

On the Numerical Analysis of Overdetermined Linear Partial Differential Systems^{*}

Marcus Hausdorf and Werner M. Seiler

Lehrstuhl für Mathematik I, Universität Mannheim,
68131 Mannheim, Germany,
{hausdorf,werner.seiler}@math.uni-mannheim.de

Summary. We discuss the use of the formal theory of differential equations in the numerical analysis of general systems of partial differential equations. This theory provides us with a very powerful and natural framework for generalising many ideas from differential algebraic equations to partial differential equations. We study in particular the existence and uniqueness of (formal) solutions, the method of an underlying system, various index concepts and the effect of semi-discretisations.

Key words: overdetermined linear system, partial differential equation, involution, completion, index, semi-discretisation, underlying equation

1 Introduction

The majority of the literature on differential equations is concerned with normal systems or systems in Cauchy–Kovalevskaya form. But many important systems arising in applications are not of this form. As examples we mention Maxwell’s equations of electrodynamics, the incompressible Navier–Stokes equations of fluid dynamics, the Yang–Mills equations describing fundamental particle interactions or Einstein’s equations of general relativity.

For ordinary differential equations, the importance of non-normal systems has been recognised for about twenty years; one usually speaks of *differential algebraic equations*. Introductions into their theory can be found e. g. in [2, 8]. Recently, the extension to partial differential systems has found some interest, see e. g. [5, 15]. However, this is non-trivial, as new phenomena appear.

About a century ago the first methods for the analysis of general partial differential systems were designed; by now, a number of different approaches exist. Some of them have already been applied in a numerical context [14, 16, 20, 25]. We use the formal theory [17, 23] with its central notion of an *involution system*. In contrast to our earlier works [9, 22, 24], we take a more algebraic point of view closely related to the theory of involutive bases [3, 6]. For simplicity, we concentrate on linear systems, although many results remain valid in the non-linear case.

^{*} Supported by Deutsche Forschungsgemeinschaft, Landesgraduiertenförderung Baden-Württemberg and INTAS grant 99-1222.

2 Involutive Systems

We consider differential equations in n independent variables $\mathbf{x} = (x_1, \dots, x_n)$ and m dependent variables $\mathbf{u}(\mathbf{x}) = (u_1(\mathbf{x}), \dots, u_m(\mathbf{x}))$, using a multi index notation for the derivatives: $p_{\alpha,\mu} = \partial^{|\mu|} u_\alpha / \partial x_\mu = \partial^{|\mu|} u_\alpha / \partial x_{\mu_1} \cdots \partial x_{\mu_n}$ for each multi index $\mu = [\mu_1, \dots, \mu_n] \in \mathbb{N}_0^n$. The dependent variable u_α is identified with the derivative $p_{\alpha,[0,\dots,0]}$. We fix the following ranking \prec on the set of all derivatives: $p_{\alpha,\mu} \prec p_{\beta,\nu}$ if $|\mu| < |\nu|$ or if $|\mu| = |\nu|$ and the rightmost non-vanishing entry in $\nu - \mu$ is negative; if $\mu = \nu$, we set $p_{\alpha,\mu} \prec p_{\beta,\nu}$, if $\alpha < \beta$. The *class* of a derivative is the leftmost non-vanishing entry of its multi index: $\text{cls}(p_{\alpha,\mu}) := \min\{i \mid \mu_i > 0\}$ and $\text{cls}(u_\alpha) := n$. The *order* of $p_{\alpha,\mu}$ is the length of its multi index $|\mu| = \sum \mu_i$. The ranking defined above respects classes: if the derivatives $p_{\alpha,\mu}, p_{\beta,\nu}$ are of the same order but $\text{cls}(p_{\alpha,\mu}) < \text{cls}(p_{\beta,\nu})$, then $p_{\alpha,\mu} \prec p_{\beta,\nu}$. The independent variables x_i with $i \leq \text{cls}(p_{\alpha,\mu})$ are called *multiplicative* for $p_{\alpha,\mu}$, the remaining ones *non-multiplicative*.

We consider a *linear (homogeneous) differential system* $\Phi(\mathbf{x}, \mathbf{p}) = 0$ where each of the p component functions Φ_τ has the form

$$\Phi_\tau(\mathbf{x}, \mathbf{p}) = \sum_{\alpha=1}^m \sum_{0 \leq |\mu| \leq q} a_{\tau\alpha\mu}(\mathbf{x}) p_{\alpha,\mu} = 0. \quad (1)$$

The *leader* of an equation is the highest occurring derivative with respect to the ranking \prec . Concepts like class and (non-)multiplicative variables are transferred to equations by defining them in terms of their leaders. We denote by $\beta_j^{(k)}$ the number of equations of order j and class k contained in the system. Equations obtained by differentiating with respect to a (non-)multiplicative variable are called *(non-)multiplicative prolongations*.

We introduce *involutive systems* via a normal form and its properties. For simplicity, we present it only for a first order system. This poses no real restriction, as every system can be transformed into an equivalent first order one. We may thus write down (after some algebraic manipulations and, possibly, coordinate transformations) each system in its *Cartan normal form*:

$$p_{\alpha,n} - \phi_{\alpha,n}(\mathbf{x}, \mathbf{u}, p_{\gamma,j}, p_{\delta,n}) = 0, \quad \begin{cases} 1 \leq \alpha \leq \beta_1^{(n)}, \\ 1 \leq j < n, \quad \beta_1^{(n)} < \delta \leq m, \end{cases} \quad (2a)$$

$$p_{\alpha,n-1} - \phi_{\alpha,n-1}(\mathbf{x}, \mathbf{u}, p_{\gamma,j}, p_{\delta,n-1}) = 0, \quad \begin{cases} 1 \leq \alpha \leq \beta_1^{(n-1)}, \\ 1 \leq j < n-1, \quad \beta_1^{(n-1)} < \delta \leq m, \end{cases} \quad (2b)$$

⋮

$$p_{\alpha,1} - \phi_{\alpha,1}(\mathbf{x}, \mathbf{u}, p_{\delta,1}) = 0, \quad \begin{cases} 1 \leq \alpha \leq \beta_1^{(1)}, \\ \beta_1^{(1)} < \delta \leq m, \end{cases} \quad (2c)$$

$$u_\alpha - \phi_\alpha(\mathbf{x}, \mathbf{u}_\beta) = 0, \quad \begin{cases} 1 \leq \alpha \leq \beta_0^{(n)} \leq m, \\ \beta_0^{(n)} < \beta \leq m. \end{cases} \quad (2d)$$

Here, the functions $\phi_{\alpha,k}$ and ϕ_α are linear in the dependent variables and derivatives. The system is in a triangular form where the subsystem in the first line comprises all equations of class n , the one in the second line all of class $n - 1$ and so on. The derivatives on the right hand side have always a class lower than or equal to the one on the left hand side. The subsystem in the last line collects all algebraic constraints; their number is denoted by $\beta_0^{(n)}$. For an involutive system, we must have $0 \leq \beta_0^{(n)} \leq \beta_1^{(1)} \leq \dots \leq \beta_1^{(n)} \leq m$, so that the subsystems may be empty below a certain class.

We deal with a *normal system* or a system in *Cauchy–Kovalevskaya form*, if all equations are of class n , i. e. if $\beta_1^{(n)} = m$ and all other $\beta_j^{(k)}$ vanish. An existence and uniqueness theory (in the real analytic category) for such systems is provided by the famous Cauchy–Kovalevskaya theorem [21]. More generally, the system is underdetermined, if and only if $\beta_1^{(n)} < m$. If the system is not underdetermined, then the subsystem (2a) is always normal. In the sequel we will exclusively study such systems.

These purely structural aspects of the normal form (2) do not yet capture that our differential system is involutive; they only express that we have chosen a local representation in triangular form.¹ Any differential system can be brought into such a form. The important point about the Cartan normal form is that involution implies certain relations between prolonged equations. First of all, we require that any *non-multiplicative* prolongation can be written as a linear combination of *multiplicative* ones. Thus, if D_k denotes the total differentiation with respect to x_k , then functions $A_{\beta ij}(\mathbf{x})$, $B_{\beta i}(\mathbf{x})$, and $C_\beta(\mathbf{x})$ must exist such that whenever $1 \leq \ell < k \leq n$

$$D_k(p_{\alpha,\ell} - \phi_{\alpha,\ell}) = \sum_{i=1}^k \sum_{\beta=1}^{\beta_1^{(i)}} \left\{ \sum_{j=1}^i A_{\beta ij} D_j(p_{\beta,i} - \phi_{\beta,i}) + B_{\beta i}(p_{\beta,i} - \phi_{\beta,i}) \right\} + \sum_{\beta=1}^{\beta_0^{(n)}} C_\beta(u_\beta - \phi_\beta). \quad (3)$$

Furthermore, no prolongation of the algebraic equations in (2d) may lead to a new equation. This implies the existence of functions $\bar{C}_\beta(\mathbf{x})$ such that

$$\frac{\partial \phi_\alpha}{\partial x_k} - \phi_{\alpha,k} + \sum_{\beta=\beta_0^{(n)}+1}^m \frac{\partial \phi_\alpha}{\partial u_\beta} \phi_{\beta,k} = \sum_{\beta=1}^{\beta_0^{(n)}} \bar{C}_\beta(u_\beta - \phi_\beta). \quad (4)$$

We cannot go into details here, but this second set of conditions has a geometric interpretation and is only partially present in purely algebraic approaches like involutive bases. We will see in Sect. 5 that this geometric

¹ In a more algebraic terminology one may say that the equations are head autoreduced.

ingredient is crucial for obtaining correct index values. Involution is that it comprises *formal integrability*: an involutive system is always consistent and possesses at least formal power series solutions. They can be computed order by order in a straightforward manner via the Cartan normal form.

The analysis of the *compatibility conditions* will later play an important role. They define (differential) relations between the equations in (2) and correspond to *syzygies* in commutative algebra. We introduce for each equation in our system (2) an inhomogeneity $\epsilon(\mathbf{x})$ on the right hand side. The relations (3) and (4) imply that the inhomogeneous system possesses solutions, if and only if the functions ϵ satisfy the homogeneous linear system

$$\frac{\partial \epsilon_{\alpha, \ell}}{\partial x_k} - \sum_{i=1}^k \sum_{\beta=1}^{\beta_1^{(i)}} \left\{ \sum_{j=1}^i A_{\beta ij} \frac{\partial \epsilon_{\beta, i}}{\partial x_j} + B_{\beta i} \epsilon_{\beta, i} \right\} - \sum_{\beta=1}^{\beta_0^{(n)}} C_{\beta} \epsilon_{\beta} = 0, \quad (5a)$$

$$\frac{\partial \epsilon_{\alpha}}{\partial x_k} - \epsilon_{\alpha, k} + \sum_{\beta=\beta_0^{(n)}+1}^m \frac{\partial \phi_{\alpha}}{\partial u_{\beta}} \epsilon_{\beta, k} - \sum_{\beta=1}^{\beta_0^{(n)}} \bar{C}_{\beta} \epsilon_{\beta} = 0. \quad (5b)$$

Example 1. In vacuum, *Maxwell's equations* of electrodynamics are

$$\mathbf{E}_t = \text{curl } \mathbf{B}, \quad \mathbf{B}_t = -\text{curl } \mathbf{E}, \quad (6a)$$

$$0 = \text{div } \mathbf{E}, \quad 0 = \text{div } \mathbf{B}. \quad (6b)$$

We have six dependent variables $(E_1, E_2, E_3, B_1, B_2, B_3)$, the components of the electric and magnetic field, and four independent variables (x, y, z, t) . The system is almost in Cartan normal form; only the Gauß laws (6b) have not been solved for their leaders. (6a) corresponds to (2a); (6b) to (2b). Thus we find $\beta_0^{(4)} = \beta_1^{(1)} = \beta_1^{(2)} = 0$, $\beta_1^{(3)} = 2$ and $\beta_1^{(4)} = 6$.

The relations (3) follow from adding the (non-multiplicative) t -prolongation of (6b) to the (multiplicative) divergence of (6a). Due to the identity $\text{div} \circ \text{curl} = 0$, they are satisfied and Maxwell's equations are involutive. Introducing inhomogeneities $\mathbf{j}_e, \mathbf{j}_m, \rho_e, \rho_m$, we get as compatibility equations the familiar continuity equations relating charge density and current

$$(\rho_e)_t + \text{div } \mathbf{j}_e = 0, \quad (\rho_m)_t + \text{div } \mathbf{j}_m = 0. \quad (7)$$

Example 2. Involution is more than the absence of integrability conditions. Consider the following system for two unknown functions $v(x, t), w(x, t)$:

$$v_t = w_x, \quad w_t = 0, \quad v_x = 0. \quad (8)$$

It arises, if we transform the second order system $u_{tt} = u_{xx} = 0$ to a first order one. Obviously, no integrability conditions are hidden in this simple system. Thus it is formally integrable, but it is *not* involutive. Differentiating the last equation with respect to the non-multiplicative variable t yields (after a simplification) a new second order equation: $w_{xx} = 0$. Such an equation is called *obstruction to involution*; we will see later why it is rather important.

3 Existence and Uniqueness of Solutions

An important notion in the theory of differential algebraic equation is that of an *underlying equation*. It refers to an unconstrained ordinary differential system such that any solution of the given system is also a solution of it. We may straightforwardly extend this notion to partial differential equations.

Definition 1. *An underlying system of a given differential system is any normal system that is solved by any solution of the original system.*

An underlying system exists, if and only if the given system is not underdetermined. It is of course not unique. For an involutive system in Cartan normal form (2), an underlying system is given by the subsystem (2a). Thus an underlying system for Maxwell’s equations is (6a). Although (8) is not involutive, the first two equations form an underlying system.

The generalisation of the Cauchy–Kovalevskaya theorem from normal to arbitrary involutive systems is provided by the *Cartan–Kähler theorem*. It guarantees the existence of a unique analytic solution for the system (2) provided the functions $\phi_{\alpha,k}$ and the initial data are analytic. For a proof and for more information on the choice of the initial conditions we refer to [23].

In order to sketch some of the basic ideas behind the proof of the Cartan–Kähler theorem and to demonstrate how they may be used to prove the existence and uniqueness of more general solutions than only analytic ones, we consider a special class of linear systems.

Definition 2. *An involutive differential system with Cartan normal form (2) is weakly overdetermined, if $\beta_1^{(n)} = m$, $\beta_1^{(n-1)} > 0$ and $\beta_1^{(k)} = 0$ else.*

In the sequel, we are only interested in equations that can be interpreted in some sense as evolution equations, as we concentrate on initial value problems. We slightly change our notation and denote the independent variables by (x_1, \dots, x_n, t) , i. e. we have $n + 1$ variables and write x_{n+1} as t . We study linear systems with smooth coefficients of the following form:

$$\mathbf{u}_t = \sum_{i=1}^n A_i(\mathbf{x}, t) \mathbf{u}_{x_i} + B(\mathbf{x}, t) \mathbf{u}, \tag{9a}$$

$$0 = \sum_{i=1}^n C_i(\mathbf{x}, t) \mathbf{u}_{x_i} + D(\mathbf{x}, t) \mathbf{u}. \tag{9b}$$

Here \mathbf{u} is again the m -dimensional vector of dependent variables. The square matrices $A_i(\mathbf{x}, t)$ and $B(\mathbf{x}, t)$ have m rows and columns; the rectangular matrices $C_i(\mathbf{x}, t)$ and $D(\mathbf{x}, t)$ have r rows and m columns. The system is weakly overdetermined, if for at least one i we have $\text{rank } C_i(\mathbf{x}, t) = r$; without loss of generality, we may assume that this is the case for C_n .

A straightforward computation shows that (9) is involutive, if and only if $r \times r$ matrices $H_i(\mathbf{x}, t)$, $K(\mathbf{x}, t)$ exist such that for all values $1 \leq i, j \leq n$

$$H_i C_j + H_j C_i = C_i A_j + C_j A_i, \quad (10a)$$

$$H_i D + K C_i + \sum_{k=1}^n H_k C_{i,x_k} = C_i B + D A_i + \sum_{k=1}^n C_k A_{i,x_k} + C_{i,t}, \quad (10b)$$

$$K D + \sum_{k=1}^n H_k D_{x_k} = D B + \sum_{k=1}^n C_k B_{x_k} + D_t. \quad (10c)$$

Because of our assumption on rank C_n , it is not difficult to see that if such matrices H_i , K exist, they are uniquely determined by (10). We derive the compatibility conditions of the linear system (9) under the assumption that it is involutive. We add on the right hand side of (9a) a “perturbation” δ and on the right hand side of (9b) a “perturbation” $-\epsilon$. The inhomogeneous system admits (at least formal) solutions, if and only if these functions satisfy the compatibility conditions

$$\epsilon_t - \sum_{i=1}^n H_i \epsilon_{x_i} - K \epsilon = \sum_{i=1}^n C_i \delta_{x_i} + D \delta. \quad (11)$$

Recall that the system (9a) is *hyperbolic in t -direction*, if for any vector $\xi \in \mathbb{R}^n$ the matrix $A_\xi(\mathbf{x}, t) := \sum \xi_i A_i(\mathbf{x}, t)$ has at every point (\mathbf{x}, t) only real eigenvalues and an eigenbasis. It is *strongly hyperbolic*, if there exists for any $A_\xi(\mathbf{x}, t)$ a symmetric, positive definite matrix $P_\xi(\mathbf{x}, t)$, a *symmetriser*, depending smoothly on ξ , \mathbf{x} and t such that $P_\xi A_\xi - A_\xi^t P_\xi = 0$. The system (9b) is *elliptic*, if the matrix $C_\xi(\mathbf{x}, t) := \sum \xi_i C_i(\mathbf{x}, t)$ defines for any vector $0 \neq \xi \in \mathbb{R}^n$ a surjective mapping.

Example 3. Maxwell’s equations are weakly overdetermined. Their evolution part (6a), corresponding to (9a), forms a symmetric hyperbolic system, whereas the constraints (6b), corresponding to (9b), are elliptic. The compatibility conditions (11) are given by (7) with $\epsilon = (\rho_e, \rho_m)$ and $\delta = (\mathbf{j}_e, \mathbf{j}_m)$.

A classical result in the theory of hyperbolic systems [12] states that the smooth, strongly hyperbolic system

$$\mathbf{u}_t = \sum_{i=1}^n A_i(\mathbf{x}, t) \mathbf{u}_{x_i} + B(\mathbf{x}, t) \mathbf{u} + \mathbf{F}(\mathbf{x}, t) \quad (12)$$

possesses for periodic boundary conditions and smooth initial conditions $\mathbf{u}(\mathbf{x}, 0) = \mathbf{f}(\mathbf{x})$ a unique smooth solution. Furthermore, at any time $t \in [0, T]$ this solution can be estimated in the weighted Sobolev norm $\|\cdot\|_{P, H^p}$ (the weight depends on the symmetriser P_ξ , see [12] for details) with $p \geq 0$ by

$$\|\mathbf{u}(\cdot, t)\|_{P, H^p} \leq K_p \left[\|\mathbf{f}\|_{P, H^p} + \int_0^t \|\mathbf{F}(\cdot, \tau)\|_{P, H^p} d\tau \right]. \quad (13)$$

We exploit this to obtain an existence and uniqueness theorem for smooth solutions of (9). If the subsystem (9a) is strongly hyperbolic, then the cited

theorem ensures that it has a unique smooth solution for arbitrary smooth initial conditions. Let us assume that the initial data $\mathbf{f}(\mathbf{x})$ has been chosen such that they satisfy the constraints (9b) for $t = 0$. The question is whether our solution of (9a) satisfies the constraints also for $t > 0$.

The answer to this question lies in the compatibility condition (11). Entering our solution into (9b) yields residuals $\epsilon(\mathbf{x}, t)$. As we are dealing with an exact solution of (9a), the residuals ϵ satisfy (11) with $\delta \equiv 0$ so that the system becomes homogeneous. Furthermore, the choice of our initial data implies $\epsilon(\mathbf{x}, 0) = 0$, hence $\epsilon \equiv 0$ is obviously a smooth solution. We are done, if we can show that it is the only one.

Thus we need a uniqueness result for (11). If the matrices H_i, K were analytic, we could apply Holmgren's theorem [21]. However, our coefficients are only smooth. Some linear algebra shows that, provided the constraints (9b) are elliptic, the compatibility system (11), viewed as system for ϵ only, inherits the strong hyperbolicity of the underlying system (9a) with a symmetriser Q_ξ determined by $Q_\xi^{-1} = C_\xi P_\xi^{-1} C_\xi^t$ [24]. Thus we may again apply the above cited theorem to prove the needed uniqueness.

These considerations lead to a simple approach to the numerical integration of the overdetermined system (9): we consider the equations (9b) only as constraints on the initial data and otherwise ignore them, i. e. we simply solve numerically the initial value problem for (9a) with initial data satisfying (9b). This integration is a standard problem in numerical analysis.

We must expect a *drift* off the constraints, i. e. the numerical solution ceases to satisfy the constraints (9b). Some discussions of this problem for Maxwell's equations are contained in [11]. For a numerical solution the residuals δ do not vanish but lead to a "forcing term" in the compatibility condition (11). Thus the residuals ϵ do not vanish either. Their growth depends on the properties of (11). In the particular case of a strongly hyperbolic system with elliptic constraints we may estimate the size of the drift via (13):

$$\|\epsilon(\cdot, t)\|_{Q, H^p} \leq K_p \int_0^t \left\| \sum_{i=1}^n C_i \delta_{x_i}(\cdot, \tau) + D\delta(\cdot, \tau) \right\|_{Q, H^p} d\tau \quad (14)$$

This estimate depends not only on the residuals δ but also on their spatial derivatives. While any reasonable numerical method controls the size of δ , it is difficult to control the size of the derivatives.

Example 4. In the case of Maxwell's equations, the estimate (13) depends only on the divergence of δ and not on δ itself. Thus a good numerical method for them should be constructed such that this divergence vanishes.

4 Completion to Involution

If a given system is not involutive, one should *complete* it to an equivalent involutive one. "Equivalent" means here that both systems possess the same

(formal) solution space. In the case of a linear system with constant coefficients, the completion is in fact equivalent to the determination of a Gröbner basis for the corresponding polynomial module. We present now a completion algorithm for linear systems that combines algebraic and geometric ideas. More details on the algorithm can be found in [10]; an implementation in the computer algebra system *MuPAD* is briefly described in [1].

In order to formulate our algorithm, we need some further notations. We suppress the zero in (1) and call the remaining left hand side a *row*. The basic idea is to restrict computations as much as possible to non-multiplicative prolongations; multiplicative prolongations are only used for determining a triangular form of the system. In order to indicate which multiplicative prolongations are present, we introduce a *global level* λ , initialised to zero and denoting the number of the current iteration, and assign to each row one or two non-negative integers, its *initial level* and its *phantom level*. The set of such indexed rows is the *skeleton* \mathcal{S}_λ of the system and we reproduce the full system $\bar{\mathcal{S}}_\lambda$ by replacing each indexed row by a set of multiplicative prolongations determined by its indices. For a single indexed row $f_{(k)}$ these are

$$\{D_\mu f \mid 0 \leq |\mu| \leq (\lambda - k); \forall i > \text{cls}(f) : \mu_i = 0\} \quad (15)$$

and for a double indexed row $f_{(k,l)}$

$$\{D_\mu f \mid (\lambda - l) < |\mu| \leq (\lambda - k); \forall i > \text{cls}(f) : \mu_i = 0\}. \quad (16)$$

Without loss of generality, we assume that the given system of order q is already in triangular form. We turn it into the skeleton \mathcal{S}_0 by setting the initial level of each row to 0. Furthermore, the numbers $e_{\mathcal{S}_\lambda, i}$ count how many rows of order i are present in the system $\bar{\mathcal{S}}_\lambda$. Since $\mathcal{S}_0 = \bar{\mathcal{S}}_0$, the starting values $e_{\mathcal{S}_0, i}$ are obtained at once. Finally, we initialise the counter $r := 0$. Each iteration step with \mathcal{S}_λ being the current skeleton proceeds as follows.

Prolongation. The global level λ is increased by one. This automatically adds all new multiplicative prolongations to the system $\bar{\mathcal{S}}_\lambda$ defined by the skeleton \mathcal{S}_λ . Concerning the non-multiplicative prolongations, only those single indexed rows $f_{(k)}$ with $k = \lambda - 1$ (these are have been created in the last iteration) are computed and become part of the new skeleton with initial level λ . These changes necessitate to recompute the numbers $e_{\mathcal{S}_\lambda, i}$: for each row in $\mathcal{S}_{\lambda-1}$ they are modified as follows:

- If $f_{(k)}$ is a single indexed row of order t with $k = \lambda - 1$, prolongations in the direction of all independent variables are computed, so $e_{\mathcal{S}_\lambda, t+1} = e_{\mathcal{S}_{\lambda-1}, t+1} + n$.
- If the initial level of $f_{(k)}$ with order t and class j is less than λ , $\binom{(\lambda-k)+j}{j-1}$ new rows of order $t + (\lambda - k) + 1$ enter the system.
- If $f_{(k,l)}$ is a double indexed row of order t and class j , there are $\binom{(\lambda-k)+j}{j-1}$ new rows of order $t + (\lambda - k) + 1$ and $\binom{(\lambda-l)+j}{j-1}$ rows of order $t + (\lambda - l) + 1$ are removed.

Triangulation. Next, the skeleton \mathcal{S}_λ is algebraically transformed such that $\bar{\mathcal{S}}_\lambda$ is in triangular form, i. e. all rows possess different leaders. Starting with the row with the highest ranked leader, one searches through the rows of the skeleton and the allowed multiplicative prolongations for a row with the same leader. If this is the case, reductions are carried out until none are possible. Then the process is repeated with the next leader. One slight subtlety has to be watched: if a row is reduced which has already produced multiplicative prolongations, removing it would mean to lose all these rows. This is the whole reason behind the introduction of double indexed rows: by adding a phantom level $l = \lambda$ to $f_{(k)}$, we ensure that in future iterations new prolongations are still taken into account and prolongations becoming reducible are removed.

The changes of the $e_{\mathcal{S}_\lambda, i}$ are trivial: If a row is reduced, the corresponding value is decreased by 1 and, if the reduction has not yielded zero, the value at the appropriate order is increased by 1.

Involution Analysis. We now check whether we have reached an involutive system. Two conditions must be fulfilled for this:

- No single indexed rows of order $q + r + 1$ may exist in \mathcal{S}_λ .
- The values of $e_{\mathcal{S}_{\lambda-1}, i}$ and $e_{\mathcal{S}_\lambda, i}$ coincide for $1 \leq i \leq q + r$.

If the first condition is violated, we set $r := r + 1$ and proceed with the next iteration. Otherwise, the new value for r is $q' - q$, where q' is the highest order at which a row occurs in \mathcal{S}_λ . If the second condition was not satisfied, we continue with the next iteration. Otherwise, we have finished.

Note that the last iteration only serves to check the involution of the system obtained in the last but one step. Obviously, nothing “new” can happen here, since otherwise our termination conditions would not hold.

Given the triangulised skeleton \mathcal{S}_λ at the end of each iteration step, one easily determines the numbers $\beta_{q+r}^{(i)}$ of the highest order part of the Cartan normal form of the corresponding system. A single indexed row $f_{(k)}$ of order t and class c contributes to them, if and only if $\lambda - k + t \geq q + r$ and $1 \leq i \leq c$; for a double indexed row with phantom level l , it is additionally required that $\lambda - l + t < q + r$. The contribution is then given by $B(c - i + 1, q + r - t, 1)$, where $B(n, c, q) = \binom{n-c+q-1}{q-1}$ denotes the number of multi indices of length n , order q and class c .

Example 5. For the system (8) our algorithm needs only two iterations. In the first one the equation $w_{xx} = 0$ is added to the skeleton, as it cannot be reduced by a multiplicative prolongation. After the second iteration the algorithm stops with $r = 1$ as the final system contains a second order equation.

Our algorithm also works for quasi-linear systems provided they remain quasi-linear during the completion. The application to fully non-linear systems leads to a number of serious issues which we cannot discuss here.

Example 6. In order to demonstrate an important difference between our combined algebraic-geometric completion algorithm and purely algebraic methods like involutive bases we analyse the planar pendulum. The dependent variables are the positions (x, y) , the velocities (u, v) and a Lagrange multiplier λ . We set all parameters like mass or length to 1. This yields the equations of motion:

$$\dot{x} = u, \quad \dot{y} = v, \quad \dot{u} = -x\lambda, \quad \dot{v} = -y\lambda - 1, \quad 0 = x^2 + y^2 + 1. \quad (17)$$

Our algorithm needs four iterations. The first two produce the new equations $0 = xu + yv$ and $0 = u^2 + v^2 - y - \lambda$. In the third step, we find no algebraic constraint, but an equation for $\dot{\lambda}$ enters the system. Finally, in the fourth step nothing further happens, so we have arrived at an involutive system.

A purely algebraic approach would need one iteration less. It would not determine an equation for the derivative $\dot{\lambda}$, as the system contains already an algebraic equation for λ . The next section shows that this additional iteration of geometric origin is important for obtaining the correct index values.

5 Indices for Differential Equations

Indices play an important role for differential algebraic equations. They serve as indicators for the difficulties one has to expect in the numerical integration of the given system: the higher the index, the more problems arise. Typical problems are that classical methods exhibit suddenly a lower order of convergence or the already mentioned drift off the constraint manifold.

The literature abounds with definitions of indices; a partial survey may be found in [4]. A recent trend is their extension to partial differential equations. One often speaks of *partial differential algebraic equations*, but this terminology is misleading, as in many cases (like Maxwell's equations) such systems do not contain any algebraic equations: the non-normality is due to the presence of equations of lower class and not of lower order.

We distinguish two classes of indices. *Differentiation indices* count the number of prolongations needed until the system possesses certain properties. In the case of "the" differentiation index, the property is that an underlying system has been found. *Perturbation indices* are based on estimates on the difference of solutions of the original system and of a perturbation of it.

Two differentiation indices follow naturally from our completion algorithm. Recall that it produces a sequence of linear systems in triangular form² $\bar{S}_0 \longrightarrow \bar{S}_1 \longrightarrow \dots \longrightarrow \bar{S}_{\lambda_f}$.

Definition 3. *The determinacy index ν_D of a differential system in n independent and m dependent variables is the first λ such that we have $\beta_{q+r}^{(n)} = m$ for the system \bar{S}_λ and the corresponding value of the counter r . The involutive index ν_I is the first λ for which \bar{S}_λ is involutive and thus $\lambda_f - 1$.*

² The systems are taken at the *end* of each iteration step.

The determinacy index corresponds to “the” differentiation index. Its name reflects that \mathcal{S}_λ is not an underdetermined system. ν_D has no finite value for an underdetermined system, as then $\beta_{q+r}^{(n)} < m$ for any value of λ . In contrast to the involution index, the determinacy index cannot be used as a basis of an existence theory for solutions, as it does not require the full completion. Thus it is well possible that the system is in fact inconsistent. The involution index is equivalent to the strangeness index [13]; one also obtains the same result with the geometric approaches presented in [18, 19].

Example 7. It follows from Ex. 6 that both the determinacy and the involution index of the pendulum are $\nu_I = \nu_D = 3$, as an underlying system is obtained only after an equation for λ is present. A purely algebraic completion would yield too low values for the indices.

In system (8), we find different values for ν_I and ν_D . The determinacy index is obviously 0, as the first two equations form already an underlying system. But $\lambda_f = 2$ and hence $\nu_I = 1$.

The perturbation index was introduced by Hairer et al. [7] for differential algebraic equations and extended to partial differential equations by Campbell and Marszalek [5]. We assume that the system lives on a compact domain $\Omega \subset \mathbb{R}^n$ and choose a norm $\|\cdot\|$ on some function space \mathcal{F} over Ω which is usually either the maximum norm $\|\cdot\|_{L^\infty}$ or the uniform norm $\|\cdot\|_{L^1}$. In addition, we define on \mathcal{F} for each integer $k > 0$ a kind of Sobolev norm

$$\|\mathbf{f}\|_k = \sum_{0 \leq |\mu| \leq k} \left\| \frac{\partial^{|\mu|} \mathbf{f}}{\partial x_\mu} \right\|, \quad (18)$$

i. e. we sum the norms of \mathbf{f} and its partial derivatives up to order k . Of course, the function space \mathcal{F} must be such that these norms make sense on it.

Partial differential equations are usually accompanied by initial or boundary conditions. In order to accommodate for this, we take the following simple approach. Let $\Omega' \subseteq \partial\Omega$ be a subdomain of the boundary of Ω and introduce on the restriction \mathcal{F}' of \mathcal{F} to Ω' similar norms denoted by $\|\cdot\|'_k$. The conditions are assumed to be of the form $\Psi(\mathbf{x}, \mathbf{p})|_{\Omega'} = 0$. This comprises most kinds of initial or boundary conditions in applications. The highest derivative in Ψ determines the order ℓ of the conditions.

We consider again a linear partial differential system with p equations (1). We do not require that $p = m$ or that the Jacobian $\partial\Phi/\partial\mathbf{p}$ has any special properties. Assume that we are given a smooth solution $\mathbf{u}(\mathbf{x})$ defined on the whole domain Ω and satisfying our initial or boundary conditions on Ω' . We compare it with solutions of the perturbed equation $\Phi(\mathbf{x}, \mathbf{p}) = \delta(\mathbf{x})$ with a smooth right hand side δ .

Definition 4. Let $\mathbf{u}(\mathbf{x})$ be a smooth solution satisfying some initial or boundary conditions of order ℓ on Ω' . The system has perturbation index ν_P along

this solution, if ν_P is the smallest integer such that for any solution $\tilde{\mathbf{u}}(\mathbf{x})$ of the perturbed equation at every point $\mathbf{x} \in \Omega$ the estimate

$$|\mathbf{u}(\mathbf{x}) - \tilde{\mathbf{u}}(\mathbf{x})| \leq C(\|\mathbf{f} - \tilde{\mathbf{f}}\|'_\ell + \|\delta\|_{\nu_P-1}) \quad (19)$$

holds, whenever the right hand side is sufficiently small. Here \mathbf{f} and $\tilde{\mathbf{f}}$ represent the restrictions of the solutions \mathbf{u} and $\tilde{\mathbf{u}}$, respectively, to Ω' . The constant C may depend only on the domains Ω , Ω' and on the function Φ .

In the case of initial value problems for differential algebraic equations this definition coincides with the one given by Hairer et al. The term $\|\mathbf{f} - \tilde{\mathbf{f}}\|'_k$ takes there the simple form $|\mathbf{u}(\mathbf{x}_0) - \tilde{\mathbf{u}}(\mathbf{x}_0)|$ where $\Omega' = \{\mathbf{x}_0\}$. For partial differential systems our definition is almost identical with the one given by Campbell and Marszalek; the only difference is that we do not include the order ℓ of the initial or boundary conditions in the index.

Whereas it is not so clear why differentiation indices should indicate the difficulty of the numerical integration of the system, this is rather obvious for the perturbation index. If we take for $\tilde{\mathbf{u}}(\mathbf{x})$ an approximate solution, we may interpret $\delta(\mathbf{x})$ as the residual obtained by entering it into the system. The estimate (19) tells us that for an equation with $\nu_P > 1$ it does not suffice to keep this residual as small as possible, since also some of its derivatives enter it. Strictly speaking, this implies that the considered initial or boundary value problem is ill-posed in the sense of Hadamard!

Obviously, it is much harder to obtain estimates of the form (19) than to compute a differentiation index, but the perturbation index contains more useful information. Hence there is much interest in relating the two concepts.

Conjecture 1. For any linear differential system $\nu_D \leq \nu_P \leq \nu_D + 1$.

For differential algebraic equations, a rigorous proof of this conjecture can be found in [9].³ One first shows that a normal ordinary differential system has the perturbation index $\nu_P \leq 1$ (a consequence of Gronwall's Lemma). For a general system, one follows the completion algorithm until an underlying system is reached. One can prove that its right hand side contains derivatives of the perturbations of order ν_D . This yields the estimate above.

For partial differential systems, the situation is much more complicated, as the perturbation index will depend in general on the chosen norm. A simple case arises, if the underlying equation can be treated with semi-group theory [21]. Then we may consider our overdetermined system as a differential algebraic equation on an infinite-dimensional Banach space and the same argument as above may be applied. Thus we obtain the same estimate.

³ In that article non-linear systems are treated where one must introduce *perturbed* differentiation indices. For linear systems they are identical with the above defined indices.

6 Semi-Discretisation

For simplicity, we restrict to first order linear homogeneous systems with constant coefficients in $n + 1$ independent variables of the form

$$\sum_{i=1}^{n+1} M_i \mathbf{u}_{x_i} + N \mathbf{u} = 0. \quad (20)$$

The matrices M_i and N are here completely arbitrary. Now we discretise the derivatives with respect to n of the independent variables by some finite difference method. This yields a differential algebraic equation in the one remaining independent variable. We are interested in the relation between the involution indices of the original partial differential system and the obtained differential algebraic equation, respectively.

Instead of (20) we complete to involution a *perturbed* system with a generic right hand side $\gamma(\mathbf{x})$. For a linear constant coefficients system, the perturbation does not affect the completion. It only serves as a convenient mean of “book-keeping” which prolongations are needed during the completion.

Definition 5. *The involution index ν_ℓ in direction x_ℓ of the system (20) is the maximal number of differentiations with respect to x_ℓ in a γ -derivative contained in the involutive completion of the perturbed system.*

For the semi-discretisation we proceed as follows. Assume that x_ℓ is the “surviving” independent variable; in other words, afterwards we are dealing with an ordinary differential system containing only x_ℓ -derivatives. In order to simplify the notation, we rewrite (20). We denote x_ℓ by t and renumber the remaining independent variables as x_i with $1 \leq i \leq n$. Then we solve as many equations as possible for a t -derivative, as we consider these as the derivatives of highest class. This yields a system of the form

$$E \mathbf{u}_t = \sum_{i=1}^n A_i \mathbf{u}_{x_i} + B \mathbf{u} + \delta, \quad 0 = \sum_{i=1}^n C_i \mathbf{u}_{x_i} + D \mathbf{u} + \epsilon. \quad (21)$$

Here we assume that the (not necessarily square) matrix E is of maximal rank, i. e. every equation in the first subsystem really depends on a t -derivative. But we do not pose any restrictions on the ranks of the matrices C_i .

In (21) we have introduced perturbations δ and ϵ which are related to the original perturbations γ by a linear transformation with constant coefficients. As we are only interested in the number of differentiations applied to them during the completion, such a transformation has no effect.

We discretise the “spatial” derivatives, i. e. those with respect to the x_i , on a grid where the points $\mathbf{x}_\mathbf{k}$ are labelled by integer vectors $\mathbf{k} = [k_1, \dots, k_n]$. $\mathbf{u}_\mathbf{k}(t)$ denotes the value of the function \mathbf{u} at the point $\mathbf{x}_\mathbf{k}$ at time t . We approximate in (21) the spatial derivative $\mathbf{u}_{x_i}(\mathbf{x}_\mathbf{k}, t)$ by the finite difference

$$\delta_i \mathbf{u}_{\mathbf{k}}(t) = \sum_{\ell_i = -a_i}^{b_i} \alpha_{\ell_i}^{(i)} \mathbf{u}_{[k_1, \dots, k_i + \ell_i, \dots, k_n]}(t) \quad (22)$$

with some real coefficients $\alpha_{\ell_i}^{(i)}$. Thus different discretisations are allowed for different values of i , but \mathbf{u}_{x_i} is everywhere discretised in the same way. Entering the approximations (22) into (21) yields

$$E \dot{\mathbf{u}}_{\mathbf{k}} = \sum_{i=1}^n A_i \delta_i \mathbf{u}_{\mathbf{k}} + B \mathbf{u}_{\mathbf{k}}, \quad 0 = \sum_{i=1}^n C_i \delta_i \mathbf{u}_{\mathbf{k}} + D \mathbf{u}_{\mathbf{k}}. \quad (23)$$

Theorem 1. *The involution index of the differential algebraic equation (23) obtained in the described semi-discretisation of (20) with respect to x_ℓ is ν_ℓ .*

The proof is given in [24]; we outline here only the basic idea. We compare what happens during the completion of the perturbed system (21) and the differential algebraic equation (23). One can show that each new equation in the discretised system corresponds to either an integrability condition or an obstruction to involution in the original system. Further examination shows that this only happens for prolongations in t -direction; mere spatial differentiations do not lead to new equations in the differential algebraic system. Since ν_ℓ counts only the number of differentiations with respect to t , the involution index of (23) must coincide with ν_ℓ .

This result is somewhat surprising. While one surely expects that integrability conditions of the original partial differential system induce integrability conditions in the differential algebraic equation obtained by semi-discretisation, Theorem 1 says that obstructions to involution also turn into integrability conditions upon semi-discretisation. Thus even if the original partial differential system is formally integrable, the differential algebraic equation might contain integrability conditions.

Example 8. A semi-discretisation with backward differences for the spatial derivatives of the linear system (8) leads to the differential algebraic equation

$$\dot{v}_n = (w_n - w_{n-1})/\Delta x, \quad \dot{w}_n = 0, \quad v_n - v_{n-1} = 0. \quad (24)$$

It hides the integrability condition $w_n - 2w_{n-1} + w_{n-2} = 0$ obviously representing a discretisation of the obstruction to involution $w_{xx} = 0$ by centred differences. The involution index of (8) in direction t is one which is also the involution index of (24) in agreement with Theorem 1.

For weakly overdetermined systems, Theorem 1 can be strengthened. For them the differential algebraic equation obtained by a semi-discretisation with respect to t is formally integrable, *if and only if* the original partial differential system is involutive [24]. This result does not only hold for semi-discretisations by finite differences but also for spectral methods: assuming periodic boundary conditions, we make the Fourier ansatz

$$\mathbf{u}(\mathbf{x}, t) \approx \sum_{\mathbf{k} \in \mathcal{G}} [\mathbf{a}_{\mathbf{k}}(t) + i\mathbf{b}_{\mathbf{k}}(t)] e^{i\mathbf{k}\mathbf{x}}. \quad (25)$$

Here \mathcal{G} is a finite grid of wave vectors, and we split the complex Fourier coefficients into their real and imaginary part. Entering this ansatz into (9) yields the following differential algebraic equation where the vectors \mathbf{A} and \mathbf{C} consist of the matrices A_i and C_i , respectively:

$$\begin{pmatrix} \dot{\mathbf{a}}_{\mathbf{k}} \\ \dot{\mathbf{b}}_{\mathbf{k}} \end{pmatrix} = \begin{pmatrix} B & -\mathbf{k} \cdot \mathbf{A} \\ \mathbf{k} \cdot \mathbf{A} & B \end{pmatrix} \begin{pmatrix} \mathbf{a}_{\mathbf{k}} \\ \mathbf{b}_{\mathbf{k}} \end{pmatrix}, \quad 0 = \begin{pmatrix} D & -\mathbf{k} \cdot \mathbf{C} \\ \mathbf{k} \cdot \mathbf{C} & D \end{pmatrix} \begin{pmatrix} \mathbf{a}_{\mathbf{k}} \\ \mathbf{b}_{\mathbf{k}} \end{pmatrix}. \quad (26)$$

One can show that this differential algebraic equation is formally integrable, *if and only if* the weakly overdetermined system (9) is involutive [24].

References

- [1] J. Belanger, M. Hausdorf, and W.M. Seiler. A *MuPAD* library for differential equations. In V.G. Ghanza, E.W. Mayr, and E.V. Vorozhtsov, editors, *Computer Algebra in Scientific Computing — CASC 2001*, pages 25–42. Springer-Verlag, Berlin, 2001.
- [2] K.E. Brenan, S.L. Campbell, and L.R. Petzold. *Numerical Solution of Initial-Value Problems in Differential-Algebraic Equations*. Classics in Applied Mathematics 14. SIAM, Philadelphia, 1996.
- [3] J. Calmet, M. Hausdorf, and W.M. Seiler. A constructive introduction to involution. In R. Akerkar, editor, *Proc. Int. Symp. Applications of Computer Algebra — ISACA 2000*, pages 33–50. Allied Publishers, New Delhi, 2001.
- [4] S.L. Campbell and C.W. Gear. The index of general nonlinear DAEs. *Numer. Math.*, 72:173–196, 1995.
- [5] S.L. Campbell and W. Marszalek. The index of an infinite dimensional implicit system. *Math. Model. Syst.*, 1:1–25, 1996.
- [6] V.P. Gerdt and Yu.A. Blinkov. Involutive bases of polynomial ideals. *Math. Comp. Simul.*, 45:519–542, 1998.
- [7] E. Hairer, C. Lubich, and M. Roche. *The Numerical Solution of Differential-Algebraic Equations by Runge-Kutta Methods*. Lecture Notes in Mathematics 1409. Springer-Verlag, Berlin, 1989.
- [8] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II*. Springer Series in Computational Mathematics 14. Springer-Verlag, Berlin, 1996.
- [9] M. Hausdorf and W.M. Seiler. Perturbation versus differentiation indices. In V.G. Ghanza, E.W. Mayr, and E.V. Vorozhtsov, editors, *Computer Algebra in Scientific Computing — CASC 2001*, pages 323–337. Springer-Verlag, Berlin, 2001.

- [10] M. Hausdorf and W.M. Seiler. An efficient algebraic algorithm for the geometric completion to involution. *Appl. Alg. Eng. Comm. Comp.*, 13: 163–207, 2002.
- [11] B.N. Jiang, J. Wu, and L.A. Povelli. The origin of spurious solutions in computational electrodynamics. *J. Comp. Phys.*, 125:104–123, 1996.
- [12] H.-O. Kreiss and J. Lorenz. *Initial-Boundary Value Problems and the Navier-Stokes Equations*. Pure and Applied Mathematics 136. Academic Press, Boston, 1989.
- [13] P. Kunkel and V. Mehrmann. Canonical forms for linear differential-algebraic equations with variable coefficients. *J. Comp. Appl. Math.*, 56: 225–251, 1994.
- [14] G. Le Vey. Some remarks on solvability and various indices for implicit differential equations. *Num. Algo.*, 19:127–145, 1998.
- [15] W. Lucht, K. Strehmel, and C. Eichler-Liebenow. Indexes and special discretization methods for linear partial differential algebraic equations. *BIT*, 39:484–512, 1999.
- [16] Y.O. Macutan and G. Thomas. Theory of formal integrability and DAEs: Effective computations. *Num. Algo.*, 19:147–157, 1998.
- [17] J.F. Pommaret. *Systems of Partial Differential Equations and Lie Pseudogroups*. Gordon & Breach, London, 1978.
- [18] P.J. Rabier and W.C. Rheinboldt. A geometric treatment of implicit differential algebraic equations. *J. Diff. Eq.*, 109:110–146, 1994.
- [19] S. Reich. On an existence and uniqueness theory for nonlinear differential-algebraic equations. *Circ. Sys. Sig. Proc.*, 10:343–359, 1991.
- [20] G.J. Reid, P. Lin, and A.D. Wittkopf. Differential elimination-completion algorithms for DAE and PDAE. *Stud. Appl. Math.*, 106: 1–45, 2001.
- [21] M. Renardy and R.C. Rogers. *An Introduction to Partial Differential Equations*. Texts in Applied Mathematics 13. Springer-Verlag, New York, 1993.
- [22] W.M. Seiler. Indices and solvability for general systems of differential equations. In V.G. Ghanza, E.W. Mayr, and E.V. Vorozhtsov, editors, *Computer Algebra in Scientific Computing — CASC 1999*, pages 365–385. Springer-Verlag, Berlin, 1999.
- [23] W.M. Seiler. Involution — the formal theory of differential equations and its applications in computer algebra and numerical analysis. Habilitation thesis, Dept. of Mathematics, Universität Mannheim, 2001.
- [24] W.M. Seiler. Completion to involution and semi-discretisations. *Appl. Num. Math.*, 42:437–451, 2002.
- [25] J. Tuomela and T. Arponen. On the numerical solution of involutive ordinary differential systems. *IMA J. Num. Anal.*, 20:561–599, 2000.