

## 4

# Two Topics in Nonlinear Stability

J. M. Sanz-Serna

### 4.1 Introduction

There are several ways in which the word stability may be understood in a numerical analysis context. First of all, stability is used in expressions like 'stability and consistency imply convergence'. In this first sense, stability refers to dependence of a numerical result on the data, and, as such, applies to most numerical computations, including numerical linear algebra, quadrature, differential equations, etc. This first notion is akin to the idea of 'well posedness'. In fact, in many cases, a numerical procedure is said to be stable in this sense if it is well posed uniformly with respect to the relevant parameters, such as the dimension of the problem, grid-size, etc. In a second use, the term stability applies in connection with the long time behaviour of discretizations of time-dependent problems in ordinary or partial differential equations. The stability of discretizations of nonlinear differential equations, ordinary or partial, is the unifying theme of the present work. In Sections 4.2 to 4.8, we deal with the question of how best to define 'stability' of a nonlinear discretization so that the familiar 'stability and consistency imply convergence' holds. We present a definition of stability, that when combined with an important lemma due to Stetter and with a suitable linearization theorem, has revealed itself to be very helpful in proving the convergence of nonlinear finite-difference, spectral and Galerkin methods. In Sections 4.9–4.11 we are concerned with stability in the second sense. We show how to employ dynamical system results to investigate the stability of numerical methods that in a linear analysis are neutrally stable. Our study leads in a natural way to the consideration of symplectic or canonical integrators, a subject briefly surveyed in the final Section 4.12.

### 4.2 Discretizations

In order to study the idea of stability, it is advisable to work in an abstract framework, so that the attention may be focused on the key issues

and a number of application fields can be covered at once. Here we work within the framework used in Sanz-Serna[29], that, while being general enough to treat most application areas, is not unnecessarily abstract. In this section we summarize some of the basic ideas of Sanz-Serna[29]. The reader is referred to the original paper for a more detailed treatment and for information on other available general frameworks.

We consider a given problem involving a linear or nonlinear differential equation supplemented by suitable side conditions, such as boundary conditions, initial conditions, etc. Let  $u$  denote a solution of this problem (well-behaved nonlinear problems may of course possess more than one solution). For the numerical solution we use the notation  $U_h$ , where the subscript  $h$  reflects the dependence of  $U_h$  on a (small) parameter such as a mesh-size, element diameter, reciprocal of number of terms retained when truncating a series, etc. We always assume that  $h$  takes values in a set  $H$  of positive numbers with  $\inf H = 0$ . The numerical approximation  $U_h$  is reached by solving a *discretized problem*

$$\Phi_h(U_h) = 0, \quad (2.1)$$

where, for each  $h$  in  $H$ ,  $\Phi_h$  is a mapping with domain  $D_h$  and taking values in  $Y_h$ . Here  $Y_h$  is a vector space and  $D_h$  is a subset of a vector space  $X_h$  with

$$\dim(X_h) = \dim(Y_h) < \infty. \quad (2.2)$$

It is typical of nonlinear situations that  $\Phi_h$  cannot be defined everywhere in  $X_h$ : the analytic expression of  $\Phi_h$  may involve logarithms, square roots, etc., which only make sense for vectors in a set  $D_h$  smaller than  $X_h$ . As  $h$  ranges in  $H$ , the family of discrete problems (2.1) is called a *discretization*.

Most discretizations used in practice for stationary and time-dependent problems may readily be cast in the format (2.1). This applies not only to finite-difference techniques, but also to Galerkin, collocation, spectral methods, etc., (more on this later).

When a solution  $U_h$  of (2.1) has been obtained, the question arises as to what extent does  $U_h$  provide a good approximation to  $u$ . A first difficulty in answering this question stems from the fact that  $U_h$  can be completely dissimilar to  $u$ . Typically, in finite differences,  $U_h$  is a vector with, say,  $d$  entries, while  $u$  is a function of one or several continuous variables. This difficulty is circumvented as follows. Since any solution  $U_h$  of (2.1) is bound to be an element in  $D_h$ , we first make up our minds as to which element  $u_h$  in  $D_h$  should be regarded as the most desirable numerical result. Typically in finite differences  $u_h$  contains  $d$  nodal values of  $u$ . Once  $u_h$  has been chosen, the vector  $e_h = u_h - U_h$  is defined to be the (*global*) *error* in  $U_h$ . In order to measure the size of  $e_h$ , we introduce, for each  $h$  in  $H$ , a norm  $\|\cdot\|$  in  $X_h$ . (Norms in different spaces will simply be denoted by  $\|\cdot\|$  without

mention of the space.) We say that the discretization (2.1) is *convergent* if there exists  $h^* > 0$ , such that for  $h$  in  $H$ ,  $h < h^*$ , (2.1) possesses a solution  $U_h$  such that  $\lim \|u_h - U_h\| = 0$ ,  $h \rightarrow 0$ . If, furthermore,  $\|e_h\| = O(h^p)$ ,  $h \rightarrow 0$ , then the convergence is said to be of order  $p$ .

The (local) *discretization error*  $\tau_h$  in  $u_h$  is, by definition, the element  $\tau_h = \Phi_h(u_h) \in Y_h$ , i.e., the residual by which the element  $u_h$  fails to satisfy the discrete equations. The measurement of the size of  $\tau_h$  requires, therefore, the introduction, for each  $h$  in  $H$ , of a norm  $\|\cdot\|$  in  $Y_h$ . When these norms have been chosen, (2.1) is said to be consistent (resp. consistent of order  $p$ ) if, as  $h \rightarrow 0$ ,  $\|\tau_h\| \rightarrow 0$  (resp.  $\|\tau_h\| = O(h^p)$ ).

Before we review how convergence is obtained from consistency and stability, it is convenient to illustrate the set-up described above with an example. In the interest of clarity, the example refers to a simple finite-difference scheme. More interesting finite-difference discretizations have been treated within our framework in López-Marcos and Sanz-Serna[20], Frutos and Sanz-Serna[10], Ortega and Sanz-Serna[23]. For examples of the application to Galerkin methods see López-Marcos and Sanz-Serna [21], Süli[44], Murdoch and Budd[22] and for spectral and pseudospectral techniques see Frutos and Sanz-Serna[9], Frutos, Ortega and Sanz-Serna[11], Abia and Sanz-Serna[1].

**Example A.** Consider the following periodic initial-value reaction-diffusion problem

$$u_t = u_{xx} + f(u), \quad -\infty < x < \infty, \quad 0 \leq t \leq T < \infty, \quad (2.3)$$

$$u(x+1, t) = u(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T < \infty, \quad (2.4)$$

$$u(x, 0) = u^0(x), \quad -\infty < x < \infty. \quad (2.5)$$

In (2.3),  $f$  is a smooth real function of the real variable  $u$ ,  $-\infty < u < \infty$ . In (2.5),  $u^0$  is a given real 1-periodic function and it is assumed that  $f$ ,  $T$  and  $u^0$  are such that (2.3)–(2.5) possesses a unique smooth solution up to  $t = T$ .

To set up the numerical scheme, choose a positive constant  $r$  (the mesh-ratio) and an integer  $J > 2$ . Set  $h = 1/J$  and consider the grid-points  $x_j = jh$ ,  $j$  integer, and the time levels  $t_n = nk$ ,  $k = rh^2$ ,  $n = 0, \dots, N = [T/k]$ . For  $j = 1, \dots, J$  and  $n = 0, \dots, N-1$  set

$$\frac{U_j^{n+1} - U_j^n}{k} - \frac{U_{j-1}^n - 2U_j^n + U_{j+1}^n}{h^2} - f(U_j^n) = 0, \quad (2.6)$$

where it is obviously understood that  $U_0^n = U_J^n$  and  $U_{J+1}^n = U_1^n$ . For  $j = 1, \dots, J$  set

$$U_j^0 - u_0(x_j) = 0. \quad (2.7)$$

Formulae (2.6)–(2.7) are cast in the format (2.1) as follows. Let  $Z_h$  denote the vector space of grid functions  $\mathbf{U} = [U_1, \dots, U_J]$  defined on  $\{x_j : 1 \leq j \leq J\}$ . For each  $n$ , all the numerical approximations  $U_j^n$  associated with the time level  $t_n$  form a vector  $\mathbf{U}^n$  in  $Z_h$ . Thus (2.6)–(2.7) may be rewritten

$$\mathbf{U}^0 - \mathbf{u}^0 = \mathbf{0}, \quad \mathbf{u}^0 = [u^0(x_1), \dots, u^0(x_J)], \quad (2.8)$$

$$\frac{\mathbf{U}^{n+1} - \mathbf{U}^n}{k} - D^2 \mathbf{U}^n - \mathbf{f}(\mathbf{U}^n) = \mathbf{0}, \quad n = 0, \dots, N-1, \quad (2.9)$$

where  $D^2$  is the standard matrix replacement of the second derivative operator with periodic boundary conditions and the notation  $\mathbf{f}(\mathbf{U}^n)$  is self-explanatory. Next choose  $X_h = D_h = Y_h$  equal to the product of  $N+1$  copies  $Z_h \times \dots \times Z_h$ . Thus  $U_h := [\mathbf{U}^0, \dots, \mathbf{U}^N]$  is a vector in  $X_h$  and (2.8)–(2.9) are clearly of the form (2.1) for a suitable choice of  $\Phi_h$ . For  $u_h$  the obvious choice is given by the vector of grid restrictions  $[\mathbf{u}_0, \dots, \mathbf{u}_N]$  of  $u$ . In  $Z_h$  we use the maximum norm and in  $X_h$  we use the  $L_\infty(L_\infty)$  norm

$$\|[\mathbf{V}^0, \dots, \mathbf{V}^N]\| = \max_n \|\mathbf{V}^n\|, \quad [\mathbf{V}^0, \dots, \mathbf{V}^N] \in X_h. \quad (2.10)$$

With this norm, convergence in the sense of the abstract framework means uniform convergence in time of the maximum spatial norm. The local discretization error  $\tau_h$  is of the form  $[\tau^0, \dots, \tau^N]$ , where, according to (2.8),  $\tau^0 = \mathbf{0}$ , while for  $n = 0, \dots, N-1$ ,  $\tau^{n+1}$  contains the familiar truncation errors of the formula (2.6), i.e., the value of the left hand side of (2.6) after replacing each ‘numerical’  $U_j^n$  by its ‘theoretical’ counterpart  $u(x_j, t_n)$ . Standard Taylor expansions show that each component of  $\tau^n$  is  $O(h^2 + k)$ . If in  $Y_h$  we use the  $L_1(L_\infty)$  norm

$$\|[\rho^0, \dots, \rho^N]\| = \|\rho^0\| + \sum_{n=1}^N k \|\rho^n\|, \quad [\rho^0, \dots, \rho^N] \in Y_h, \quad (2.11)$$

the  $O(h^2 + k)$  behaviour of the components of  $\tau^n$  yields

$$\|\tau_h\| = O(h^2 + k) = O(h^2 + \tau h^2) = O(h^2),$$

i.e., second order of consistency.

### 4.3 Stability in the linear case

The general idea of stability will first be presented in the context of linear discretizations. Let us suppose that (2.1) takes the linear form

$$\Phi_h(U_h) \equiv \Psi_h U_h - g_h = 0, \quad (3.1)$$

i.e.

$$\Psi_h U_h = g_h, \quad (3.2)$$

where, for each  $h$  in  $H$ ,  $\Psi_h$  is a linear operator (matrix) mapping  $X_h$  into  $Y_h$  and  $g_h$  is a fixed vector in  $Y_h$ . The discretization (3.1)–(3.2) is stable with respect to the chosen norms in  $X_h$  and  $Y_h$ , if there exist positive constants  $S$  (the stability constant) and  $h_0$  such that for each  $h$  in  $H$ ,  $h \leq h_0$ , and for each  $V_h$  in  $X_h$

$$\|V_h\| \leq S \|\Psi_h V_h\|. \quad (3.3)$$

In other words stability represents an *a priori* bound for the solutions of (3.2) with a constant  $S$  that must be independent of  $h$ . The bound (3.3) can be used in two ways:

- (i) For  $h$  fixed,  $h \leq h_0$ , (3.3) implies that the kernel of  $\Psi_h$  is trivial. This combined with (2.2) reveals that the solution  $U_h$  of (3.1)–(3.2) exists and is unique. The unique existence of the numerical solution is obvious in some cases, e.g., explicit algorithms in initial value problems, but of course cannot in general be taken for granted. We emphasize that, in the linear case, existence and uniqueness of  $U_h$  follow from stability.
- (ii) For  $h \leq h_0$ , we may write

$$\|u_h - U_h\| \leq S \|\Psi_h(u_h - U_h)\| = S \|\Phi_h(u_h) - \Phi_h(U_h)\|, \quad (3.4)$$

which according to (2.1) and the definition of local discretization error, implies

$$\|e_h\| = \|u_h - U_h\| \leq S \|\Phi_h(u_h) - \Phi_h(U_h)\| = S \|\tau_h\|. \quad (3.5)$$

This bounds the global error in terms of local discretization error and shows that “consistent (of order  $p$ ) + stable  $\implies$  convergent (of order  $p$ )”.

**Example A** (revisited). Consider (2.3)–(2.7) with  $f \equiv 0$  and let us check whether (3.3) holds. If  $V_h = [\mathbf{V}^0, \dots, \mathbf{V}^N] \in X_h$  and  $\Phi_h(V_h) = [\boldsymbol{\rho}^0, \dots, \boldsymbol{\rho}^N] \in Y_h$ , then

$$\mathbf{V}^0 = \boldsymbol{\rho}^0, \quad (3.6)$$

$$\frac{\mathbf{V}^{n+1} - \mathbf{V}^n}{k} - D^2 \mathbf{V}^n = \boldsymbol{\rho}^{n+1}, \quad n = 0, \dots, N-1. \quad (3.7)$$

It is expedient to rewrite (3.7) in the form

$$\mathbf{V}^{n+1} = C_h \mathbf{V}^n + k \boldsymbol{\rho}^{n+1}, \quad (3.8)$$

where

$$C_h = I + kD^2, \quad (3.9)$$

is the 'transition matrix', i.e. the matrix that in (2.9) with  $f \equiv 0$  maps  $U^n$  into  $U^{n+1}$ . Recursion in (3.8) shows that for  $n = 0, \dots, N$ ,

$$V^n = C_h^n V^0 + kC_h^{n-1} \rho^1 + kC_h^{n-2} \rho^2 + \dots + k\rho^n,$$

so that, by (2.10)–(2.11) and (3.6)

$$\|V_h\| \leq \left\{ \max_{0 \leq n \leq N} \|C_h^n\| \right\} \|\Psi_h(V_h)\|, \quad (3.10)$$

where  $\|C_h^n\|$  is the operator norm for linear transformations in  $Z_h$ . It is easy to check that the constant in brackets in the bound (3.10) is the best possible. Therefore stability in the sense of the framework with (best) stability constant  $S$  is equivalent to the following Lax-stability requirement (boundedness of powers of the transition matrix)

$$S := \max\{\|C_h^n\| : h \in H, h \leq h_0, 0 \leq n \leq N\}, \quad (3.11)$$

(cf. Palencia and Sanz-Serna[24]). It is well known that (3.11) holds if and only if  $r \leq 1/2$ , and then  $\|C_h^n\| = 1$  for all  $h$  and  $n$ . Note that the equivalence between stability in the abstract sense and Lax stability is quite general in that it holds for (3.8) independently of the specific nature of  $Z_h$  and  $C_h$ . (The equivalence hinges on the  $L_\infty(Z_h)$  and  $L_1(Z_h)$  choice of norms for  $X_h$  and  $Y_h$ , see Sanz-Serna and Palencia[35].)

#### 4.4 Nonlinear stability

When trying to decide how to modify the linear stability definition to cater for nonlinear cases, the key observation is of course that we would like (3.5) to be valid, i.e., in the linear case (3.3) is applied to *differences* of vectors in  $X_h$ . Therefore, it is natural to define (2.1) to be stable if positive constants  $h_0$  and  $S$  exist such that for each  $h$  in  $H$ ,  $h \leq h_0$ , and each pair  $V_h, W_h$  of vectors in  $D_h$

$$\|V_h - W_h\| \leq S \|\Phi_h(V_h) - \Phi_h(W_h)\|. \quad (4.1)$$

Clearly, in linear cases, this is equivalent to (3.3). We shall say that discretizations stable in the sense of this definition are N-stable (N for natural or for naive, according to your preferences). If the existence of  $U_h$  is obvious or has been proved in some way, then "consistent (of order  $p$ ) + N-stable  $\implies$  convergent (of order  $p$ )". The trivial proof is again given by (3.5).

**Example A** (revisited). Consider (2.3)–(2.7) with  $r \leq 1/2$  and  $f$  globally Lipschitz, i.e.

$$|f(v) - f(w)| \leq L|v - w|, \quad (4.2)$$

for all  $v, w \in \mathbb{R}$ . Then, if  $V_h = [\mathbf{V}^0, \dots, \mathbf{V}^N]$ ,  $W_h = [\mathbf{W}^0, \dots, \mathbf{W}^N]$ ,  $\Phi_h(V_h) = [\rho^0, \dots, \rho^N]$ ,  $\Phi_h(W_h) = [\sigma^0, \dots, \sigma^N]$ ,

$$\mathbf{V}^0 - \mathbf{u}^0 = \rho^0, \quad (4.3)$$

$$\frac{\mathbf{V}^{n+1} - \mathbf{V}^n}{k} - D^2 \mathbf{V}^n - \mathbf{f}(\mathbf{V}^n) = \rho^{n+1}, \quad n = 0, \dots, N-1. \quad (4.4)$$

$$\mathbf{W}_0 - \mathbf{u}^0 = \sigma^0, \quad (4.5)$$

$$\frac{\mathbf{W}^{n+1} - \mathbf{W}^n}{k} - D^2 \mathbf{W}^n - \mathbf{f}(\mathbf{W}^n) = \sigma^{n+1}, \quad n = 0, \dots, N-1. \quad (4.6)$$

Subtract (4.6) from (4.4) and use (3.9), to obtain, for  $n = 0, \dots, N-1$ ,

$$\mathbf{V}^{n+1} - \mathbf{W}^{n+1} = C_h(\mathbf{V}^n - \mathbf{W}^n) + k[\mathbf{f}(\mathbf{V}^n) - \mathbf{f}(\mathbf{W}^n)] + k[\rho^{n+1} - \sigma^{n+1}], \quad (4.7)$$

(cf. (3.8)). Note that (4.2) implies

$$\|\mathbf{f}(\mathbf{V}^n) - \mathbf{f}(\mathbf{W}^n)\| \leq L\|\mathbf{V}^n - \mathbf{W}^n\|, \quad (4.8)$$

so that, since  $\|C_h\| = 1$ , (4.7) yields

$$\|\mathbf{V}^{n+1} - \mathbf{W}^{n+1}\| \leq (1 + kL)\|\mathbf{V}^n - \mathbf{W}^n\| + k\|\rho^{n+1} - \sigma^{n+1}\|, \quad (4.9)$$

and a standard recursion leads to (4.1) with  $S = \exp(LT)$ . Thus, for  $r \leq 1/2$ , and  $f$  globally Lipschitz, the scheme is N-stable and hence convergent.

The argument used in this example to show N-stability is essentially identical to the argument frequently used to show that one or multistep ordinary differential equation discretizations for  $y' = f(y)$  are N-stable if they are N-stable as applied to  $y' = 0$  (0-stability see e.g. Hairer et al.[14], Section III.4).

## 4.5 Stability thresholds

As pointed out before, the abstract N-stability definition (4.1) includes, as particular cases, the notions of Lax-stability of linear initial-value problems in partial differential equations and of 0-stability in numerical ordinary differential equations. In these two application areas the idea of stability is well understood and frequently invoked. However, for nonlinear partial differential equations there is no general notion of stability that is commonly invoked to prove convergence. This may be due to the fact that the theory in Section 4.4 suffers from some drawbacks. Firstly, the question of existence of discrete solutions, which in the linear case is implied by stability, was not addressed. Secondly the scope of application of (4.1) is too restrictive. Let us illustrate this second point in the context of Example A.

**Example A** (revisited). The global Lipschitz condition used in the proof of N-stability is so demanding that few functions of interest satisfy it. In fact in reaction-diffusion problems  $f(u)$  is often a polynomial and (4.2) does not hold. Now it is not difficult to see that, when  $f$  is not globally Lipschitz, the scheme (2.8)–(2.9) is not N-stable. In fact, set  $f(u) = u^2$ ;  $T = 1$ ;  $\mathbf{V}^n = \mathbf{0}$ ,  $n = 0, \dots, N$ ;  $\mathbf{W}^0 = [1, \dots, 1]$  and  $\mathbf{W}^{n+1}$ ,  $n = 0, \dots, N-1$ , determined by (4.6) with  $\sigma^{n+1} = 0$ . Then

$$\|\Phi_h(V_h) - \Phi_h(W_h)\| = \|\mathbf{V}^0 - \mathbf{W}^0\| = \|\mathbf{W}^0\| = 1,$$

while

$$\|V_h - W_h\| \geq \|\mathbf{V}^N - \mathbf{W}^N\| = \|\mathbf{W}^N\|,$$

a quantity that, according to Sanz-Serna and Verwer[38], grows like

$$1/(h|\log h|).$$

This behaviour of the difference scheme should not be surprising. The vectors  $\{\mathbf{V}^n\}$  and  $\{\mathbf{W}^n\}$  are the numerical solutions of the scheme (2.8)–(2.9) when  $u^0 \equiv 0$  and  $u^0 \equiv 1$ , respectively. For these initial conditions the solutions of (2.3)–(2.5) are, respectively,  $u \equiv 0$  and  $u^* = 1/(1-t)$ . Clearly  $u(x, t) - u^*(x, t)$  cannot be bounded in terms of the difference in initial condition and it cannot be hoped that, in the numerical scheme,  $\|\mathbf{V}^n - \mathbf{W}^n\|$  can be bounded in terms of  $\|\mathbf{V}^0 - \mathbf{W}^0\|$ .

Since, for  $f$  not globally Lipschitz, (2.8)–(2.9) is not N-stable, it follows that the convergence of the scheme cannot be derived by invoking the ‘consistency + N-stability’ result of Section 4.4. In practice the convergence of (2.8)–(2.9) can be proved by using various well-known ad hoc techniques (tricks). Probably the quickest trick is the following (Shampine and Gordon[39], p.24). Introduce a globally Lipschitz function  $f^*$  that coincides with  $f$  in a neighbourhood  $\Omega$  of  $\{u(x, t) : 0 \leq x \leq 1, 0 \leq t \leq T\}$ . The scheme for  $f^*$  is convergent by the material in Section 4.4. Now the theoretical solution  $u^*$  corresponding to  $f^*$  is the solution  $u$  corresponding to  $f$ . Hence the numerical solution  $U_h^*$  corresponding to  $f^*$ , for  $h$  small, takes values in  $\Omega$ , which in turn shows that  $U_h^*$  is the numerical solution  $U_h$  corresponding to  $f$ . Therefore, the convergence of  $U_h^*$  to  $u^*$  actually implies the convergence of  $U_h$  to  $u$ . This technique is often used for terms of the form  $f(u)$ , but does not apply to more general nonlinearities such as the common  $uu_x$ . A second trick, of a wider applicability, is as follows. Since  $f$  is smooth, (4.2) holds when  $v$  and  $w$  are restricted to belong to the neighbourhood  $\Omega$  introduced above. Assume *a priori* that the numerical solution takes values in  $\Omega$ . Then (4.8) holds for the particular choice  $\mathbf{V}^n = \mathbf{u}^n$ ,  $\mathbf{W}^n = \mathbf{U}^n$ . Consequently, recursion from (4.9), taking into account that  $\rho^{n+1} = O(h^2)$ ,  $\sigma^{n+1} = \mathbf{0}$ , leads to a  $O(h^2)$  estimate for

$\|\mathbf{u}^{n+1} - \mathbf{U}^{n+1}\|$ , as long as the vectors  $\mathbf{U}^n, \mathbf{U}^{n-1}, \dots, \mathbf{U}^0$  have their components in  $\Omega$ . The estimate shows that for  $h$  small enough  $\mathbf{U}^{n+1}$  will also have its components in  $\Omega$ . By induction, the estimate holds for  $n$  up to  $N$ .

The lack of applicability of the N-definition to the case at hand may be attributed to its *global* character: all  $V_h$  and  $W_h$  are allowed in (4.1). On the other hand the tricks used to prove convergence attract the attention to the point that in error estimation we are only interested in the local behaviour of  $\Phi_h$  near the theoretical solution. The question arises of whether the definition in (4.1) cannot be made local in some way, so that with the new local definition common convergent schemes like (2.8)–(2.9), are again classified as stable. Local versions of the idea of stability have been introduced by Stetter[42] and H. B. Keller[16]. The advantages and drawbacks of the two versions have been compared by López-Marcos and Sanz-Serna[20]. It turns out that Keller's approach should be favoured. Keller's definition is as follows. The discretization (2.1) is said to be K-stable (K for Keller) if there exist constants  $S > 0$ ,  $h_0 > 0$  and  $R$ ,  $0 < R \leq \infty$ , such that for each  $h$  in  $H$ ,  $h \leq h_0$ , the open ball

$$B(u_h, R) = \{V_h \in X_h : \|V_h - u_h\| < R\}$$

is contained in the domain  $D_h$ , and, for each  $V_h$  and  $W_h$  in  $B(u_h, R)$ , (4.1) holds.

The quantity  $R$  is called the *stability threshold*. It is easy to show that when (2.1) takes the linear form (3.1), a discretization that is stable with threshold  $R < \infty$  is also stable for the threshold  $R = \infty$ . For this choice (4.1), is asked to hold for all  $V_h, W_h$  in  $X_h$  and we recover the standard linear definition. Note that K-stability is a local notion that explicitly refers to the 'theoretical' vectors  $u_h$ . This is different from the naive situation whether the stability or otherwise of a discretization does not relate to the theoretical solution being approximated.

For each real  $R > 0$ , the scheme (2.8)–(2.9), with  $r \leq 1/2$  and  $f$  smooth, is K-stable with threshold  $R$  and stability constant  $S = \exp(LT)$ , where  $L = L(R)$  is the Lipschitz constant of  $f$  in

$$\Omega = \{v : |v - u(x, t)| < R \text{ for some } (x, t), 0 \leq x \leq 1, 0 \leq t \leq T\}.$$

This is shown by the argument in (4.3)–(4.9). Note that now (4.1) has only to be proved for  $V_h$  and  $W_h$  in  $B(u_h, R)$  and that such vectors have their components in  $\Omega$ , so that (4.8) holds for them.

Now that we know that, in the new sense, (2.8)–(2.9) is stable, it is appropriate to ask whether this knowledge is of any help in proving the convergence of that discretization. Actually, it is not evident that, for general discretizations (2.1), consistency and K-stability imply convergence. To begin with, the question of the existence of  $U_h$  must be answered. Then,

it is doubtful that the argument in (3.5) can be applied as it is not clear that  $U_h$  lies in the ball  $B(u_h, R)$  where the stability bound (4.1) holds. At first glance it may seem that the general result would read 'consistency + K-stability + existence of  $U_h$  + a priori estimate  $\|u_h - U_h\| < R'$  imply convergence'. In the context of the scheme (2.8)–(2.9), where the existence of  $U_h$  is obvious, we would need, on top of stability and convergence, the *a priori* bounds  $|u(x_j, t_n) - U_j^n| < R$ . This more or less sends us back to the tricks mentioned above, so that apparently little has been gained from the introduction of the notion of K-stability.

However the considerations just outlined are unduly pessimistic. It is true that consistency and K-stability *on their own* imply convergence. The key ingredient of the proof is the following Lemma due to Stetter[42], whose usefulness in this context was first shown by López-Marcos[18].

**Lemma 5.1.** *Let  $\Phi$  be a  $Y$ -valued mapping defined and continuous in an open ball  $B(v^*, R) = \{v \in X : \|v - v^*\| < R\}$ , where  $X$  and  $Y$  are finite-dimensional normed spaces with  $\dim(X) = \dim(Y)$ . Assume that a positive constant  $S$  exists such that for all  $v$  and  $w$  in  $B_R$*

$$\|v - w\| \leq S\|\Phi(v) - \Phi(w)\|.$$

*Then the inverse mapping  $\Phi^{-1}$  exists (uniquely) in the open ball of radius  $R/S$  centered at  $\Phi(v^*)$ .*

Now assume that in (2.1)  $\Phi_h$  is a continuous mapping (a hypothesis usually satisfied in the applications) and suppose that (2.1) is consistent and K-stable. The application of the lemma to  $\Phi_h$  in the ball centered at  $u_h$  with radius equal to the stability threshold shows that for  $h$  sufficiently small  $\Phi^{-1}(0)$  exists in the open ball  $B(u_h, R)$ , i.e., there is a solution  $U_h$  of the discrete equations that satisfies the bound  $\|u_h - U_h\| < R$ . Then (3.5) can be used to bound the global error in terms of the local discretization error and convergence follows.

The discrete solution  $U_h$  is unique in  $B(u_h, R)$ , so that any other solution  $U_h^*$  of (2.1) is away from  $u_h$  in the sense that  $\|u_h - U_h^*\| \geq R$ . However global uniqueness of  $U_h$  cannot be expected. On the one hand, the original problem being solved is likely to possess several solutions  $u$ . In this situation, often found in nonlinear stationary problems, (2.1) will typically have a solution approximating each possible  $u_h$ . Also nonlinear discretizations should be expected to have spurious discrete solutions, i.e., solutions with no 'theoretical counterpart'. As a simple example take the backward Euler discretization  $u^{n+1} = u^n + h[-u^{n+1} + (u^{n+1})^2]$ , of the well behaved ordinary differential equation Cauchy problem  $u' = -u + u^2$ ,  $0 \leq t \leq 1$ ,  $u(0) = 1/2$ . The equation to be solved *at each step* is quadratic and, for

$u^n$  near the theoretical solution has two real roots one of which is spurious. Hence the discretized equations (2.1) that embrace the computation of  $u^n$  at all time levels  $t_n$ ,  $n = 1, \dots, [1/h]$  certainly possess many solutions. This multiplicity of solutions is very frequent in real-life nonlinear discretizations (for ordinary differential equation problems see Hairer et al. [13], Iserles[15]). Yet it is found to be surprising by some numerical analysts that have been brought up with the concept of naive stability. (Note that N-stability clearly implies global uniqueness of solutions.)

#### 4.6 $h$ -dependent stability thresholds

As shown above, the notion of stability due to Keller, when combined with Stetter's lemma, provides a very convenient method for the analysis of discretizations. Unfortunately some interesting numerical schemes for partial differential equation problems are not K-stable (see e.g., Frutos and Sanz-Serna[10]). In this section we present a useful extension of Keller's definition. It is expedient to study first the following example.

**Example B.** Consider the periodic initial-value hyperbolic problem

$$u_t + uu_x = 0, \quad -\infty < x < \infty, \quad 0 \leq t \leq T < \infty, \quad (6.1)$$

$$u(x+1, t) = u(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T < \infty, \quad (6.2)$$

$$u(x, 0) = u^0(x), \quad -\infty < x < \infty. \quad (6.3)$$

In (6.3)  $u^0$  is a given smooth, 1-periodic function and it is assumed  $T$  is small enough so that the solution of (6.1)–(6.3) is smooth up to  $t = T$ , i.e., the first crossing of characteristics occurs after  $t = T$  (see e.g., Whitham[47]). The equation (6.1) has often been used in the study of nonlinear stability issues, see e.g., Richtmyer and Morton[26], Fornberg[8], Vardillo and Sanz-Serna[46].

The notation for the numerical scheme is similar to that employed in Example A. Choose a positive constant  $r$  (the mesh-ratio) and an integer  $J > 2$ . Set  $h = 1/J$ ;  $x_j = jh$ ,  $j$  integer;  $t_n = nk$ ,  $k = rh$ ,  $n = 0, \dots, N = [T/k]$ . For  $j = 1, \dots, J$  and  $n = 0, \dots, N-1$  set

$$\frac{U_j^{n+1} - U_j^n}{k} + U_j^{n+1/2} \frac{U_{j+1}^{n+1/2} - U_{j-1}^{n+1/2}}{2h} = 0, \quad (6.4)$$

where the index  $j$  must be understood in the obvious periodic way and  $U_j^{n+1/2}$  stands for the average  $(U_j^{n+1} + U_j^n)/2$ . For  $j = 1, \dots, J$  set

$$U_j^0 - u^0(x_j) = 0. \quad (6.5)$$

Denote by  $Z_h$  the space of grid functions  $\mathbf{U} = [U_1, \dots, U_J]$ , endowed with the standard discrete  $L_2$ -norm. The scheme (6.4)–(6.5) may be rewritten

$$\mathbf{U}^0 - \mathbf{u}^0 = \mathbf{0}, \quad (6.6)$$

$$\frac{\mathbf{U}^{n+1} - \mathbf{U}^n}{k} + \mathbf{Q}(\mathbf{U}^{n+1/2}) = \mathbf{0}, \quad n = 0, \dots, N-1, \quad (6.7)$$

where  $\mathbf{Q}$  is a nonlinear operator in  $Z_h$ . Next set again  $u_h$  equal to the grid restriction of  $u$  and  $X_h = D_h = Y_h$  equal to the product of  $N+1$  copies  $Z_h \times \dots \times Z_h$ . The norms in  $X_h$  and  $Y_h$  are derived from the  $L_2$ -norm in  $Z_h$  via (2.10) and (2.11). The scheme is easily seen to be consistent of the second order.

The scheme (6.6)–(6.7) is not Keller stable. In fact consider the case  $u^0 \equiv 0$ ,  $T = 1$ . Introduce the vectors  $V_h = 0 \in X_h$ ,  $W_h = [\mathbf{W}^0, \dots, \mathbf{W}^N] \in X_h$ , where  $\mathbf{W}^n = \mathbf{0}$ ,  $n = 0, 1, \dots, N-1$  and  $\mathbf{W}^N$  has components  $W_j^N = 0$ ,  $j = 3, 4, \dots, J$ ,  $W_1^N = 8/r$ ,  $W_2^N = -8/r$ . It is trivial to check that  $\Phi_h(V_h) = \Phi_h(W_h) = 0 \in Y_h$ . If the scheme were K-stable with threshold  $R$  then, by the local uniqueness of the discrete solution,  $\|W_h - u_h\| \geq R$  ( $h$  sufficiently small). On the other hand,

$$\|W_h - u_h\| = \|W_h\| = \|\mathbf{W}^N\| = 8\sqrt{2h}/r,$$

and we have reached a contradiction.

This example shows that local uniqueness does not take place in an open ball around  $u_h$  of radius larger than  $8\sqrt{2h}/r$ . In this way we are led to the idea of  $h$ -dependent stability thresholds. In his thesis, J. C. López-Marcos[18] introduced the following definition.

**Definition 6.1.** Suppose that, for each  $h$  in  $H$ ,  $R_h$  is a value  $0 < R_h \leq \infty$ . The discretization (2.1) is said to be stable restricted to the thresholds  $R_h$  if there exist positive constants  $h_0$  and  $S$  such that for  $h$  in  $H$ ,  $h \leq h_0$ , the open ball  $B(u_h, R_h)$  is contained in the domain  $D_h$  and for any  $V_h$  and  $W_h$  in  $B(u_h, R_h)$  the bound (4.1) holds.

For linear problems, this notion reduces to the standard (3.3). In general, this definition is weaker than that of Keller, so that it considers as stable schemes that are not K-stable. However the definition is strong enough to prove in some cases that consistency and stability lead to convergence. In fact the following theorem is a direct consequence of Stetter's lemma.

**Theorem 6.2.** Assume that (2.1) is consistent and stable with thresholds  $R_h$ . If  $\Phi_h$  is continuous in  $B(u_h, R_h)$  and  $\|\tau_h\| = o(R_h)$  as  $h \rightarrow 0$ , then:

- (i) For  $h$  sufficiently small the equations (2.1) possess a unique  $U_h$  solution in  $B(u_h, R_h)$ .
- (ii) The global errors in the solutions considered in (i) have a bound  $\|u_h - U_h\| \leq S\|\tau_h\|$ , where  $S$  is the stability constant. In particular the discretization is convergent with an order not smaller than the order of consistency.

**Example B** (revisited). Let us apply this theorem to the discretization (6.6)–(6.7). We first show stability restricted to thresholds  $\mu h^{3/2}$ , where  $\mu$  is any positive constant. If  $V_h = [\mathbf{V}^0, \dots, \mathbf{V}^N]$ ,  $W_h = [\mathbf{W}^0, \dots, \mathbf{W}^N]$ ,  $\Phi_h(V_h) = [\rho^0, \dots, \rho^N]$ ,  $\Phi_h(W_h) = [\sigma^0, \dots, \sigma^N]$ , then

$$\mathbf{V}^0 - \mathbf{u}^0 = \rho^0, \quad (6.8)$$

$$\frac{\mathbf{V}^{n+1} - \mathbf{V}^n}{k} + \mathbf{Q}(\mathbf{V}^{n+1/2}) = \rho^{n+1}, \quad n = 0, \dots, N-1, \quad (6.9)$$

$$\mathbf{W}^0 - \mathbf{u}^0 = \sigma^0, \quad (6.10)$$

$$\frac{\mathbf{W}^{n+1} - \mathbf{W}^n}{k} + \mathbf{Q}(\mathbf{W}^{n+1/2}) = \sigma^{n+1}, \quad n = 0, \dots, N-1. \quad (6.11)$$

We need the following estimate, valid for any  $\mathbf{V}, \mathbf{W}$  in  $Z_h$ :

$$| \langle \mathbf{Q}(\mathbf{V}) - \mathbf{Q}(\mathbf{W}), \mathbf{V} - \mathbf{W} \rangle | \leq M \|\mathbf{V} - \mathbf{W}\|^2, \quad (6.12)$$

where

$$M = M(\mathbf{V}, \mathbf{W}) = \frac{3}{4} \max_j \left( \frac{|V_{j+1} - V_j|}{h} + \frac{|W_{j+1} - W_j|}{h} \right). \quad (6.13)$$

In (6.12), angular brackets denote the standard  $L_2$ -inner product. The proof of (6.12), not given here, follows standard finite-difference energy method manipulations, see e.g., López-Marcos[19]. Subtract (6.11) from (6.9) and take the inner product with  $\mathbf{V}^{n+1/2} - \mathbf{W}^{n+1/2}$ , to obtain

$$\frac{(e^{n+1})^2 - (e^n)^2}{2k} \leq M(\mathbf{V}^{n+1/2}, \mathbf{W}^{n+1/2}) \|\mathbf{V}^{n+1/2} - \mathbf{W}^{n+1/2}\|^2 + \|\rho^{n+1} - \sigma^{n+1}\| \|\mathbf{V}^{n+1/2} - \mathbf{W}^{n+1/2}\|,$$

where we have used the abbreviation  $e_n = \|\mathbf{V}_n - \mathbf{W}_n\|$ . Next

$$\|\mathbf{V}^{n+1/2} - \mathbf{W}^{n+1/2}\| \leq \frac{e^{n+1} + e^n}{2},$$

so that

$$\frac{e^{n+1} - e^n}{k} \leq M(\mathbf{V}^{n+1/2}, \mathbf{W}^{n+1/2}) \frac{e^{n+1} + e^n}{2} + \|\rho^{n+1} - \sigma^{n+1}\|. \quad (6.14)$$

The key observation is now that if  $V_h$  and  $W_h$  satisfy the threshold condition  $V_h, W_h \in B(u_h, R_h)$ , then  $M(\mathbf{V}^{n+1/2}, \mathbf{W}^{n+1/2})$  can be bounded uniformly. This is because bounds  $\|\mathbf{V} - \mathbf{u}^n\| = O(h^{3/2})$ ,  $\|\mathbf{W} - \mathbf{u}^n\| = O(h^{3/2})$  imply that the components of  $\mathbf{V}$ ,  $\mathbf{W}$  are  $O(h)$  away from the components of  $\mathbf{u}^n$  and hence  $M(\mathbf{V}, \mathbf{W})$  is  $O(1)$  away from  $M(\mathbf{u}^n, \mathbf{u}^n) = O(1)$ . This is just an example of the use of inverse inequalities. The stability of the scheme restricted to the thresholds  $R_h = \mu h^{3/2}$  may now be proved by a standard recursion in (6.14). Convergence of the second order (including the existence of solutions of the implicit equations for  $h$  small) follows from the theorem above. We emphasize that the study of the solvability of the equations (6.7) has required no separate proof.

To end the section we show that the restriction to thresholds of the form  $\mu h^{3/2}$  is tight. Again we focus our attention on the case  $u^0 \equiv 0$ ,  $T = 1$ . Fix a number  $\eta$ ,  $0 < \eta < 2$ , and for each  $h$  in  $H$ , consider the recursion

$$\frac{\beta_{n+1} - \beta_n}{k} = \frac{(\beta_{n+1} + \beta_n)^2}{8}, \quad k = rh, \quad n = 0, \dots, N-1, \quad N = [1/k],$$

with initial condition  $\beta_0 = \eta$ . This is the implicit mid-point discretization of the problem  $d\beta/dt = \beta^2/2$ ,  $\beta(0) = \eta$ , with solution  $2\eta/(2 - \eta t)$ , and therefore

$$\beta_N \rightarrow 2\eta/(2 - \eta) \quad \text{as } h \rightarrow 0. \quad (6.15)$$

Now consider vectors  $V_h = 0 \in X_h$ ,  $W_h = [\mathbf{W}^0, \dots, \mathbf{W}^N] \in X_h$ , where  $\mathbf{W}^n = h\beta_n \mathbf{E}$ , with  $\mathbf{E} = [E_1, \dots, E_J]$ ,  $E_1 = 1$ ,  $E_2 = -1$ ,  $E_j = 0$ ,  $j = 3, \dots, J$ . The vector  $\mathbf{E}$  is an 'eigenfunction' of the quadratic operator  $\mathbf{Q}$  in (6.7) (see Vadiello and Sanz-Serna[46]). It is readily checked that  $W_h$  satisfy the equations (6.7). Thus  $\|\Phi_h(W_h)\| = \|\mathbf{W}^0\| = \sqrt{2}h^{3/2}\eta$ . On the other hand  $\|W_h\| = \|\mathbf{W}^N\| = \sqrt{2}h^{3/2}\beta_N$ , so that, on taking into account (6.15),

$$\frac{\|V_h - W_h\|}{\|\Phi_h(V_h) - \Phi_h(W_h)\|} = \frac{\|W_h\|}{\|\Phi_h(W_h)\|} \rightarrow \frac{2}{2 - \eta}, \quad h \rightarrow 0.$$

By implication, if the discretization is stable for some thresholds  $R_h > \|W_h - u_h\|$ , then the corresponding stability constant  $S$  satisfies  $S \geq 2/(2 - \eta)$ . Since  $S$  must be finite,  $\eta$  cannot be allowed to become arbitrarily close to 2. In other words, for  $\eta$  close to 2, and  $h$  small,  $W_h$  must violate the threshold condition, i.e.,

$$R_h \leq \|W_h - u_h\| = \|W_h\| = \sqrt{2}h^{3/2}\beta_N = O(h^{3/2}).$$

#### 4.7 Linear investigation of nonlinear stability

In most practical situations, the mapping  $\Phi_h$  in (2.1) is smooth, so that the (Fréchet) derivative (i.e., roughly the Jacobian matrix)  $\Phi'_h(u_h)$  of  $\Phi_h$  at  $u_h$  exists. Furthermore, if the discretization (2.1) is successful, a solution  $U_h$  of (2.1), close to  $u_h$ , exists. Thus

$$0 = \Phi_h(U_h) \approx \Phi_h(u_h) + \Phi'_h(u_h)(U_h - u_h)$$

and one is led to consider the linearized discretization

$$\Xi_h(U_h) \equiv \Phi_h(u_h) - \Phi'_h(u_h)u_h + \Phi'_h(u_h)U_h = 0. \quad (7.1)$$

Our aim is to study the relation between the stability of (2.1) and that of its linearization (7.1). The main motivation for this sort of research is of course that the stability or otherwise of linear discretizations is more easily investigated than that of their nonlinear counterparts. Our presentation in this section follows closely that in López-Marcos and Sanz-Serna[21]. The following result shows that, under suitable technical assumptions, the stability of (7.1) and (2.1) are equivalent.

**Theorem 7.1.** (i) Assume that for  $h$  in  $H$ ,  $h$  sufficiently small, the mapping  $\Phi_h$  in (2.1) is (Fréchet) differentiable at  $u_h$ . If (2.1) is stable restricted to some thresholds  $R_h$ , then (7.1) is stable.

(ii) Assume that, for each  $h$  in  $H$ ,  $h$  sufficiently small, the mapping  $\Phi_h$  in (2.1) is (Fréchet) differentiable at each point  $v_h$  in an open ball  $B(u_h, R_h)$ . Suppose that (7.1) is stable with stability constant  $L$  and that there exists a constant  $Q$ , with  $0 \leq Q < 1$ , such that, for  $h$  in  $H$ , sufficiently small, and for each  $v_h$  in  $B(u_h, R_h)$ , we have

$$\|\Phi'_h(v_h) - \Phi'_h(u_h)\| \leq Q/L. \quad (7.2)$$

Then (2.1) is stable with thresholds  $R_h$  and stability constant  $L(1 - Q)$ .

Let us look closer at part (ii) above in the common case where  $\Phi_h$  in (2.1) is (Fréchet) differentiable at each point  $v_h$  in an open ball  $B(u_h, R_h)$  and

$$\|\Phi'_h(v_h) - \Phi'_h(u_h)\| \leq K_h \|v_h - u_h\|, \quad \text{for } v_h \in B(u_h, R_h). \quad (7.3)$$

Choose a number  $S$  larger than the stability constant  $L$  of (7.1). If  $h$  is sufficiently small and  $\|v_h - u_h\| \leq \min\{R_h, (L^{-1} - S^{-1})K_h^{-1}\}$ , then

$$\|\Phi'_h(v_h) - \Phi'_h(u_h)\| \leq K_h(L^{-1} - S^{-1})K_h^{-1} = [(S - L)/S]/L,$$

so that (7.2) holds with  $Q = (S - L)/S$ . We conclude that, for discretizations (2.1) satisfying (7.3), linearized stability with constant  $L$  implies stability with constant  $S > L$  and thresholds

$$\min\{R_h, (L^{-1} - S^{-1})K_h^{-1}\}. \quad (7.4)$$

In particular (7.3) holds if the  $\Phi_h$  are continuously differentiable. Thus for smooth discretizations linearized stability is equivalent to stability with some suitable ( $h$ -dependent) thresholds. According to (7.4), the size of the thresholds decreases when  $K_h$  in (7.3) increases, i.e., when the problem becomes more nonlinear.

It is perhaps useful to reconsider the stability of (2.8)–(2.9) with  $f$  smooth. If  $V_h = [\Phi'_h(v_h) - \Phi'_h(u_h)]W_h$ , where  $W_h$  is any vector in  $X_h$ ;  $V_h$  has components

$$\begin{aligned} \mathbf{V}^0 &= \mathbf{0} \\ \mathbf{V}^n &= -(\text{diag}[\mathbf{f}'(\mathbf{v}^n)] - \text{diag}[\mathbf{f}'(\mathbf{u}^n)])\mathbf{W}^n, \quad n = 1, \dots, N. \end{aligned} \quad (7.5)$$

Note that only the nonlinear term in (2.9) has contributed here; the contributions of the linear terms in (2.9) to the Fréchet derivative cancel when subtracting  $\Phi'_h(u_h)$  from  $\Phi'_h(v_h)$ . It is an easy task to see that (7.5) implies (7.3), where  $R_h$  can be chosen to be any ( $h$ -independent) positive number  $R$  and  $K_h$  can also be chosen to be independent of  $h$  ( $K_h$  is essentially the Lipschitz constant of  $f'$  in an  $R$ -neighbourhood of  $\{u(x, t) : 0 \leq x \leq 1, 0 \leq t \leq T\}$ ).

When working similarly in the case of Example B, the components of  $[\Phi'_h(v_h) - \Phi'_h(u_h)]W_h$  have negative powers of  $h$ , a fact that results, via (7.4) in thresholds that decrease with  $h$ .

#### 4.8 Discussion: some open problems

The theory outlined in Sections 4.6 and 4.7 has been proved to be useful in the analysis of convergence of nonlinear discretizations, see the list of references before Example A in Section 4.2. The theory provides a systematic approach to the study of stability and convergence and replaces a number of problem-dependent tricks. Furthermore, with the methodology presented here, the existence of solutions is a consequence of stability and consistency and many *a priori* bounds can be done away with (Frutos and Sanz-Serna[10]).

The issue of whether stability in the senses of Keller or López-Marcos implies stability has not been discussed here. For consistent discretizations, equivalence theorems between convergence and stability can be proved. The reader is referred to the discussions in Sanz-Serna[29] and Palencia and Sanz-Serna[25]).

The theorem in Section 4.6, requires the hypothesis  $\|\tau_h\| = o(R_h)$  as  $h \rightarrow 0$ . In practice, the thresholds are often of the form  $R_h = \mu h^s$ , and this hypothesis means that the order  $p$  of consistency should satisfy  $p > s$ . It is still possible, via an idea of Strang's[43], to prove convergence in some cases where  $s \geq p$ . A systematic treatment of Strang's idea can be seen in Spijker[41]; see also Sanz-Serna[28].

Let us examine some open questions. Consider again the vectors  $W_h = W_h(\eta)$  introduced at the end of Section 4.6. We saw that any threshold condition that renders (6.7) stable ( $u^0 \equiv 0, T = 1$ ) must exclude these vectors when  $\eta$  is large. On recalling that  $\mathbf{W}^0 = h\eta\mathbf{E}$ , we see that any smooth function  $v^0$  whose restriction to the  $x_j$  nodes coincides with  $\mathbf{W}^0$  satisfies  $\max(dv^0/dx) \geq 2\eta$ . Therefore, in the solution of (6.2)–(6.3) with initial condition  $v^0$  the characteristics cross (Whitham[47]) before  $t = 1/(2\eta)$ , so that we should not expect the vectors  $\mathbf{W}^n$  to behave smoothly if  $\eta \geq 1/2$ : they cannot be conceived as approximations of a smooth solution. Even though, as measured in the  $L_2$ -norm (in which convergence is proved),  $\mathbf{W}^0$  is a small  $O(h^{3/2})$  perturbation of  $\mathbf{u}^0 = 0$ ,  $\mathbf{W}^0$  is an  $O(1)$  perturbation in the seminorm  $\max|(dv^0/dx)|$ . It is the latter seminorm which is relevant in deciding the fate of an initial profile. This suggests that the  $O(h^{3/2})$  thresholds in the  $L_2$ -norm could be replaced by more meaningful  $h$ -independent thresholds in a Sobolev norm including the maximum of the first derivative (see also (6.13)). A theory could be considered where a norm in  $X_h$  is used to measure global errors and write the stability bound and a different norm is used to impose the threshold condition. Stability in this alternative setting could then be related to the well-posedness of the continuous problem whose solution  $u$  is being numerically approximated. Such a relation is not possible with the definition in Section 4.6, where it is difficult to find a continuous analogue to the  $h$ -dependency of the thresholds. It would also be interesting to investigate continuous versions of the linearization results in Section 4.7.

#### 4.9 Stability of equilibria via Moser's twist theorem

We now leave the idea of stability in connection with convergence/error bounds and turn our attention to the idea of stability in connection with the long-time qualitative behaviour of discretizations of evolutionary problems.

We consider the familiar systems of ordinary differential equations that describe the motion of a pendulum

$$dp/dt = -\sin q, \quad dq/dt = p. \quad (9.1)$$

The system (9.1) is Hamiltonian with the Hamiltonian function (energy)  $H = p^2/2 + (1 - \cos q)$ . Let  $\phi_t$  denote the flow of (9.1), i.e., the mapping that associates with each point  $(p^0, q^0)$  the value at time  $t$ ,  $(p(t), q(t))$ , of the

solution of (9.1) that satisfies  $p(0) = p^0$ ,  $q(0) = q^0$ . Then, it is well known (see e. g. Arnold[3], Section 16) that  $\phi_t$  is an *area-preserving* mapping, so that for each bounded open set  $\Omega$  in the phase  $(p, q)$ -plane, the sets  $\Omega$  and  $\phi_t(\Omega)$  have the same area. Recall also that (9.1) has periodic solutions, whose trajectories in the phase plane are closed curves which fill the region  $0 < H(p, q) < 2$ . Hence  $H$  is a Liapunov function in the neighbourhood of the origin and the equilibrium  $p = q = 0$  is stable.

The system (9.1) is discretized by the formulae

$$p^{n+1} - p^n = -k \sin q^n, \quad q^{n+1} - q^n = kp^{n+1}. \quad (9.2)$$

Our choice of numerical method has been determined by the fact that the transformation in phase space

$$p^+ = p - k \sin q, \quad q^+ = q + kp - k^2 \sin q \quad (9.3)$$

that maps  $(p^n, q^n)$  into  $(p^{n+1}, q^{n+1})$  is area-preserving (its Jacobian determinant is 1).

We are interested in ascertaining whether the origin is a stable equilibrium of (9.2), i.e., whether numerical solutions remain in the neighbourhood of the origin provided that  $|p^0|$ ,  $|q^0|$  are small enough.

The linearization

$$p^+ = p - kq, \quad q^+ = q + kp - k^2q \quad (9.4)$$

of (9.3) near the origin has eigenvalues

$$(1 - k^2/2) \pm \sqrt{-k^2 + k^4/4}.$$

For  $k > 2$ , the eigenvalues are real and one of them is greater than 1. The origin is therefore unstable, both for (9.4) and for the original (9.3). For  $0 < k < 2$ , we find unimodular complex conjugate eigenvalues  $\lambda$  and  $\lambda^*$ , with

$$\lambda = (1 - k^2/2) + ik\sqrt{1 - k^2/4}. \quad (9.5)$$

In this range of values of  $k$ , the linearization is neutrally stable and no conclusion on the stability of (9.3) can be obtained from the linear analysis. We are going to prove that, for  $0 < k < 2$ , and  $k \neq \sqrt{2}, \sqrt{3}$ , the origin is in fact a stable equilibrium of (9.2)–(9.3). The excluded values  $\sqrt{2}, \sqrt{3}$  are those for which  $\lambda$  in (9.5) takes respectively the values  $i$  and  $(-1 + i\sqrt{3})/2$  (roots of unity). Our first step is to find the *normal form* of (9.3) (Guckenheimer and Holmes[12], Sections 3.3 and 3.5; Arnold[2], Chapter 5, Arnold [3], Appendix 7). Essentially, this means changing variables in (9.3) so as to rewrite the problem in a form better suited for the analysis. The theory of normal forms shows that, because the eigenvalue  $\lambda$  satisfies

$\lambda^3 \neq 1$ ,  $\lambda^4 \neq 1$ , there exists an origin-preserving, invertible, cubic change of variables  $P = P(p, q)$ ,  $Q = Q(p, q)$ , so that in the new variables  $P$ ,  $Q$  the mapping (9.3) is given by

$$(P^+ + iQ^+) = \lambda \exp[i\gamma(P^2 + Q^2)](P + iQ) + O(4), \quad (9.6)$$

where  $O(4)$  denotes terms of order four and higher in the variables  $P$  and  $Q$ , and

$$\gamma = -(k/256)\sqrt{1 - k^2/4} \frac{16 + 4k^2 + 2k^4 - k^4(1 - k^2/2)}{(1 - k^2/4)}. \quad (9.7)$$

In (9.6), we have used a complex form for convenience. A real form may clearly be obtained by separating real and imaginary parts. As  $k$  increases from 0 to 2,  $\gamma$  in (9.7) decreases monotonically from 0 to  $-\infty$ . The variables  $P$ ,  $Q$  have been normalized so that, for  $k = 0$ ,  $P = p$ ,  $Q = q$ .

Let us now discard the  $O(4)$  terms in (9.6) and introduce polar coordinates  $(R, \Theta)$  with  $P = R \cos \Theta$ ,  $Q = R \sin \Theta$ . We obtain

$$R^+ = R, \quad \Theta^+ = \Theta + \arg(\lambda) + \gamma R^2. \quad (9.8)$$

This mapping clearly leaves invariant all circles  $R = \text{constant}$  in the  $(P, Q)$ -plane. On each circle, the mapping acts as a rotation by an angle  $\omega = \arg(\lambda) + \gamma R^2$  that, because  $\gamma \neq 0$ , varies with the radius of the circle. Such mappings are called *twists*. Clearly the origin is a stable equilibrium of (9.8), surrounded by invariant circles. If we change back to the original  $(p, q)$  variables, we obtain invariant closed curves of equation  $P(p, q)^2 + Q(p, q)^2 = \text{constant}$  that surround the origin and hence stability can be concluded.

Note however that the preceding argument applies to (9.8), i.e., to the mapping (9.6) *after removal of the  $O(4)$  terms*. The question arises as to whether the stability of the full (9.6) can be concluded from the stability of the truncated (9.8). At first glance, one may guess that the answer should be negative: since (9.8) is only neutrally stable, arbitrarily small perturbations may render the origin either unstable or asymptotically stable. Nevertheless the answer is positive: (9.6) is an *area-preserving* perturbation of a twist mapping, and, by restricting the attention to a sufficiently small neighbourhood of the origin, the size of the perturbation can be made arbitrarily small. Now according to Moser's twist theorem (see e.g. Siegel and Moser[40], Sections 31-34), if an area-preserving mapping is a sufficiently small perturbation of a twist, then it has an invariant curve surrounding the origin. Thus, in the announced range of values of  $k$ , the origin is a stable equilibrium of the numerical method.

The above technique is very general. Assume that we have a numerical method described by an *area-preserving* mapping in  $\mathbb{R}^2$  for which the origin

is an elliptic equilibrium (i.e., the eigenvalues are complex conjugate  $\lambda$  and  $\lambda^*$  with unit modulus). Then, except in the resonant cases where  $\lambda$  is a cubic or fourth root of unity, the mapping can be brought into the form (9.6) for a suitable real  $\gamma$ . The origin is therefore a stable equilibrium, except, perhaps, in the 'degenerate' case  $\gamma = 0$ . (Note that degeneracy means, that, except for  $O(4)$  terms, the mapping acts as a *linear* rotation, i.e., a rotation where the angular velocity and hence the period are independent of the amplitude.)

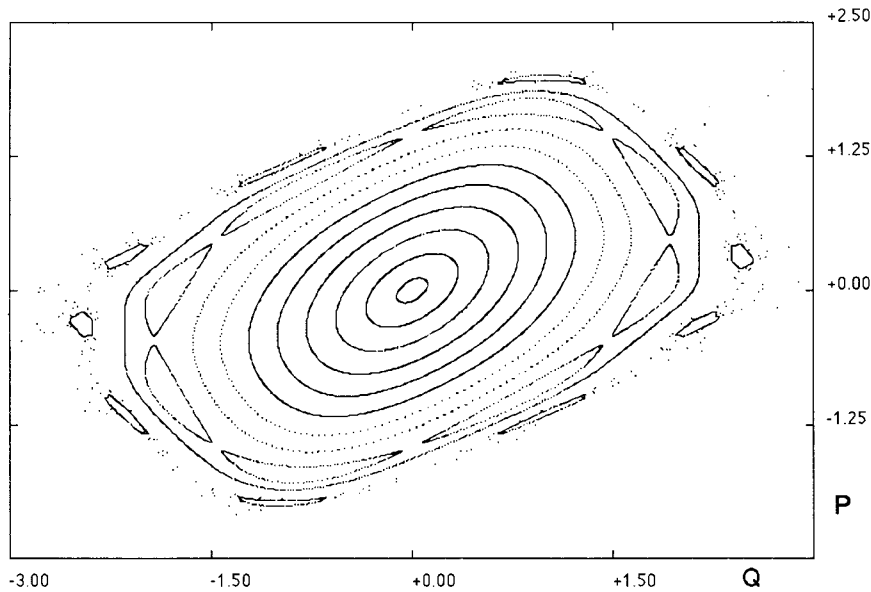
#### 4.10 Behaviour of numerical methods near elliptic equilibria

Let us return to the example (9.2). Now that we know that the origin is a stable equilibrium, we may ask ourselves about the qualitative behaviour of the computed points  $(p^n, q^n)$  near the origin. Judging by the situation in the truncation (9.8), one may guess that there exists a neighbourhood of the origin that is made up of invariant curves of (9.3). If this guess were true, near the origin, the qualitative behaviour of (9.2) would be the same as that of the problem (9.1) being discretized (cf. Beyn[4], Sanz-Serna[32]). However, it turns out that the dynamics of (9.2) is far more complicated than the dynamics of the flow of (9.1).

Let us reconsider the twist  $T$  in (9.8). Assume that, for a given radius  $R$ ,  $\omega = \arg(\lambda) + \gamma R^2$  is rational with respect to  $2\pi$ , i.e., of the form  $2\pi p/q$  with  $p, q$  integers,  $q > 1$ . Then, for each point in the circle of radius  $R$ ,  $q$  applications of the twist send the point back to its initial position after having completed  $p$  revolutions around the origin. Hence  $T^q$  restricted to such a circle is the identity, a structurally unstable mapping whose dynamics is highly sensitive to perturbations. Accordingly, for such values of  $R$ , it can be shown that the invariant circle of the twist  $T$  disappears under the addition of the  $O(4)$  perturbation leading to the full mapping (9.6). On the other hand, the circles of the twist for which  $\omega$  is very irrational with respect to  $2\pi$  are not destroyed under the perturbation, and give rise to invariant curves of (9.6). Here very  $\omega/2\pi$  irrational means (Arnold[2], Sections 12-13) that for some suitable positive numbers  $K$  and  $s$ ,

$$\left| \frac{\omega}{2\pi} - \frac{p}{q} \right| \geq \frac{K}{q^{2+s}}, \quad (10.1)$$

for all integers  $p$  and  $q$ ,  $q > 0$ . (Recall that for any irrational number  $\nu$  there exist infinitely many rational approximations whose error is less than the reciprocal value of the square of the denominator  $|\nu - p/q| < 1/q^2$ .) The set of invariant curves of (9.6) obtained in this way form a subset of the  $(P, Q)$ -plane with positive measure. In actual fact the complement of this subset in the circle  $P^2 + Q^2 < \rho$  has an area  $o(\pi\rho^2)$  so that the majority of points near the equilibrium belong to the set of invariant curves. In



**Fig. 4.1.** Solutions of 9.2 with varying initial conditions

the gaps between the invariant curves the dynamics of (9.6) is exceedingly complicated (see e.g., Guckenheimer and Holmes[12], Section 4.8).

An illustration is provided in Figure 4.1, where the solutions of (9.2) corresponding to 12 different initial conditions have been presented. We have used  $k = 1$  and for each solution 2000 points have been displayed. For the initial conditions  $(p^0, q^0)$  given by  $(0, 0.1)$ ,  $(0, 0.3)$ ,  $(0, 0.5), \dots, (0, 1.5)$  we observe that the solutions correspond to invariant curves. The initial condition  $(0, 1.7)$  gives rise to an orbit consisting of eight suborbits (islands), so that after eight time-steps the point returns to the original suborbit. Thus, if only every eighth iterate were plotted or in other words the 8-th power  $M^8$  of the mapping  $M$  in (9.3) were considered, then only one suborbit would be seen. In fact this suborbit is a twist theorem invariant curve of  $M^8$  around an elliptic equilibria (i.e. around an elliptic 8-periodic equilibrium of (9.3)). For the initial condition  $(0, 2.1)$  we see again an invariant curve and for  $(.26, 2.55)$  a structure of 10 islands is found. Finally for  $(.33, 2.55)$  the computed points behave erratically and eventually leave the plotting window.

Note that the dynamics depicted in Figure 4.1 cannot be the dynamics of the  $k$ -flow of a differential system: in the latter, all trajectories would

lie on curves. This shows that, in general, numerical trajectories cannot be viewed as exact trajectories of a system of ordinary differential equations close to that being integrated.

#### 4.11 An alternative application of Moser's twist theorem

It is useful to emphasize that the analysis above has only employed the properties of the numerical method (9.2), without taking into account any feature of the system (9.1) being approximated. Furthermore the analysis is valid for any fixed  $k$  in the announced range and  $(p, q)$  *sufficiently small*. There is an alternative way in which Moser's twist theorem could have been used (Sanz-Serna and Vadiello[36,37]). The starting point for the alternative approach is that in the region  $H(p, q) < 2$ , it is possible to change the dependent variables  $(p, q)$  of the continuous problem (9.1) into the so-called *action/angle* variables  $(I, \alpha)$  (Arnold[3], Chapter 10). (The abstract angle  $\alpha$  should not be confused with the physical angle  $q$  by which the pendulum deviates from the vertical axis.) Among the properties of  $(I, \phi)$  we need the following three: (i)  $p = p(I, \alpha)$ ,  $q = q(I, \alpha)$  are  $2\pi$ -periodic in  $\alpha$ , i.e.,  $\alpha$  behaves as a genuine angle. (ii)  $I$  takes the value 0 at the origin and increases away from it. (iii) In the new variables (9.1) takes the simple form

$$dI/dt = 0, \quad d\alpha/dt = J(I), \quad (11.1)$$

where  $J$  is a known function of the action  $I$ . It is possible to give closed form expressions for  $I(p, q)$ ,  $\alpha(p, q)$ ,  $J(I)$  in terms of elliptic integrals, but such expressions are not actually needed. The main advantage of the new variables is that now (11.1) is readily integrated to yield

$$I(t) = I(0), \quad \alpha(t) = J[I(0)]t + \alpha(0). \quad (11.2)$$

As a consequence, the  $k$ -flow of the system is now simply given by  $(I, \alpha) \rightarrow (I, \alpha + J[I(0)]k)$ , i.e., by a twist. By consistency, (9.3), when written in action/angle variables, is an  $O(k^2)$  perturbation of this  $O(k)$  twist. For  $0 < H(p, q) < 2$ ,  $k$  *small enough*, Moser's theorem can be invoked to prove the existence of invariant curves.

We make a final remark. It is well known that the function  $J(I)$  in (11.1) is decreasing so that the period  $2\pi/J(I)$  of the oscillations increases with  $I$  (i.e., the larger the amplitude  $q_{max}$  of the swing of the pendulum the larger the period). This matches the fact that  $\gamma$  in (9.7) is negative. Moreover the classical theory of the motion of the pendulum shows that the period is given by  $2\pi(1 + (1/16)q_{max}^2) + o(q_{max}^2)$ . This should be compared with the value  $2\pi(1 + (1/16)q_{max}^2) + o(k + q_{max}^2)$  in the numerical (9.7)–(9.8). This agreement is of course no coincidence and could have been anticipated by the convergence of (9.2).

#### 4.12 Symplectic numerical integrators

The material in the previous sections has used the twist theorem, a result restricted to dynamical systems with two dependent variables. In numerical analysis we are of course interested in multidimensional systems, including those resulting from the space discretization of evolutionary partial differential equations. The KAM theory (see e.g. Arnold[3], Appendix 8, Guckenheimer and Holmes[12], Section 4.8) provides the multidimensional extension of the twist theorem. The role played above by area-preserving mappings is now played by *symplectic* mappings. A mapping  $T: \mathbf{p}^+ = \mathbf{p}^+(\mathbf{p}, \mathbf{q})$ ,  $\mathbf{q}^+ = \mathbf{q}^+(\mathbf{p}, \mathbf{q})$ ,  $\mathbf{p} \in \mathbb{R}^g$ ,  $\mathbf{q} \in \mathbb{R}^g$  is called symplectic or canonical (Arnold[3] Chapters 8-9) if it preserves the differential form

$$d\mathbf{p} \wedge d\mathbf{q} = dp^{(1)} \wedge dq^{(1)} + \dots + dp^{(g)} \wedge dq^{(g)}, \quad (12.1)$$

(superscripts denote components). In plain terms, this means that if we choose an open bounded set  $\Omega$  in phase space  $\mathbb{R}^g \times \mathbb{R}^g$ , project it onto the  $g$  two-dimensional planes of the variables  $(p^{(i)}, q^{(i)})$ ,  $1 \leq i \leq g$  and sum the two-dimensional areas of the resulting projections, we obtain the same result for  $\Omega$  and for the image  $T(\Omega)$ .

The flow  $\phi$  of a differential system is canonical if and only if the system is Hamiltonian, i.e. there exists a real function  $H(\mathbf{p}, \mathbf{q})$  in phase space so that

$$dp^{(i)}/dt = -\partial H/\partial q^{(i)}, \quad dq^{(i)}/dt = \partial H/\partial p^{(i)}, \quad 1 \leq i \leq g.$$

Consequently the numerical integration of Hamiltonian systems is the natural setting in which symplectic mappings can occur in numerical analysis. A one-step method for the numerical integration of Hamiltonian system is said to be symplectic or canonical if when applied to any Hamiltonian system with any step-length it gives rise to a symplectic transformation in phase-space (Ruth[27], Feng[7], Sanz-Serna[31,33], Lasagni[17], Suris[45]).

For symplectic schemes, the KAM theory can be applied to obtain results like those in Sections 4.9–4.11. Furthermore, since the symplectic character of the flow characterizes Hamiltonian systems, the qualitative features of Hamiltonian dynamics derive from the conservation of (12.1). Hence the dynamics of symplectic schemes should be expected to mimic the qualitative features of the Hamiltonian flow. This point has been discussed in Sanz-Serna[33].

Let us present some examples of symplectic methods.

**Example (A).** An  $s$ -stage Runge-Kutta method of the form

$$\begin{aligned} \mathbf{Y}_i &= \mathbf{y}^n + k \sum_{1 \leq j \leq s} a_{ij} \mathbf{F}(\mathbf{Y}_j), \quad 1 \leq i \leq s, \\ \mathbf{y}^{n+1} &= \mathbf{y}^n + k \sum_{1 \leq i \leq s} b_i \mathbf{F}(\mathbf{Y}_i), \end{aligned}$$

is canonical provided that for each  $i, j = 1, \dots, s$ ,

$$b_i a_{ij} + b_j a_{ji} - b_i b_j = 0, \quad (12.2)$$

(see Sanz-Serna[31], Lasagni[17], Suris[45]). The Gauss-Legendre methods with  $s$  stages and order  $2s$  are canonical. This includes the standard mid-point rule. It is interesting to note that, as shown by Sanz-Serna and Abia[34]), when condition (12.2) holds, the conditions for the RK method to be of order  $p$  are notably simplified. Rather than a condition per rooted tree with  $p$  or fewer nodes (Butcher[6]), we have a condition per so-called non-superfluous tree. For instance, for order  $p = 5$ , 17 order conditions are required in general and only 6 when (12.2) holds.

**Example (B).** The explicit mid-point rule is often used to advance in time the solution of systems of ordinary differential equations arising from the space-discretization of time-dependent partial differential equations (leap-frog schemes). Since this rule is not a one-step method, it does not give rise directly to a mapping in phase space and the definition of canonicity is not applicable. Sanz-Serna[30] noticed that it is possible to rewrite the rule as one-step convergent discretization. The technique is as follows. Write two consecutive steps of the explicit mid-point rule

$$\mathbf{y}^{2n+2} = \mathbf{y}^{2n} + 2k\mathbf{F}(\mathbf{y}^{2n+1}), \quad \mathbf{y}^{2n+3} = \mathbf{y}^{2n+1} + 2k\mathbf{F}(\mathbf{y}^{2n+2}), \quad (12.3)$$

as applied to a general system of ordinary differential equations

$$d\mathbf{y}/dt = \mathbf{F}(\mathbf{y}), \quad (\mathbf{y} \in \mathbb{R}^d). \quad (12.4)$$

On denoting  $\mathbf{y}^{2n} = \mathbf{u}^n$ ,  $\mathbf{y}^{2n+1} = \mathbf{v}^n$ ,  $n$  integer, (12.3) becomes

$$\mathbf{u}^{n+1} = \mathbf{u}^n + 2k\mathbf{F}(\mathbf{v}^n), \quad \mathbf{v}^{n+1} = \mathbf{v}^n + 2k\mathbf{F}(\mathbf{u}^{n+1}). \quad (12.5)$$

Now (12.5) is a convergent one-step method for the integration of the  $2d$ -dimensional system

$$d\mathbf{u}/dt = \mathbf{F}(\mathbf{v}), \quad d\mathbf{v}/dt = \mathbf{F}(\mathbf{u}). \quad (12.6)$$

This is called the augmented system corresponding to the original system (12.4). While (12.3) and (12.5) differ only in notation, the interpretation associated with (12.5) should be preferred in the study of the dynamics of the solution. In fact (12.3) is a linear two-step method violating the strong root condition and hence gives rise to a dynamics widely different from that of the original (12.4). On the other hand, (12.5), being of a one-step nature, inherits many features of the dynamics of the system (12.6) it approximates (see Beyn[4,5], Sanz-Serna[32] and the references therein). Sanz-Serna and Vadiello[37] show that, if the original system is Hamiltonian, then the augmented system is also Hamiltonian and that, in such a case, (12.5) is a symplectic scheme for the approximation of the augmented system.

For other symplectic methods see Sanz-Serna[33] and references therein. There is much work to be done in constructing, implementing and testing symplectic formulae.

## Acknowledgement

The author has been partly supported by project CICYT PB-86-0313. He is most thankful to Dr. L. Abia for his help in the preparation of Section 4.9.

## References

- [1] L. Abia, and J. M. Sanz-Serna, *The spectral accuracy of a fully-discrete scheme for a nonlinear third order equation*. Computing, 44, 1990, 187-196.
- [2] V. I. Arnold, *Geometrical methods in the theory of ordinary differential equations*. Springer-Verlag, New York, 1988.
- [3] V. I. Arnold, *Mathematical methods of classical mechanics* (2nd ed.). Springer-Verlag, Berlin, 1989.
- [4] W.-J. Beyn, *On the numerical approximation of phase portraits near stationary points*. SIAM J. Numer. Anal., 24, 1987, 1095-1113.
- [5] W.-J. Beyn, *On invariant closed curves for one-step methods*. Numer. Math., 51, 1987, 103-122.
- [6] J. C. Butcher, *The numerical analysis of ordinary differential equations*. John Wiley, Chichester, 1987.
- [7] K. Feng, *Difference schemes for Hamiltonian formalism and symplectic geometry*. J. Comput. Mat., 4, 1986, 279-289.
- [8] B. Fornberg, *On the instability of leap-frog and Crank-Nicolson approximations of a nonlinear partial differential equation*. Math. Comput., 27, 1973, 45-57.
- [9] J. de Frutos and J. M. Sanz-Serna, *Split-step spectral schemes for nonlinear Dirac systems*. J. Comput. Phys., 83, 1989, 407-423.
- [10] J. de Frutos and J. M. Sanz-Serna *h-dependent stability thresholds avoid the need for a priori bounds in nonlinear convergence proofs*. In *Computational Mathematics III, Proceedings of the Third International Conference on Numerical Mathematics and its Applications*, January 1988, Benin City, Nigeria, (in the press).
- [11] J. de Frutos, T. Ortega and J. M. Sanz-Serna, *A Hamiltonian, explicit algorithm with spectral accuracy for the 'good' Boussinesq system*. Comput. Methods Appl. Mech. Engrg., 80, 1990, 417-423.

- [12] J. Guckenheimer and P. Holmes, *Nonlinear oscillations, dynamical systems, and bifurcations of vector fields*. Springer-Verlag, New York, 1983.
- [13] E. Hairer, A. Iserles and J. M. Sanz-Serna, *Equilibria of Runge-Kutta methods*. University of Cambridge, Numerical Analysis Report DAMTP/1989NA 4, 1989.
- [14] E. Hairer, S. P. Nørsett and G. Wanner, *Solving ordinary differential equations I. Nonstiff problems*. Springer-Verlag, Berlin, 1987.
- [15] A. Iserles, *Stability and dynamics of numerical methods for nonlinear ordinary differential equations*. IMA J. Numer. Anal., 10, 1990, 1-30.
- [16] H. B. Keller, *Approximation methods for nonlinear problems with application to two-point boundary value problems*. Math. Comput., 29, 1975, 464-474.
- [17] F. Lasagni, *Canonical Runge-Kutta methods*. ZAMP, 39, 1988, 952-953.
- [18] J. C. López-Marcos, *Estabilidad de discretizaciones no lineales*. Ph. D. Thesis, Universidad de Valladolid, Valladolid, 1985.
- [19] J. C. López-Marcos, *A difference scheme for a nonlinear partial integrodifferential equation*. SIAM J. Numer. Anal., 27, 1990, 20-31.
- [20] J. C. López-Marcos and J. M. Sanz-Serna, *A definition of stability for nonlinear problems*. In *Numerical treatment of differential equations* (Strehmel, K. ed.), pp. 216-226. Teubner, Leipzig, 1988.
- [21] J. C. López-Marcos and J. M. Sanz-Serna, *Stability and convergence in Numerical Analysis III: Linear investigation of nonlinear stability*. IMA J. Numer. Anal., 8, 1988, 71-84.
- [22] T. Murdoch and C. Budd, *Convergent and spurious solutions of nonlinear elliptic equations*. Preprint, 1989.
- [23] T. Ortega and J. M. Sanz-Serna, *Nonlinear stability and convergence of finite-difference methods for the 'good' Boussinesq equation*. Numer. Math., (in the press).
- [24] C. Palencia and J. M. Sanz-Serna, *An extension of the Lax-Richtmyer theory*. Numer. Math., 44, 1984, 279-283.
- [25] C. Palencia and J. M. Sanz-Serna, *Equivalence theorems for incomplete spaces: an appraisal*. IMA J. Numer. Anal., 4, 1984, 109-115.
- [26] R. D. Richtmyer and K. W. Morton, *Difference methods for initial-value problems*. John Wiley, New York, 1967.
- [27] R. Ruth, *A canonical integration technique*. IEEE Trans. Nucl. Sci., 30, 1983, 2669-2671.
- [28] J. M. Sanz-Serna, *Convergence of the Lambert-McLeod trajectory solver and of the CELF method*. Numer. Math., 45, 1984, 173-182.

- [29] J. M. Sanz-Serna, *Stability and convergence in numerical analysis I: Linear problems, a simple comprehensive account*. In *Nonlinear differential equations* (Hale, J. K. and Martínez Amores, P. eds.), pp. 64-113. Pitman, Boston, 1985.
- [30] J. M. Sanz-Serna, *Studies in numerical nonlinear instability I. Why do leapfrog schemes go unstable?* SIAM J. Sci. Stat. Comput., 6, 1985, 923-938.
- [31] J. M. Sanz-Serna, *Runge-Kutta schemes for Hamiltonian systems*. BIT, 28, 1988, 877-883.
- [32] J. M. Sanz-Serna, *Numerical ordinary differential equations vs. dynamical systems*. Universidad de Valladolid, Applied Mathematics and Computation Report 1990/3, 1990.
- [33] J. M. Sanz-Serna, *The numerical integration of Hamiltonian systems*. To appear.
- [34] J. M. Sanz-Serna and L. Abia, *Order conditions for canonical Runge-Kutta schemes*. Universidad de Valladolid, Applied Mathematics and Computation Report 1990/1, 1990.
- [35] J. M. Sanz-Serna and C. Palencia, *A general equivalence theorem in the theory of discretization methods*. Math. Comput., 45, 1985, 143-152.
- [36] J. M. Sanz-Serna and F. Vadillo, *Nonlinear instability, the dynamic approach*. In *Numerical Analysis* (Griffiths, D. F. and Watson, G. A. eds.), pp. 187-199. Longman, London, 1986.
- [37] J. M. Sanz-Serna and F. Vadillo, *Studies in nonlinear instability III: augmented Hamiltonian systems*. SIAM J. Appl. Math., 47, 1987, 92-108.
- [38] J. M. Sanz-Serna and J. G. Verwer, *A study of the recursion  $y_{n+1} = y_n + \tau y_n^m$* . J. Math. Anal. Appl., 116, 1986, 456-464.
- [39] L. F. Shampine and M. K. Gordon, *Computer solution of ordinary differential equations*. W. H. Freeman and Co., San Francisco, 1975.
- [40] C. L. Siegel and J. K. Moser, *Lectures on Celestial Mechanics*. Springer, Berlin-Heidelberg-New York, 1971.
- [41] M. N. Spijker, *Equivalence theorems for nonlinear finite difference methods*. In *Numerische Behandlung nichtlinearer Integrodifferential und Differential Gleichungen* (Ansorge, R. and Törnig, W., eds.), pp. 109-122. Springer, Berlin, 1974.
- [42] H. J. Stetter, *Analysis of discretization methods for ordinary differential equations*. Springer Verlag, Berlin-Heidelberg-New York, 1973.
- [43] G. Strang, *Accurate partial difference methods II, nonlinear problems*. Numer. Math., 6, 1964, 37-46.

- [44] E. E. Süli, *Convergence and nonlinear stability of the Lagrange Galerkin method for the Navier-Stokes equations*. Numer. Math., 53, 1988, 459-483.
- [45] Y. B. Suris, *Canonical transformations generated by methods of Runge-Kutta type for the numerical integration of the system  $x'' = -\partial U/\partial x$* . Zh. Vychisl. Mat. i Mat. Fiz., 29, 1989, 202-211 (in Russian).
- [46] F. Vadillo and J. M. Sanz-Serna, *Studies in nonlinear instability II. A new look at  $u_t + uu_x = 0$* . J. Comput. Phys., 66, 1986, 225-238.
- [47] G. B. Whitham, *Linear and nonlinear waves*. Wiley-Interscience, New York, 1974.

Professor J. M. Sanz-Serna  
Dpto. Matematica Aplicada y Computacion  
Facultad de Ciencias  
Universidad de Valladolid  
Valladolid  
Spain