix} I & & \\ & I - H^T(HH^T)^{-1}H & \\ & & 0 \end{bmatrix}$$

Now, the assumption for $\Pi := PP_1$ to be smooth appears to be a smoothness condition for the nullspace of $H(u, t)$. Since $H(u, t)$ has constant rank, a sufficient condition for this nullspace to be smooth is that $H(u, t)$ itself is continuously differentiable. Then, also P_1 becomes continuously differentiable.

Further we compute $Q_1 Q = 0$,

$$G_2 = \begin{bmatrix} I & -T & \\ & I + g'_v T H^T & \\ & & 0 \end{bmatrix},$$

$$T := H^T(HH^T)^{-1}H, \quad S := I - T,$$

$$A_2 = \begin{bmatrix} I & -T & \\ & I + (g'_v + T'_u \beta + T'_t) T H^T & \\ & & 0 \end{bmatrix},$$

$$f'_x PP_1 = \begin{bmatrix} 0 & -I & 0 \\ g'_u + h''_{uu}{}^T w & g'_v & 0 \\ H & 0 & 0 \end{bmatrix}$$

Finally,

$$A_2(y, x, t)z = 0, \quad f'_x(y, x, t)P(t)P_1(x, t)z \in \text{im}(A_2(y, x, t)), \quad z = \begin{bmatrix} a \\ b \\ c \end{bmatrix},$$

imply

$$a = Tb, \quad b + (\dots)Tb + H^T c = 0, \quad Ha = 0,$$

thus $a = Ta$, $b = Tb$, further $a = 0$, $b = 0$, $H^T c = 0$, hence $c = 0$. In consequence, (5.12)–(5.14) is index-3-tractable everywhere.

Following the idea of [6], Theorem 1.1, where systems (5.12)–(5.14) with autonomous constraints are considered, we turn from (5.12)–(5.14) to the equivalent system

$$(5.15) \quad \begin{aligned} u'(t) - v(t) + h'_u(u(t), t)^T z(t) &= 0, \\ v'(t) + g(u(t), v(t), t) + h'_u(u(t), t)^T w(t) &= 0, \\ h(u(t), t) &= 0, \\ h'_u(u(t), t)v(t) + h'_t(u(t), t) &= 0. \end{aligned}$$

Surely, if $u_*, v_* \in C^1, w_* \in C$ form a solution of (5.12)–(5.14) then $(u_*, v_*, w_*, \mathcal{O})$ satisfies (5.15). Conversely, for each solution (u_*, v_*, w_*, z_*) of (5.15) with $z_* = \mathcal{O}$ the first components (u_*, v_*, w_*) solve (5.12)–(5.14). The point of (5.15) is that any solution has a trivial fourth component (cf. [6]). Namely, we have $h'_v v + h'_t = 0, u' = v - h'_u{}^T z$. Differentiating the constraint in (5.15) yields

$$h''_u u' + h'_t = 0,$$

hence

$$h'_u v + h'_t - h'_u h''_u{}^T z = 0,$$

and finally $z = 0$.

The DAE (5.15) is index-2-tractable. To show this we put this system in our general form (5.1) and compute the related matrices. We have

$$f'_y = \begin{bmatrix} I & & & \\ & I & & \\ & & 0 & \\ & & & 0 \end{bmatrix}, \quad f'_x = \begin{bmatrix} h''_{uu}{}^T z & -I & 0 & H^T \\ g'_u + h''_{uu}{}^T w & g'_v & H^T & 0 \\ H & 0 & 0 & 0 \\ h''_{uu} v + h''_{tt} & H & 0 & 0 \end{bmatrix},$$

$$A_1 = \begin{bmatrix} I & & & H^T \\ & I & H^T & \\ & & 0 & \\ & & & 0 \end{bmatrix}, \quad Q_1 = \begin{bmatrix} H^T(HH^T)^{-1} H & & & \\ & H^T(HH^T)^{-1} H & & \\ & & -(HH^T)^{-1} H & 0 \\ & & & 0 \end{bmatrix}.$$

Now,

$$\tilde{z} \in \ker(A_1(y, x, t)), \quad f'_x(y, x, t) P \tilde{z} \in \text{im}(A_1(y, x, t)), \quad \tilde{z} = [a^T \ b^T \ c^T \ d^T]^T,$$

imply

$$a + H^T d = 0, \quad b + H^T c = 0, \quad Ha = 0, \quad (\dots) a + Hb = 0.$$

This yields $a = 0, d = 0, b = 0, c = 0$, i.e. the matrix $G_2 = A_1 + f'_x P Q_1$ remains nonsingular. Hence, (5.15) is index-2-tractable.

In [6], [12], the BDF's applied to index-2 systems are reported to work well. But the BDF applied to index-3 systems becomes much more unreliable. A reduction step, for instance from (5.12)–(5.14) to (5.15), is recommended in [6] (cf. also [7]).

Now, consider the variable-order variable-step BDF applied to a non-linear index-2-tractable DAE (5.1). Denote by $x_* \in C_N^1$ the solution of (5.1) to be approximated. Let starting values $x_0, \dots, x_{k-1} \in R^m$ be given. The BDF is now

$$(5.16) \quad f\left(\frac{1}{h_j} \sum_{i=0}^k a_{ji} x_{j-i}, x_j, t_j\right) = 0, \quad j = k, \dots, n.$$

Recall that convergence is proved for certain semi-explicit systems ([2], [6], [12]), in particular, for (5.15).

The related map acting in $\mathbf{R}^{m(n+1)}$ is

$$\mathcal{F}_\Pi z := \begin{bmatrix} z_0 - x_0 \\ z_{k-1} - x_{k-1} \\ f\left(\frac{1}{h_k} \sum_{i=0}^k a_{ki} z_{k-i}, z_k, t_k\right) \\ \vdots \\ f\left(\frac{1}{h_n} \sum_{i=0}^k a_{ni} z_{n-i}, z_n, t_n\right) \end{bmatrix}, \quad z := \begin{bmatrix} z_0 \\ z_1 \\ \vdots \\ z_n \end{bmatrix} \in \mathbf{R}^{m(n+1)}.$$

Surely, \mathcal{F}_Π is continuously differentiable provided assumption (A) is satisfied. More precisely, the Jacobian taken at

$$z^* = \begin{bmatrix} x_*(t_0) \\ \vdots \\ x_*(t_n) \end{bmatrix}$$

is

$$(5.17) \quad \mathcal{F}'_\Pi(z^*) = \begin{bmatrix} I \\ \dots \\ \dots \\ \frac{a_{kk}}{h_k} \tilde{A}_k \quad \dots \quad \frac{a_{k1}}{h_k} \tilde{A}_k \quad \tilde{F}_k \\ \dots \\ \dots \\ \frac{a_{nk}}{h_n} \tilde{A}_n \quad \dots \quad \frac{a_{n1}}{h_n} \tilde{A}_n \quad \tilde{F}_n \end{bmatrix}$$

where

$$\tilde{A}_j := f'_y(\eta_j), \quad \tilde{F}_j := \frac{a_{j0}}{h_j} f'_y(\eta_j) + f'_x(\eta_j),$$

$$\eta_j := \left(\frac{1}{h} \sum_{i=0}^k a_{ji} x_*(t_{j-i}), x_*(t_j), t_j \right) \in \mathcal{G}.$$

Besides the nonlinear DAE (5.1) we consider the linear DAE

$$(5.18) \quad Az' + Bz = q$$

with the coefficients $\{A, B\} = \{f'_y(M_*(\cdot)), f'_x(M_*(\cdot))\} \in \mathcal{N}$ (cf. (5.4)). This linear equation is index-2-tractable. Apply the BDF under consideration to the linear DAE (5.18) also. For $P' = 0$, the matrix \mathcal{L}_Π (cf. (3.10)) which is related to the BDF applied to (5.18) is nonsingular. Supposed the related projection $Q_1(t)$ is smooth, the inequalities (3.11), (3.18) and (3.20) hold for $\mathcal{L}_\Pi z = w$.

Now, we turn to a more restricted class of nonlinear DAE's, namely to quasilinear ones of the form

$$(5.19) \quad A(t) x'(t) + g(x(t), t) = 0$$

where the nullspace of $A(t)$ is assumed to be constant, and $P' = 0$. Moreover, let $g'_x(x, t)$ be Lipschitz with respect to x . Then, for $j = k, \dots, n$, it holds that

$$\begin{aligned} A_j - \tilde{A}_j &= f'_y(M_*(t_j)) - f'_y(\eta_j) = 0, \\ F_j - \tilde{F}_j &= f'_x(M_*(t_j)) - f'_x(\eta_j) = 0. \end{aligned}$$

This yields

$$(5.20) \quad \mathcal{L}_H = \mathcal{F}'_H(z^*),$$

further, for

$$\begin{aligned} \tilde{\mathcal{F}}_H(z, \bar{z}) &:= \int_0^1 \mathcal{F}'_H(s z + (1-s)\bar{z}) ds, \\ \tilde{\mathcal{F}}_H(z, \bar{z}) - \mathcal{F}'_H(z^*) &= \text{diag}(0, \dots, 0, \alpha_k, \dots, \alpha_n), \end{aligned}$$

where

$$\alpha_j := \int_0^1 (g'_x(s z_j + (1-s)\bar{z}_j, t_j) - g'_x(z_j^*, t_j)) ds.$$

Now, for

$$z, \bar{z} \in \mathcal{B}_H(z^*, \varrho) := \{\bar{z} \in \mathbf{R}^{m(n+1)} : \|\bar{z} - z^*\|_\infty \leq \varrho\},$$

we derive

$$\|\alpha_j\|_\infty \leq \varrho L, \quad j = k, \dots, n,$$

thus

$$(5.21) \quad \|\tilde{\mathcal{F}}_H(z, \bar{z}) - \mathcal{F}'_H(z^*)\|_\infty \leq \varrho L.$$

The relations (5.21), (5.20), (3.11) imply, by the Banach Lemma, the nonsingularity of $\tilde{\mathcal{F}}_H(z, \bar{z})$ as well as the inequality

$$(5.22) \quad \|\tilde{\mathcal{F}}_H(z, \bar{z})^{-1}\|_\infty = \frac{\underline{h}^{-1} S}{1 - \underline{h}^{-1} S \varrho L},$$

if $\varrho > 0$ is chosen so that

$$(5.23) \quad \underline{h}^{-1} S L \varrho < 1$$

is valid.

Introduce the map $E_H: \mathbf{R}^{m(n+1)} \rightarrow \mathbf{R}^{m(n+1)}$ by

$$E_H z := z - \mathcal{F}'_H(z^*)^{-1} \mathcal{F}_H z, \quad z \in \mathbf{R}^{m(n+1)}.$$

For $z, \bar{z} \in \mathcal{B}_H(z^*, \varrho)$, we obtain

$$\begin{aligned} (5.24) \quad \|E_H z - E_H \bar{z}\|_\infty &= \|z - \bar{z} - \mathcal{F}'_H(z^*)^{-1} (\mathcal{F}_H z - \mathcal{F}_H \bar{z})\|_\infty \\ &\leq \|\mathcal{F}'_H(z^*)^{-1}\|_\infty \|\mathcal{F}'_H(z^*) - \tilde{\mathcal{F}}_H(z, z)\|_\infty \|z - \bar{z}\|_\infty \\ &\leq \underline{h}^{-1} S \varrho L \|z - \bar{z}\|_\infty, \end{aligned}$$

$$\begin{aligned}
(5.25) \quad \|E_{\Pi} z - z^*\|_{\infty} &= \|z - z^* - \mathcal{F}'_{\Pi}(z^*)^{-1} (\mathcal{F}_{\Pi} z - \mathcal{F}_{\Pi} z^* + \mathcal{F}_{\Pi} z^*)\|_{\infty} \\
&\leq \|\mathcal{F}'_{\Pi}(z^*)^{-1}\|_{\infty} \|\mathcal{F}'_{\Pi}(z^*) - \tilde{\mathcal{F}}_{\Pi}(z, z^*)\|_{\infty} \|z - z^*\|_{\infty} + \|\mathcal{F}'_{\Pi}(z^*)^{-1} \mathcal{F}_{\Pi} z^*\|_{\infty} \\
&\leq \underline{h}^{-1} S_{\varrho} L \|z - z^*\|_{\infty} + \|\mathcal{F}'_{\Pi}(z^*)^{-1} \mathcal{F}_{\Pi} z^*\|_{\infty}.
\end{aligned}$$

Denote $\mathcal{F}_{\Pi} z^* =: \tau$. The values $\tau_j = x_*(t_j) - x_j$, $j = 0, \dots, k-1$, represent the errors of the starting values while

$$\tau_j = A(t_j) \frac{1}{h_j} \sum_{i=0}^k a_{ji} x_*(t_{j-i}) + g(x_*(t_j), t_j), \quad j = k, \dots, n,$$

are the local discretization errors.

Taking into account that (cf. (3.22), (5.7))

$$\begin{aligned}
\tau_j &\in \text{im}(A(t_j)), \\
G_2(t_j)^{-1} \tau_j &= G_2(t_j)^{-1} A(t_j) A(t_j)^+ \tau_j = P_1(t_j) P A(t_j)^+ \tau_j, \\
G_2(t_j) &:= G_2(M_*(t_j)),
\end{aligned}$$

hold for $j = k, \dots, n$, we derive the inequality

$$\begin{aligned}
(5.26) \quad \|\mathcal{F}'_{\Pi}(z^*)^{-1} \mathcal{F}'_{\Pi} z^*\|_{\infty} \\
\leq S_7 \|\tau\|_{\infty} + S_8 \max_{j=k, \dots, 2k-1} \frac{1}{h_j} \left| \mathcal{Q} \mathcal{Q}_1(t_j) \sum_{i=j-k+1}^k a_{ji} \mathcal{Q}_1(t_{j-i}) \tau_{j-i} \right|
\end{aligned}$$

from (3.18), (3.19), (3.20).

Now, for a fixed value $\alpha \in (0, 1)$ we choose a refined grid $\pi \in \Pi(\eta_1, \eta_2, h_{\max})$ and sufficiently accurate starting values so that

$$(5.27) \quad \|\mathcal{F}'_{\Pi}(z^*) \mathcal{F}_{\Pi} z^*\|_{\infty} \leq (1 - \alpha).$$

This is always possible by (5.26). Practically, this means a choice $x_j = x_*(t_j) + O(h_j^{\mu+1})$, $j = 0, \dots, k-1$, applying the μ -order BDF.

Then, we choose ϱ small enough so that

$$(5.28) \quad \underline{h}^{-1} S L \varrho < \alpha.$$

In this way, the map E_{Π} is contractive with $\alpha < 1$ on the closed ball $\mathcal{B}_{\Pi}(z^*, \varrho)$, and $E_{\Pi}(\mathcal{B}_{\Pi}(z^*, \varrho)) \subseteq \mathcal{B}_{\Pi}(z^*, \varrho)$. Consequently, by Banach's Fixed Point Theorem, the map E_{Π} has a unique fixed point onto $\mathcal{B}_{\Pi}(z^*, \varrho)$, i.e. the nonlinear equation

$$\mathcal{F}_{\Pi} z = 0,$$

representing our integration method has a locally unique solution, i.e. there exist locally unique values $x_k, \dots, x_n \in \mathbf{R}^m$, satisfying (5.16).

Since, for $z, \bar{z} \in \mathcal{B}_{\Pi}(z^*, \varrho)$,

$$z - \bar{z} = \tilde{\mathcal{F}}_{\Pi}(z, \bar{z})^{-1} (\mathcal{F}_{\Pi} z - \mathcal{F}_{\Pi} \bar{z})$$

holds, the inequality

$$(5.29) \quad \|z - \bar{z}\|_\infty \leq \frac{h^{-1} S}{1 - h^{-1} S \varrho L} \|\mathcal{F}_\Pi z - \mathcal{F}_\Pi \bar{z}\|_\infty$$

follows (cf. (5.22)). Thus, (5.16) is shown to be only weakly unstable.

Using (5.25) for the true numerical solution $z = x$, i.e. the fixed point of E_Π , we obtain the error estimation

$$(5.30) \quad \max_{j=0, \dots, n} |x_*(t_j) - x_j| \leq \frac{1}{1 - \alpha} \|\mathcal{F}'_\Pi(z^*)^{-1} \mathcal{F}_\Pi z^*\|_\infty.$$

Consequently, if exact starting values are used, i.e. $x_*(t_j) = x_j, j = 0, \dots, k - 1$, we have simply

$$(5.31) \quad \max_{j=0, \dots, n} |x_*(t_j) - x_j| \leq \frac{1}{1 - \alpha} \max_{j=k, \dots, n} |\tau_j|$$

that means the variable-order variable-step BDF applied to the quasilinear DAE (5.19) is convergent with the same order as in the case of regular explicit ODE's.

THEOREM 5.1. *Let $x_* \in C_N^1$ be a solution of the quasilinear DAE (5.19). Let (5.19) be index-2-tractable in the neighbourhood of x_* . Further, let g'_x be Lipschitz there with respect to x , and let $A(t)$ have a constant nullspace.*

Moreover, let the projection $Q_1(t)$ onto the nullspace $N_1(t) := \ker(A_1(t))$ along $S_1(t) := \{z \in \mathbf{R}^m : g'_x(x_(t), t)Pz \in \text{im}(A_1(t))\}$ be continuously differentiable, where*

$$(5.32) \quad A_1(t) := A_1(M_*(t)) = A(t) + g'_x(x_*(t), t)Q$$

Let the variable-step variable order BDF be stable on $\Pi(\eta_1, \eta_2, h_{\max})$ in the regular ODE-case, h_{\max} sufficiently small. Then this method is only weakly unstable on $\Pi(\eta_1, \eta_2, h_{\max})$. It is convergent with the same order as in the case of regular ODE's.

Moreover, for each grid $\pi \in \Pi(\eta_1, \eta_2, h_{\max})$, for sufficiently accurate starting values x_0, \dots, x_{k-1} , the nonlinear equations (5.16) are uniquely solvable with respect to x_j on balls $\{u \in \mathbf{R}^m : |x_(t_j) - u| \leq \varrho\}, j = k, \dots, n$, where ϱ depends on Π ,*

$$(5.33) \quad \lim_{h \rightarrow 0} \varrho = 0.$$

The behaviour of ϱ (cf. e.g. (5.23)) seems not to be only a consequence of the technique used for the proof but reflects as matters stand in the nonlinear equations to be solved. In any case, we are not able to compute the values x_j exactly. Instead of x_j we generate only certain $\tilde{x}_j \in \mathbf{R}^m$ satisfying (5.16) approximately, i.e. we have

$$(5.34) \quad f\left(\frac{1}{h_j} \sum_{i=0}^k a_{ji} \tilde{x}_{j-i}, \tilde{x}_j, t_j\right) = \delta_j, \quad j = k, \dots, n,$$

$$(5.35) \quad \tilde{x}_j := x_j, \quad j = 0, \dots, k-1.$$

Recall that $\tau_j := x_*(t_j) - x_j = x_*(t_j) - \tilde{x}_j$, $j = 0, \dots, k-1$, denote the errors in the starting values.

THEOREM 5.2. *Let the assumptions of Theorem 5.1 be satisfied. Then, for sufficiently small defects δ_j in (5.34), (5.35), the error estimation*

$$(5.36) \quad \max_{j \geq k} |x_*(t_j) - \tilde{x}_j| \leq S_9 \max \left\{ \max_{j \leq k-1} |\tau_j|, \max_{j \geq k} |\tau_j - \delta_j| \right\} \\ + \max_{j \geq k} \left| QQ_1(t_j) \frac{1}{h_j} \sum_{i=0}^k a_{ji} Q_1(t_{j-i}) \omega_{j-i} \right|$$

holds, where

$$\omega_j := \begin{cases} \tau_j, & j \leq k-1, \\ G_2(t_j)^{-1} \delta_j, & j \geq k, \end{cases}$$

$$G_2(t_j) := A(t_j) + g'_x(x_*(t_j), t_j)(Q + PQ_1(t_j)).$$

Proof. With

$$\tilde{x} = \begin{bmatrix} \tilde{x}_0 \\ \vdots \\ \tilde{x}_n \end{bmatrix}, \quad z^* = \begin{bmatrix} x_*(t_0) \\ \vdots \\ x_*(t_n) \end{bmatrix}, \quad \delta = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \delta_k \\ \vdots \\ \delta_n \end{bmatrix}, \quad \tau = \begin{bmatrix} \tau_0 \\ \vdots \\ \tau_n \end{bmatrix}$$

we have $\mathcal{F}_H \tilde{x} = \delta$, $\mathcal{F}_H z^* = \tau$,

$$\tilde{x} - z^* = \tilde{x} - z^* - \mathcal{F}'_H(z^*)^{-1} \{ \mathcal{F}_H \tilde{x} - \mathcal{F}_H z^* \} + \mathcal{F}'_H(z^*)^{-1} \{ \mathcal{F}_H \tilde{x} - \mathcal{F}_H z^* \} \\ = \mathcal{F}'_H(z^*)^{-1} \{ \mathcal{F}'_H(z^*) - \mathcal{F}'_H(\tilde{x}, z^*) \} (\tilde{x} - z^*) + \mathcal{F}'_H(z^*)^{-1} (\delta - \tau),$$

trivially. Hence (cf. (5.21), (3.11), (5.28))

$$(5.37) \quad \|\tilde{x} - z^*\|_\infty \leq h^{-1} S_9 L \|\tilde{x} - z^*\|_\infty + \|\mathcal{F}'_H(z^*)^{-1} (\delta - \tau)\|_\infty,$$

further

$$(5.38) \quad \|\tilde{x} - z^*\|_\infty \leq \frac{1}{1 - \alpha} \|\mathcal{F}'_H(z^*)^{-1} (\delta - \tau)\|_\infty.$$

Due to (3.18), (3.20), we find

$$(5.39) \quad \|\mathcal{F}'_H(z^*)^{-1} (\delta - \tau)\|_\infty \\ \leq S_{10} \|\sigma - \tau\|_\infty + \max_{j=k, \dots, n} \frac{1}{h_j} \left| QQ_1(t_j) \sum_{i=0}^k a_{ji} Q_1(t_{j-i}) \omega_{j-i} \right|$$

where

$$\begin{aligned}\omega_j &:= \tau_j, \quad j = 0, \dots, k-1, \\ \omega_j &:= G_2^*(t_j)^{-1} \delta_j. \quad \blacksquare\end{aligned}$$

EXAMPLE. For the Euler backward method applied to the index-2 system (5.15) we have

$$PQ_1 = \text{diag}(T, T, 0, 0), \quad T := H^+ H, \quad H^+ := H^T (HH^T)^{-1},$$

$$\begin{aligned} & QQ_1(t_j) \frac{1}{h_j} (Q_1(t_j) \omega_j - Q_1(t_{j-1}) \omega_{j-1}) \\ &= \begin{cases} QQ_1(t_1) \frac{1}{h_1} G_2(t_1)^{-1} \delta_1 - QQ_1(t_1) \frac{1}{h_1} Q_1(t_0) (x_*(t_0) - x_0) \\ QQ_1(t_j) \frac{1}{h_j} (Q_1(t_j) G_2(t_j)^{-1} \delta_j - Q_1(t_{j-1}) G_2(t_{j-1})^{-1} \delta_{j-1}) \end{cases} \quad \text{for } j \geq 1, \end{aligned}$$

further

$$Q_1(t) G_2(t)^{-1} \begin{bmatrix} \bar{u} \\ \bar{v} \\ \bar{w} \\ \bar{z} \end{bmatrix} = \begin{bmatrix} H^+ \bar{w} \\ H^+ (z - (h''_{uu} v + h''_{iu}) H^+ \bar{w}) \\ -(HH^T)^{-1} (z - (h''_{uu} v + h''_{iu}) H^+ \bar{w}) \\ -(HH^T)^{-1} w \end{bmatrix}.$$

Thus, the defects in the third and fourth equation of the system (5.15) should be kept small in comparison with the defects in the first and second equation. This corresponds completely to the considerations in [6], [12] as well as to the practical experience reported e.g. in [6] and [19]. \blacksquare

References

- [1] Ju. E. Boyarintsev, *Regulyarnye i singulyarnye sistemy linejnyh obyknovennyh differentsialnyh uravnenij*, Nauka (Sib. otd.), Novosibirsk 1980.
- [2] K. E. Brennan and B. E. Engquist, *Backward difference approximations of nonlinear differential-algebraic equations*, California, El Segundo, The Aerospace Corporation, 1985 (ATR-85 (9990)-5).
- [3] S. L. Campbell, *Regularizations of linear time varying singular systems*, Automatica — J. IFAC 20 (1984) 365–370.
- [4] F. R. Gantmakher, *Theory of Matrices*.
- [5] C. W. Gear and L. Petzold, *ODE methods for the solution of differential/algebraic systems*, SIAM J. Numer. Anal. 21 (1984), 716–728.
- [6] C. W. Gear, B. Leimkuhler and G. K. Gupta, *Automatic integration of Euler-Lagrange equations with constraints*, J. Comput. Appl. Math. 12-13 (1985), 77–90.
- [7] E. Griepentrog, *Transformation concepts for differential-algebraic equations*, this volume, 223–231.
- [8] E. Griepentrog and R. März, *Differential-Algebraic Equations and Their Numerical Treatment*, Leipzig; BSB B. G. Teubner, 1986 (Teubner-Texte zur Mathem. 88).
- [9] M. Hanke, *On the regularization of index 2 differential-algebraic equations*, Berlin, Humboldt University, Section of Mathematics, 1987 (Preprint 137).

- [10] M. Hanke, R. März and A. Neubauer, *On the regularization of a class of nontransferable differential algebraic equations*, J. Differential Equations, 73 (1988), 119-132.
- [11] —, —, —, *On the regularization of linear differential-algebraic equation*, in H. W. Engl and C. W. Groetsch (eds.), *Inverse and Ill-Posed Problems*, Academic Press, New York 1987.
- [12] P. Lötstedt and L. Petzold, *Numerical solution of nonlinear differential equations with algebraic constraints, I*, Math. Comp. 46 (1986), 491-516.
- [13] R. März, *Index-2 differential-algebraic equations*, Results in Math. 15 (1989), 149-171.
- [14] —, *Some new results concerning index-3-differential-algebraic equations*, J. Math. Anal. Appl., 140 (1989), 177-199.
- [15] —, *Recent Results on Higher-Index Differential-Algebraic Equations*, Teubner Texte zur Mathem. (Proc. Halle 1987), 104, 104-112.
- [16] —, *On the numerical treatment of differential-algebraic equation*, Z. Angew. Math. Mech. 67 (1987), 23-34.
- [17] —, *On correctness and numerical treatment of boundary value problem in DAE's*, Zh. Vychisl. Mat. i Mat. Fiz. 26 (1986), 50-64.
- [18] —, *On tractability with index 2*, Berlin, Humboldt University, Section of Mathematics, 1987 (Preprint 109).
- [19] L. Petzold and P. Lötstedt, *Numerical solution of nonlinear differential equations with algebraic constraints, II, Practical implications*, SIAM J. Sci. Statist. Comput. 7 (1986), 720-733

*Presented to the Semester
Numerical Analysis and Mathematical Modelling
February 25 — May 29, 1987*
