

On the Solution of Nonlinear Hyperbolic Differential Equations by Finite Differences

By RICHARD COURANT, EUGENE ISAACSON, and MINA REES

Introduction

Existence, uniqueness and stability of the solution of the initial value problem for the hyperbolic system of n quasi-linear first order partial differential equations in two independent variables have been established in previous publications by various authors [1, 2, 4, 5, 6, 7]. The present paper supplements these results by proving that a rather flexible finite difference scheme provides an approximate numerical solution accurate within an error of the same order of magnitude as the mesh width of the net. Two such schemes, adaptable to computation by intelligent human effort and by automatic machines, are considered in detail. The first,* briefly described in [3], uses a curvilinear net while the second is based on a rectangular net. We determine criteria for selecting the size of mesh widths, the type of difference quotients, and the proper number of decimal places, to insure convergence of the numerical procedure. These conclusions are obtained by a straightforward analysis of the step by step growth of the "error," the difference between the true solution and the finite difference solution. It should be emphasized that our task is simplified by making use of the existence theorem in the form proved in [5] by R. Courant and P. Lax.

1. Differences along a Curvilinear Net

The general quasi-linear system of first order partial differential equations in two independent variables has the form

$$(1) \quad \sum_{i=1}^n A^{ii} u_x^i + B^{ii} u_y^i = C^i, \quad i = 1, \dots, n,$$

where A^{ii} , B^{ii} and C^i are functions of the $n + 2$ variables (x, y, u^i) . If the system is hyperbolic, n distinct linear combinations of the equations (1) may be formed to produce an equivalent system of equations in the normal form¹

*In present day computing practice, either characteristic or rectangular nets are used. Our discussion of curvilinear nets is supposedly a model for convergence proofs that may be fashioned for the characteristic nets (see footnote 4). It was thought worthwhile to include an analysis of a rectangular net scheme to clarify the role that the characteristic directions play in determining such a difference scheme.

¹Several illustrations of such reduction to normal form are to be found in [3].

$$(2) \quad \sum_{i=1}^n a^{ij}(u_v^i + c^j u_x^i) = b^j, \quad j = 1, \dots, n,$$

where a^{ij} , c^j and b^j are functions of (x, y, u^i) (assuming that the direction $y = \text{constant}$ is not characteristic).

We observe that in each equation of the normal form (2) the variables u^i are differentiated in a common direction, that is

$$\frac{du^i}{dy} = u_v^i + c^j u_x^i,$$

along a curve $x = x(y)$ for which

$$\frac{dx}{dy} = c^j \quad (\text{the } j^{\text{th}} \text{ characteristic direction}).$$

It is the direct replacement of the derivative along characteristics that forms the basis of the first finite difference method.

We consider the initial value problem for the equations (2) under the hypothesis that there exists a vector function

$$g = \{g^i(x, y)\}$$

such that along the continuously differentiable initial curve, I_0 , given by $y = y_0(x)$, $a \leq x \leq b$,

$$u_0 = \{u^i(x, y_0(x))\} = \{g^i(x, y_0(x))\} = g_0.$$

Furthermore, in some neighborhood R_1 of the initial curve we suppose that the first derivatives of g satisfy a uniform Lipschitz condition

$$(\text{e.g., } |g_x^i(x_2, y_2) - g_x^i(x_1, y_1)| \leq M(|x_2 - x_1| + |y_2 - y_1|)).$$

In addition we assume that, for (x, y) in R_1 ,

$$|u - g| = \sum_{i=1}^n |u^i - g^i| < K$$

(with some positive constant K), that $|c^j| < M$, that $\det |a^{ij}| \neq 0$, that the first derivatives of a^{ij} , c^j and b^j satisfy a Lipschitz condition with respect to the $n + 2$ variables x, y, u^i and that we have a one-parameter family² of sectionally smooth curves $I_r : y = y_r(x)$ which cover R_1 smoothly and which are not characteristic at any point (i.e. the curves I_r are "spacelike"). In [5] it is proved that in a, perhaps smaller, neighborhood R there exists a unique solution u which satisfies the initial conditions and such that its first derivatives actually satisfy a uniform Lipschitz condition. We are now led to the following procedure for computing the solution of the initial value problem for (2).

We select a mesh width h and pick net points on I_0 which are very nearly h units apart (see Figure 1). Then we pick net points on I_r , $r = 1, 2, \dots$,

²For example, such a family may be obtained by the translation of I_0 in a fixed direction.

which on each curve are very nearly h units apart but which between neighboring curves I_r, I_{r+1} are merely of the order of magnitude of h units apart. The order of the distance between I_r and I_{r+1} is chosen sufficiently small so that if P is a net point on I_{r+1} and Q its nearest neighboring net point on I_r , then the line segments drawn backwards through P with the slopes $dx/dy = c^i$ intersect I_r at points which lie between Q and its nearest neighboring net points Q' and Q'' on I_r .

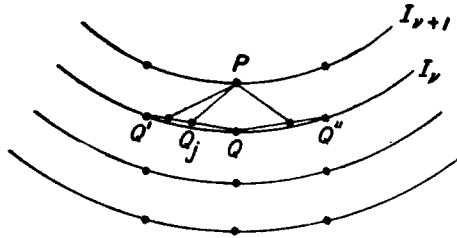


FIGURE 1

The natural step now is to replace in (2) derivatives in the characteristic directions by differences. To make these ideas clear we let $w^i(Q)$ represent the solution of the difference equation (3), which we assume has been computed at the net points on I_0, I_1, \dots, I_r . We then compute $w^i(P)$ by means of the equation³

$$(3) \quad a^{ii}(Q, w^i(Q)) \left[\frac{w^i(P) - w^i(Q_i)}{\Delta y} \right] = b^i(Q, w^i(Q)),$$

where Δy is the change in y between Q_i and P and Q_i is the point of intersection of the straight line through P with slope $dx/dy = c^i(Q, w^i(Q))$ and either the straight line segment joining Q to Q' or the one joining Q to Q'' . Let us assume that Q_i lies between Q and Q' and that in vector form

$$(4) \quad Q_i = Q + \lambda_i(Q' - Q).$$

We define the functions $w^i(Q_i)$ by linear interpolation with the same constant λ_i , in the form

$$f(Q_i) = f(Q) + \lambda_i[f(Q') - f(Q)].$$

We note that λ_i is a function of the solution $w^i(Q)$, the characteristic direction $c^i(Q, w^i(Q))$ and the geometry of the curves I_r .

The computation of $w(P)$ can be effected if $\det |a^{ii}(Q, w(Q))| \neq 0$, which is assured when Q is in the region R and $|w - g| < K$. The purpose of the present paper is to show that in a region D , independent of h and perhaps smaller than R , the functions $w(Q)$ can be defined and will differ from the solution $u(Q)$ of the

³In (3) and in subsequent formulas we shall omit the summation sign for the repeated superscript i .

original differential equation (2), by a quantity of order h . The generality of this method has the advantage⁴ that a skilled computer has freedom in choosing the net.

2. Proof of Convergence

We shall use the symbol $O(h)$ to designate any quantity of order h , i.e., $|O(h)/h| < N$ for $h < \delta$. For example, in (3), Δy designates the change in y from Q to P and may be replaced by $O(h)$ in the region R . For the solution u of (2), we observe that

$$(5) \quad a^{ij}(Q, u(Q))[u^i(P) - u^i(Q^*)] = (\Delta y^*)b^i(Q, u(Q)) + O(h^2)$$

where $Q^* = Q + \lambda_i(u(Q))(Q' - Q)$ and we have made use of the fact that the first derivative of u satisfies a Lipschitz condition. Now since $\lambda_i(Q, u(Q))$ satisfies a Lipschitz condition with respect to its last variables, $u(Q)$,

$$\begin{aligned} u^i(Q_i) - u^i(Q^*) &= O(|Q_i - Q^*|) \\ &= O(|\lambda_i(w(Q)) - \lambda_i(u(Q))| |Q' - Q|) \\ (6) \quad &= O(|w(Q) - u(Q)| |Q' - Q|) \\ &= O(hv) \end{aligned}$$

holds, where

$$(7) \quad v = w - u^*.$$

We define the function u^* by

$$\begin{aligned} u^*(Q) &= u(Q) \quad (\text{for } Q \text{ a net point}) \\ (8) \quad u^*(Q_k) &= (1 - \lambda)u(Q) + \lambda u(Q') \end{aligned}$$

when

$$Q_k = (1 - \lambda)Q + \lambda Q' \quad \text{for} \quad 0 \leq \lambda \leq 1.$$

Since u is differentiable, we note that $u^*(Q_i) - u(Q_i) = O(h^2)$.

We shall now proceed to estimate the error function v . If in (5) we replace $u(Q^*)$ by its value as given in (6) we obtain, since $\Delta y^* - \Delta y = O(hv)$,

$$(9) \quad a^{ij}(Q, u(Q))[u^i(P) - u^{i*}(Q_i)] = (\Delta y)b^i(Q, u(Q)) + O(hv) + O(h^2).$$

⁴It is probably simpler to devise more rapidly convergent approximations to the differential equations when the curvilinear net is used. For a description of such a procedure in a special gas dynamical application, see [3], page 202. Recently, some computations have been made on the Eniac at Aberdeen, for the purpose of comparing the accuracy of various difference schemes (Clippinger, Dimsdale, *et al.*), based on a characteristic net. A proof of convergence for such schemes along the above lines can probably be given.

We rearrange the left side of equation (3) to find that

$$\begin{aligned} a^{ij}(Q, w(Q))[w^i(P) - w^i(Q_i)] \\ &= a^{ij}(Q, u(Q))[w^i(P) - w^i(Q_i)] + O(v^i(Q))[u^i(P) - u^{i*}(Q_i) + v^i(P) - v^i(Q_i)] \\ &= a^{ij}(Q, u(Q))[w^i(P) - w^i(Q_i)] + O(hv(Q)) + O(v(Q)v(P)) + O(v(Q)v(Q_i)), \end{aligned}$$

provided that $|w(Q) - g(Q)| < K$ (see end of Section 4).

Finally

$$\begin{aligned} (10) \quad a^{ij}(Q, w(Q))[w^i(P) - w^i(Q_i)] \\ &= a^{ij}(Q, u(Q))[w^i(P) - w^i(Q_i)] + O(hv(Q)) + O(v(Q)v(P)) \\ &\quad + O(v^2(Q)) + O(v^2(Q_i)). \end{aligned}$$

By subtracting (9) from (3) and using (10), we find

$$\begin{aligned} (11) \quad a^{ij}(Q, u(Q))[v^i(P) - v^i(Q_i)] \\ &= O(v(Q)v(P)) + O(v^2(Q)) + O(v^2(Q_i)) + O(hv(Q)) + O(h^2). \end{aligned}$$

We introduce another measure of the error, $V = \{V^i\}$, by

$$(12) \quad a^{ij}(Q, u^i(Q))v^i(Q) = V^i(Q),$$

for net points Q , and as in (8) we set

$$V^i(Q_k) = (1 - \lambda)V^i(Q) + \lambda V^i(Q').$$

Since a^{ij} and u^* satisfy a Lipschitz condition, equation (11) may be modified⁵ to yield

$$\begin{aligned} (13) \quad V^i(P) - V^i(Q_i) &= O(hV(P)) + O(hV(Q_i)) + O(V^2(Q)) \\ &\quad + O(V^2(Q_i)) + O(V(Q)V(P)). \end{aligned}$$

Owing to the fact that $V^i(Q_k)$ is an average with positive weights of the values of V^i at the neighboring net points, we may introduce as norm on I ,

$$(14) \quad \max_{(j, Q)} |V^j(Q)| = M,$$

and use (13) to get

$$(15) \quad M_{r+1} \leq M_r + \alpha(hM_r + hM_{r+1} + M_r M_{r+1} + M_r^2 + h^2),$$

with α a positive constant.

If we show that for a fixed domain D contained in R , $M_r = O(h)$, the con-

⁵Since $a^{ij}(u(Q))v^i(Q_i) = a^{ij}(u(Q))[(1 - \lambda)v^i(Q) + \lambda v^i(Q')] = V^i(Q_i) + O(hV(Q)) + O(hV(Q'))$.

vergence theorem will have been proved. Of course we must also show that, in D , $|w - g| < K$, in order that our use of the properties of a^{ij} , c^j and b^j in the derivation of inequality (15) shall be justified.

The final argument in the convergence proof is the same for both the curvilinear and rectangular nets and therefore is given in Section 4, following our description of the finite difference scheme for a rectangular net.

3. The Rectangular Net

The finite difference scheme is based on a rectangular lattice of points, $\{(x_k, y_l)\}$ such that $x_k = k \Delta x$, $y_l = l \Delta y$, where $\Delta x = 1/m$, m an integer and $\Delta y = r \Delta x$, r a positive constant.⁶

The mesh may be refined by letting m increase, but r will be kept fixed. (The range $r \leq 1/M$ will be permitted so that $|rc^j| < 1$, see equation (22).)

We shall use the notation $U_{k,l}^i$ to designate the values of the solution of the difference equations at the net points (x_k, y_l) , $u_{k,l}^i = u^i(x_k, y_l)$, $A_{k,l}^{ij} = a^{ij}(x_k, y_l, \{U_{k,l}^i\})$, $C_{k,l}^j = c^j(x_k, y_l, \{U_{k,l}^i\})$ and $B_{k,l}^i = b^i(x_k, y_l, \{U_{k,l}^i\})$. Equations (2) may then be replaced by

$$\sum_{i=1}^n A_{k,l}^{ij} \left\{ \frac{U_{k,l+1}^i - U_{k,l}^i}{\Delta y} + C_{k,l}^j \left(\frac{U_{k+1,l}^i - U_{k,l}^i}{\Delta x} \right) \right\} = B_{k,l}^i$$

$$\text{if } C_{k,l}^j \leq 0,$$

(16) or

$$\sum_{i=1}^n A_{k,l}^{ij} \left\{ \frac{U_{k,l+1}^i - U_{k,l}^i}{\Delta y} + C_{k,l}^j \left(\frac{U_{k,l}^i - U_{k-1,l}^i}{\Delta x} \right) \right\} = B_{k,l}^i$$

$$\text{if } C_{k,l}^j \geq 0, \quad \text{for } j = 1, \dots, n$$

with

$$(I.C.) \quad U_{k,0}^i = g_{k,0}^i \quad \text{for } \begin{cases} i = 1, \dots, n, \\ k = 0, \dots, m. \end{cases}$$

The equations (16) are linear equations for the unknowns $U_{k,l+1}^i$ in terms of presumably known values on the line $y = y_l = l \Delta y$. Under our hypothesis about the coefficients (i.e. $\det |a^{ij}| \neq 0$) it is clear that we may solve (16) in some restricted neighborhood of the initial line if the values $U_{k,l}^i$ are sufficiently close to g , i.e. $|U - g| < K$. The decisive step in our reasoning will be to demonstrate that in a suitable small neighborhood independent of h , $|U - g| < K$.

We observe that as expected the characteristic directions $C_{k,l}^j$ determine the

⁶For simplicity, we have here assumed that the initial curve is the segment $y = 0$, $0 \leq x \leq 1$.

domain of dependence since the decision to use either a forward or a backward difference⁸ to replace the derivative with respect to x is based on whether $C_{k,l}^i$ is negative or positive.

We shall show that the error $v_{k,l}^i = U_{k,l}^i - u_{k,l}^i$ will approach zero uniformly in some fixed region about the initial line, as $m \rightarrow \infty$. To this end we consider the growth of the error as governed by (16), for the case that $C_{k,l}^i > 0$. Again, with omission of the summation sign for the repeated superscript i , we find

$$(17) \quad A_{k,l}^{ii} v_{k,l+1}^i = (1 - rC_{k,l}^i) A_{k,l-1}^{ii} v_{k,l}^i + rC_{k,l}^i A_{k-1,l-1}^{ii} v_{k-1,l}^i + R_{k,l}^i + S_{k,l}^i$$

where

$$(18) \quad R_{k,l}^i = (1 - rC_{k,l}^i) v_{k,l}^i [A_{k,l}^{ii} - A_{k,l-1}^{ii}] + rC_{k,l}^i v_{k-1,l}^i [A_{k,l}^{ii} - A_{k-1,l-1}^{ii}]$$

and

$$(19) \quad S_{k,l}^i = -A_{k,l}^{ii} u_{k,l+1}^i + (1 - rC_{k,l}^i) A_{k,l}^{ii} u_{k,l}^i + rC_{k,l}^i A_{k,l}^{ii} u_{k-1,l}^i + \Delta y B_{k,l}^i.$$

The addition and subtraction of auxiliary terms is motivated by our desire to use the norm

$$(20) \quad E_{l+1} = \max_{(i,k)} | A_{k,l}^{ii} v_{k,l+1}^i |$$

as a measure of the error.⁹

From (17), we find, by taking absolute values, that

$$(21) \quad | A_{k,l}^{ii} v_{k,l+1}^i | \leq (1 - rC_{k,l}^i) | A_{k,l-1}^{ii} v_{k,l}^i | + rC_{k,l}^i | A_{k-1,l-1}^{ii} v_{k-1,l}^i | + | R_{k,l}^i | + | S_{k,l}^i |.$$

The restriction $r \leq 1/M$, together with the fact that $0 \leq C_{k,l}^i < M$ if $| U_{k,l} - g_{k,l} | < K$ implies that the combination

$$(22) \quad (1 - rC_{k,l}^i) | A_{k,l-1}^{ii} v_{k,l}^i | + rC_{k,l}^i | A_{k-1,l-1}^{ii} v_{k-1,l}^i |$$

is an average formed with non-negative weights whose sum is one.

At this point our reason for choosing a backward¹⁰ difference is explained.

⁸This idea of preserving a sufficiently large domain of dependence for the solution of the difference equations has been known for some time, e.g. see [4].

⁹J. Keller and P. Lax have devised a symmetric scheme in which the x -derivative is replaced by a central difference. Since they in addition use the x -average of its neighbors, $(U_{k+1,l}^i + U_{k-1,l}^i)/2$, instead of $U_{k,l}^i$ in the forward difference formula to replace U_y^i and corresponding averages for the coefficients that appear in the equation, their symmetric scheme can be used on a lattice that is staggered (i.e. (x_k, y_l) where k and l have the same parity). Such a method may have some computational advantages over the one we describe in that it may be used directly on (1), since (1) and (2) are linearly related. But the domain in which the function V can be determined by this method may be smaller than the one discussed in this paper.

¹⁰A similar norm for measuring the growth of the functions has been introduced in [7].

¹¹If $C_{k,l}^i$ were negative, we would use the forward difference and obtain the expression

$$(22)' \quad (1 + rC_{k,l}^i) | A_{k,l-1}^{ii} v_{k,l}^i | - rC_{k,l}^i | A_{k+1,l-1}^{ii} v_{k+1,l}^i |$$

Therefore

$$(23) \quad \max_{(j,k)} |A_{k,l}^{ij} v_{k,l+1}^i| \leq \max_{(j,k)} |A_{k,l-1}^{ij} v_{k,l}^i| + \max_{(j,k)} |R_{k,l}^i| + \max_{(j,k)} |S_{k,l}^i|.$$

That is, by using the definition (20)

$$(24) \quad E_{l+1} \leq E_l + \max_{(j,k)} |R_{k,l}^i| + \max_{(j,k)} |S_{k,l}^i|.$$

In the appendix, we establish the following estimates, under the assumption that $|U_{k,l} - g_{k,l}| < K$:

Lemma 1:

$$(25) \quad |R_{k,l}^i| \leq \beta_4 E_l [\Delta y + E_l + E_{l-1}]$$

and

Lemma 2:

$$(26) \quad |S_{k,l}^i| \leq \beta_5 \Delta y E_l + \beta_6 (\Delta y)^2,$$

where the β_k are non-negative constants.

Consequently, we find that if $|U_{k,l} - g_{k,l}| < K$, the growth of the error is governed by

$$(27) \quad E_{l+1} \leq E_l + N_1 \Delta y E_l + N_2 E_l^2 + N_3 E_l E_{l-1} + N_4 (\Delta y)^2, \quad \text{for } l \geq 0,$$

$$\text{and with } E_0 = 0, \quad E_{-1} = 0.$$

The right hand side of inequality (27) differs from (15) in that a term $\alpha \Delta y E_{l+1}$ is not present and a subscript $l - 1$ rather than $l + 1$ appears. We therefore treat both simultaneously in the following section.

The method involving a rectangular net is probably the one that can be adapted¹¹ to automatic machine computation.

4. Completion of the Convergence Proof

We note that for a sufficiently large positive constant, γ , the measures of error M_l and E_l both satisfy the inequality

$$(28) \quad F_{l+1} \leq F_l + \frac{\gamma}{2} (h + F_l)(2h + F_l + F_{l+1}) \quad \text{for } l \geq 0$$

$$F_0 < h, \quad F_{-1} = 0$$

(where h should be replaced by Δy when F is replaced by E .) If we set

$$h + F_l = G_l,$$

¹¹It is also possible to work with a rectangular net in a "characteristic parameter" plane, e.g. see [3].

we find

$$(29) \quad \begin{aligned} G_{l+1} &\leq G_l + \frac{\gamma}{2} G_l (G_l + G_{l+1}) \quad \text{for } l \geq 0 \\ G_0 &< 2h, \quad G_{-1} = h. \end{aligned}$$

It is clear that $G_l \leq H_l$ if we define H_l by

$$(30) \quad \begin{aligned} H_{l+1} &= H_l + \frac{\gamma}{2} H_l (H_l + H_{l+1}) \\ H_0 &= 2h, \quad H_{-1} = h. \end{aligned}$$

Now consider

$$(31) \quad T_l = \frac{2h}{1 - 2\gamma lh}$$

for $l \geq 0$ and observe that $T_{l+1} > T_l$. In addition, we note that

$$\begin{aligned} T_{l+1} - T_l &= 2h2\gamma h \frac{1}{1 - 2\gamma lh} \cdot \frac{1}{1 - 2\gamma(l+1)h} \\ &= \gamma T_l T_{l+1} > \frac{\gamma}{2} T_l (T_l + T_{l+1}). \end{aligned}$$

Hence by (30)

$$G_l \leq H_l < T_l.$$

Therefore, with the use of (31)

$$(32) \quad G_l < \frac{2h}{1 - 2\gamma lh} < 4h$$

for

$$(33) \quad l < \frac{1}{4\gamma h}$$

or equivalently

$$lh < \frac{1}{4\gamma}.$$

The range of l thus varies with h , in such a way that in a fixed region D_1 about the initial curve,

$$F_l < 3h.$$

Consequently

$$(34) \quad M_l < 3h \quad \text{or} \quad E_l < 3\Delta y$$

in a fixed region¹² about the initial curve. We shall carry out the rest of the proof for E_l . It remains to be shown that for a fixed region D contained in D_1 and R , the inequality (27) holds. In other words, we must show that with Δy sufficiently small $|U_{k,l} - g_{k,l}| < K$ for net points in D . We therefore seek a value Y , independent of Δy , such that, for Δy sufficiently small, and for $0 \leq l \leq Y/\Delta y$, $|U_{k,l} - g_{k,l}| < K$. From (33) we must restrict Y so that $Y < 1/4\gamma$. By induction on l we shall find such a Y . Assume $|U_{k,l} - g_{k,l}| < K$ for $0 \leq l \leq p$. Then we may apply (34) to obtain

$$E_{p+1} < 3\Delta y,$$

from which it follows that

$$(35) \quad |U_{k,p+1} - u_{k,p+1}| < \beta_3 \cdot E_{p+1} < 3\beta_3 \Delta y$$

(β_3 is given in (41) of appendix). Now given a positive number $\epsilon < K$, the existence theorem states that there is a region about the initial line, $0 \leq y \leq Y_1$, in which

$$(36) \quad |u_{k,l} - g_{k,l}| < K - \epsilon.$$

We now pick $\Delta y < \epsilon/3\beta_3$ and use (35) and (36) to obtain

$$|U_{k,p+1} - g_{k,p+1}| < K - \epsilon + \epsilon = K.$$

Therefore we observe that $Y = \min(1/4\gamma, Y_1)$ defines a satisfactory fixed region.

5. Round-Off Numbers

In practice, we compute numbers \bar{U} which are the solutions of the finite difference equations, rounded off to a certain number of decimal places. If we kept the number of decimal places fixed, but decreased the mesh width, we could not expect to obtain convergence. It is clear that we should increase the number of decimal places as we decrease the mesh width, in order to get an adequate representation of the solution. We shall now explain why the round-off error¹³ should be of the order $O[(\Delta y)]^2$. Let us use the notation $\bar{\Gamma}$ to denote the quantity Γ after round-off. Therefore,

$$\bar{\Gamma} = \Gamma + O[(\Delta y)^2].$$

To determine \bar{U} we must first solve either (3) or (16) with approximate coefficients and then round-off the answer. That is, for (16)

$$(37) \quad \bar{A}_i^i U_{i+1}^i = T^i(\bar{A}_i, \bar{B}_i, \bar{C}_i, \bar{U}_i)$$

$$\bar{U}_{i+1}^i = U_{i+1}^i + O(\Delta y)^2.$$

¹²In [7], another method of estimating the solution of inequality (28) is given.

¹³We indicate the argument for the rectangular net. In the case of the curvilinear net, the round-off error should be of the order $O(h^2)$.

Both operations may be combined into the single system

$$(38) \quad \bar{A}_i^{ii} \bar{U}_{i+1}^i = T^i(\bar{A}_i, \bar{B}_i, \bar{C}_i, \bar{U}_i) + O[(\Delta y)^2].$$

From this point on the analysis of the error proceeds as before with $v = \bar{U} - u$. The variations in the estimates can all be put into the term $N_4(\Delta y)^2$, again subject to the restriction that $|\bar{U}_{k,i} - g_{k,i}| < K$. Consequently we have established convergence subject to the requirement that the round-off be of the order $O[(\Delta y)^2]$. Specifically, when the mesh width is decreased in the ratio $1/\nu$, the number of decimal places should be increased by $2 \log_{10} \nu$ digits.

6. Summary

We have shown that the mesh width ratio,¹⁴ $r = \Delta y/\Delta x$, should be chosen in such a way that the domain of dependence of any point in the mesh as given by the difference equations, is not less than the domain of dependence determined by the differential equations. We have seen that the choice of difference quotients (forward or backward) should be made with the idea of preserving the domain of dependence.

Although we have only established the above criteria as sufficient conditions for convergence, it is quite easy to see that all things being equal they are necessary. Furthermore, we have shown that the round-off error should be of the order $O[(\Delta y)^2]$. This requirement is not a necessary condition for convergence—but it is close enough for practical purposes, since it is clear that the round-off error should be smaller than $O(\Delta y)$. That is, the addition of $2k$ decimal digits instead of at least k is not an enormously wasteful operation for k small. We might remark that after a calculation of \bar{U} has been completed for a certain mesh, it is possible to estimate the optimum number of decimal digits that should have been kept, that is, the optimum number can be determined principally from an estimate of N_4 (the Lipschitz constant etc.), which should be possible after a single calculation of \bar{U} .

We may point out that the proof of convergence¹⁵ given here applies almost without change to mixed initial and boundary value problems for the system (2). Modifications are needed when free boundaries are to be determined in the problem (e.g. shocks, contact discontinuities, etc.).

¹⁴In the case of the curvilinear net, the mesh width ratio is quite variable, but subject to the same requirements regarding domains of dependence that we impose for the rectangular net.

¹⁵A method of practical utility in testing for local stability, attributed to J. von Neumann is described in [8].

Appendix

We shall establish the estimates stated in the inequalities (25) and (26). From the definition (18) we obtain readily that

$$(39) \quad \begin{aligned} R_{k,l}^i &= (1 - rC_{k,l}^i)v_{k,l}^i[(A_{k,l}^{ii} - a_{k,l}^{ii}) - (A_{k,l-1}^{ii} - a_{k,l-1}^{ii}) + (a_{k,l}^{ii} - a_{k,l-1}^{ii})] \\ &\quad + rC_{k,l}^iv_{k-1,l}^i[(A_{k,l}^{ii} - a_{k,l}^{ii}) - (A_{k-1,l-1}^{ii} - a_{k-1,l-1}^{ii}) + (a_{k,l}^{ii} - a_{k-1,l-1}^{ii})]. \end{aligned}$$

By making use of the Lipschitz condition satisfied by a^{ii} and u^i , we obtain

$$(40) \quad \begin{aligned} |R_{k,l}^i| &\leq (1 - rC_{k,l}^i)\beta_1(\sum_i |v_{k,l}^i|)[\Delta y + \sum_i |v_{k,l}^i| + \sum_i |v_{k,l-1}^i|] \\ &\quad + rC_{k,l}^i\beta_2(\sum_i |v_{k,l}^i|)[\Delta y + \sum_i |v_{k,l}^i| + \sum_i |v_{k-1,l-1}^i|], \end{aligned}$$

with suitable positive constants β_i . Now we observe that

$$(41) \quad \sum_i |v_{k,l}^i| \leq \beta_3 \sum_i |A_{k,l-1}^{ii}v_{k,l}^i| \leq \beta_3 E_l,$$

since $\det |A_{k,l-1}^{ii}| \neq 0$, for $|U_{k,l-1} - g_{k,l-1}| < K$. Therefore we have established

Lemma 1:

$$|R_{k,l}^i| \leq \beta_4 E_l [\Delta y + E_l + E_{l-1}].$$

From the definition (19) we see that

$$(42) \quad \begin{aligned} S_{k,l}^i &= (A_{k,l}^{ii} - a_{k,l}^{ii})(u_{k,l}^i - u_{k,l+1}^i) + rC_{k,l}^i(u_{k-1,l}^i - u_{k,l}^i) \\ &\quad + a_{k,l}^{ii}[(u_{k,l}^i - u_{k,l+1}^i) + (C_{k,l}^i - c_{k,l}^i)(u_{k-1,l}^i - u_{k,l}^i)] \\ &\quad + a_{k,l}^{ii}rc_{k,l}^i(u_{k-1,l}^i - u_{k,l}^i) + \Delta y(B_{k,l}^i - b_{k,l}^i) \\ &\quad + \Delta y b_{k,l}^i. \end{aligned}$$

Now if we observe that $rC_{k,l}^i$ is bounded for $|U_{k,l} - g_{k,l}| < K$ and that $\{u^i\}$ is a solution of (2), which has first derivatives that satisfy a Lipschitz condition, then by inequality (41) we obtain

Lemma 2:

$$|S_{k,l}^i| \leq \beta_5 \Delta_y E_l + \beta_6 (\Delta y)^2.$$

Let us remark that in both (25) and (26) the constants β_k depend only upon a neighborhood of the initial data that is fixed throughout our discussion, $|U_{k,l} - g_{k,l}| < K$.

BIBLIOGRAPHY

- [1] Beckert, Herbert, *Über quasilineare hyperbolische Systeme partieller Differentialgleichungen erster Ordnung mit zwei unabhängigen Variablen. Das Anfangswertproblem, die gemischte Anfangs-Randwertaufgabe, das charakteristische Problem*, Berichte Über Die Verhandlungen Der Sächsischen Akademie Der Wissenschaften Zu Leipzig, Mathematisch-naturwissenschaftliche Klasse, Volume 97, No. 5, 1950, p. 68.
- [2] Cinquini-Cibrario, Maria, *Sopra il problema di Cauchy per i sistemi di equazioni alle derivate parziali del primo ordine*, Rendiconti del Seminario Matematico dell'Università di Padova, Volume 17, 1948, pp. 75-96.
- [3] Courant, Richard, and Friedrichs, K. O., *Supersonic Flow and Shock Waves*, Interscience, New York, 1948.
- [4] Courant, Richard, Friedrichs, K. O., and Lewy, Hans, *Über die partiellen Differenzengleichungen der Mathematischen Physik*, Mathematische Annalen, Volume 100, 1928, pp. 32-74.
- [5] Courant, Richard, and Lax, Peter, *On nonlinear partial differential equations with two independent variables*, Communications on Pure and Applied Mathematics, Volume 2, No. 3, 1949, pp. 255-273.
- [6] Friedrichs, K. O., *Nonlinear hyperbolic differential equations for functions of two independent variables*, American Journal of Mathematics, Volume 70, 1948, pp. 555-589.
- [7] Friedrichs, K. O., and Lewy, Hans, *Das Anfangswertproblem einer beliebigen nichtlinearen hyperbolischen Differentialgleichung beliebiger Ordnung in zwei Variablen. Existenz, Eindeutigkeit und Abhängigkeitsbereich der Lösung*, Mathematische Annalen, Volume 99, 1928, pp. 200-221.
- [8] O'Brien, George G., Hyman, Morton A., and Kaplan, Sidney, *A study of the numerical solution of partial differential equations*, Journal of Mathematical Physics, Volume 29, 1951, pp. 223-251.