

## CHAPTER 6

# Finite difference methods for hyperbolic equations

### 1. Scalar equations

**1.1. Constant velocity advection in one dimension.** The simplest example of a hyperbolic equation is the constant velocity advection equation

$$(1.1) \quad q_t + u q_x = 0$$

with some initial condition  $q(x, t = 0) = q_0(x)$ . The equation can be solved along the entire real axis in  $x$  or some portion thereof. In numerical work we always have a finite subdomain which we shall conveniently choose as  $[0, 2\pi]$  with a view to applying Fourier analysis later on. When using a finite subdomain the question of boundary conditions arises which we shall postpone by considering periodic boundary conditions  $q(x + 2\pi, t) = q(x, t)$ .

1.1.1. *Exact solution by characteristics.* A first attack on finding the solution to (1.1) is to try to reduce it to a simpler problem. One can ask whether there is any subdomain over which the equation can be cast in a simpler form. For instance we can inquire whether there are any particular curves within the  $(x, t)$  plane over which the equation simplifies. A general curve  $\Gamma$  of curvilinear parameter is given by

$$(1.2) \quad \Gamma : x = x(s), t = t(s)$$

and the infinitesimal change in  $q$  when going along  $\Gamma$  is

$$(1.3) \quad \frac{dq}{ds} = \frac{\partial q}{\partial t} \frac{dt}{ds} + \frac{\partial q}{\partial x} \frac{dx}{ds}$$

Comparing (1.3) with (1.1) we see that if we impose

$$(1.4) \quad \frac{dt}{ds} = 1, \quad \frac{dx}{ds} = u$$

then by (1.1) we must have that

$$(1.5) \quad \frac{dq}{ds} = 0.$$

This means that  $q$  is constant along the curves  $\Gamma$  defined by (1.4) which are  $x = ut + C$ . The curves are shown in Fig. (1)

1.1.2. *Finite difference methods.* We can construct numerical methods for (1.1) by the same approaches used for the heat equation.

FIGURE 1. Characteristic curves for  $q_t + q_x = 0$ .

Semi-discretization. Define a computational grid  $x_j = jh$ ,  $h = 2\pi/(M+1)$ ,  $t^n = nk$  with step size  $h, k$  in space and time. Define  $Q_j(t)$  to be the restriction of  $q(x, t)$  to  $x = x_j$

$$(1.6) \quad Q_j(t) = q(x_j, t), \quad j = 0, 1, \dots, M+1.$$

We can choose some approximation of the  $x$  derivative. For instance approximating

$$(1.7) \quad \frac{dq(x_j, t)}{dx} \cong \frac{\delta}{h} Q_j = \frac{Q_{j+1}(t) - Q_{j-1}(t)}{2h}$$

leads to the ODE system

$$(1.8) \quad \frac{d}{dt} \mathbf{Q} = -\frac{u}{2h} B \mathbf{Q}$$

with

$$(1.9) \quad \mathbf{Q} = [Q_1 \quad Q_2 \quad \dots \quad Q_M]^T$$

$$(1.10) \quad B = \begin{bmatrix} 0 & 1 & & & -1 \\ -1 & 0 & 1 & & \\ & -1 & 0 & 1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 0 & 1 \\ 1 & & & & -1 & 0 \end{bmatrix}$$

We can now try various ODE schemes to solve (1.8). Using Euler's method would lead to a FTCS scheme

$$(1.11) \quad \mathbf{Q}^{n+1} = \mathbf{Q}^n - \frac{uk}{2h} B \mathbf{Q}^n = \left( I - \frac{uk}{2h} B \right) \mathbf{Q}^n.$$

If instead of Euler's scheme we use the midpoint method we obtain the update formula

$$(1.12) \quad \mathbf{Q}^{n+1} = \mathbf{Q}^{n-1} - \frac{uk}{h} B \mathbf{Q}^n$$

known as the *leap-frog* or *Dufort-Frankel* method.

Full discretization. Instead of the semi-discretization or method of lines approach we can also directly discretize both the space and time derivatives appearing in (1.1). A first-order forward in time, second-order centered in space discretization would lead to the FTCS scheme

$$(1.13) \quad Q_j^{n+1} = Q_j^n - \frac{uk}{2h} (Q_{j+1}^n - Q_{j-1}^n)$$

A modification of (1.13) of historical relevance is the Lax-Friedrichs scheme

$$(1.14) \quad Q_j^{n+1} = \frac{1}{2} (Q_{j+1}^n + Q_{j-1}^n) - \frac{uk}{2h} (Q_{j+1}^n - Q_{j-1}^n).$$

This scheme is obtained by replacing  $Q_j^n$  with its arithmetic average using values to the left and to the right. Since the formula has not been derived from discretizations of the derivative which we know to be consistent with the original equation it is useful to determine the truncation error. The exact advection operator is

$$(1.15) \quad D = \frac{\partial}{\partial t} + u \frac{\partial}{\partial x}$$

and we have  $Dq = 0$  according to (1.1). Our approximation of this operator is

$$(1.16) \quad \tilde{D}q(x_j, t^n) = \frac{q_j^{n+1} - \frac{1}{2}(q_{j+1}^n + q_{j-1}^n)}{k} + \frac{u}{2h}(q_{j+1}^n - q_{j-1}^n)$$

with  $q_j^n = q(x_j, t^n)$ . The truncation error is therefore

$$(1.17) \quad \tau_j^n = (\tilde{D} - D)q(x_j, t^n) = \frac{1}{k}q_j^{n+1} - \frac{1}{2k}(q_{j+1}^n + q_{j-1}^n) + \frac{u}{2h}(q_{j+1}^n - q_{j-1}^n)$$

We can now carry out a Taylor's series expansion around  $(x_j, t^n)$ . To simplify notation  $q$  and its derivatives will be understood to be evaluated at  $(x_j, t^n)$  if not explicitly shown otherwise

$$(1.18) \quad \tau_j^n = \frac{1}{k} \left( q + kq_t + \frac{k^2}{2}q_{tt} + \dots \right) - \frac{1}{2k} \left( q + hq_x + \frac{h^2}{2}q_{xx} + \dots + q - hq_x + \frac{h^2}{2}q_{xx} + \dots \right)$$

$$(1.19) \quad + \frac{u}{2h} \left( q + hq_x + \frac{h^2}{2}q_{xx} + \frac{h^3}{6}q_{xxx} + \dots - q + hq_x - \frac{h^2}{2}q_{xx} + \frac{h^3}{6}q_{xxx} + \dots \right)$$

which gives

$$(1.20) \quad \tau_j^n = q_t + \frac{k}{2}q_{tt} + \dots - \frac{h^2}{2k}q_{xx} + \dots + uq_x + \frac{uh^2}{12}q_{xxx} + \dots$$

Using (1.1) leads to the leading order error

$$(1.21) \quad \tau_j^n \cong \frac{k}{2}q_{tt} + \frac{uh^3}{12}q_{xxx} = O(k, h^2)$$

The analysis shows that the scheme is first order in time and second order in space. The Lax-Friedrichs scheme is consistent, i.e.

$$(1.22) \quad \lim_{k, h \rightarrow 0} \tau_j^n = 0.$$

Instead of centered finite differences, other approximations may be introduced. A good choice would be to use one-sided finite differences. This would take into account what we know from the exact solution to the advection equation: information travels along the characteristic lines. It would make sense to use finite differences which have a stencil that mimics this behavior. If  $u > 0$  we would use left-sided differences and for  $u < 0$  we would use right-sided differences. A first order approximation would be

$$(1.23) \quad Q_j^{n+1} = \begin{cases} Q_j^n - \frac{uk}{h}(Q_j^n - Q_{j-1}^n) & u \geq 0 \\ Q_j^n - \frac{uk}{h}(Q_{j+1}^n - Q_j^n) & u \leq 0 \end{cases}$$

This is known, naturally enough, as the *upwind scheme*.

Taylor series approach. A procedure useful in deriving higher order finite difference approximations for (1.1) is the Taylor series approach. In practical work it is economical to only store two time levels at any given stage in the computation. The general prescription for attaining higher order for the time derivative would involve keeping more terms from the operator series

$$(1.24) \quad \frac{\partial}{\partial t} = \frac{1}{k} \left( \Delta_+ - \frac{\Delta_+^2}{2} + \frac{\Delta_+^3}{3} - \dots \right).$$

This would be inconvenient since the time stencil of the scheme would become wider and we would need to store more than two time levels. We can however use (1.1) to convert time derivatives into spatial derivatives

$$(1.25) \quad q_t = -uq_x$$

$$(1.26) \quad q_{tt} = \frac{\partial}{\partial t} (q_t) = -\frac{\partial}{\partial t} (uq_x) = -u \frac{\partial}{\partial x} (q_t) = u \frac{\partial}{\partial x} (uq_x) = u^2 q_{xx}$$

The Taylor series approach can now be applied to obtain as high an order of approximation in time as needed

$$(1.27) \quad q(t+k) = q + kq_t + \frac{k^2}{2} q_{tt} + \frac{k^3}{6} q_{ttt} + \dots$$

As an example, let us construct a second order scheme by truncating

$$(1.28) \quad q(t+k) \cong q + kq_t + \frac{k^2}{2} q_{tt}$$

We now replace the time derivatives with spatial derivative

$$(1.29) \quad q(t+k) \cong q - ukq_x + \frac{u^2 k^2}{2} q_{xx}$$

and use second order accurate, centered finite differences to approximate the spatial derivatives. The resulting scheme is

$$(1.30) \quad Q_j^{n+1} = Q_j^n - \frac{uk}{2h} (Q_{j+1}^n - Q_{j-1}^n) + \frac{u^2 k^2}{2h^2} (Q_{j+1}^n - 2Q_j^n + Q_{j-1}^n)$$

This is known as the *Lax-Wendroff scheme*.

Instead of centered finite differences we might want to use one-sided formulas to take into account the direction of the characteristics of (1.1). Let us construct a second order accurate one sided approximation using (??). One-sided differences can be obtained to arbitrary order of accuracy using the series

$$(1.31) \quad \frac{\partial}{\partial x} = \frac{1}{h} \left( \Delta_{x+} - \frac{\Delta_{x+}^2}{2} + \frac{\Delta_{x+}^3}{3} - \dots \right)$$

$$(1.32) \quad = \frac{1}{h} \left( \Delta_{x-} + \frac{\Delta_{x-}^2}{2} + \frac{\Delta_{x-}^3}{3} - \dots \right)$$

with the finite difference operators defined by

$$(1.33) \quad \Delta_{x+} q(x, t) = q(x+h, t) - q(x, t)$$

$$(1.34) \quad \Delta_{x-} q(x, t) = q(x, t) - q(x-h, t)$$

Let us assume that  $u \geq 0$  and therefore that we will be using backward differences so that the computational stencil mimics the true domain of dependence. A second

order accurate approximation of  $q_x$  is given by

(1.35)

$$\begin{aligned} q_x(x_j, t^n) &= \frac{\partial q}{\partial x}(x_j, t^n) \cong \frac{1}{h} \left( \Delta_{x-} + \frac{\Delta_{x-}^2}{2} \right) q(x_j, t^n) \\ (1.36) \quad &= \frac{1}{h} \left[ q(x_j, t^n) - q(x_{j-1}, t^n) + \frac{1}{2} (q(x_j, t^n) - 2q(x_{j-1}, t^n) + q(x_{j-2}, t^n)) \right] \end{aligned}$$

$$(1.37) \quad \cong \frac{1}{h} \left[ Q_j^n - Q_{j-1}^n + \frac{1}{2} (Q_j^n - 2Q_{j-1}^n + Q_{j-2}^n) \right]$$

$$(1.38) \quad = \frac{1}{2h} (3Q_j^n - 4Q_{j-1}^n + Q_{j-2}^n)$$

The second derivative is obtained from

(1.39)

$$\begin{aligned} \frac{\partial^2}{\partial x^2} &= \frac{\partial}{\partial x} \frac{\partial}{\partial x} = \frac{1}{h} \left( \Delta_{x-} + \frac{\Delta_{x-}^2}{2} + \frac{\Delta_{x-}^3}{3} - \dots \right) \frac{1}{h} \left( \Delta_{x-} + \frac{\Delta_{x-}^2}{2} + \frac{\Delta_{x-}^3}{3} - \dots \right) \\ (1.40) \quad &= \frac{1}{h^2} \left( \Delta_{x-}^2 + \Delta_{x-}^3 + \frac{11}{12} \Delta_{x-}^3 \right) \end{aligned}$$

Note that in (??) we have neglected terms of  $O(k^3)$ . The exact solution of the advection equation is  $q(x, t) = q_0(x - ut)$ . The  $x - ut$  argument suggests that similar step sizes should be used for  $x$  and  $t$ . This will be confirmed by our stability analysis below. So let us assume that  $k = O(h)$ . Since  $q_{xx}$  already has a  $k^2$  factor in (??) we only need an  $O(h)$  approximation of  $q_{xx}$ . The leading order error term from  $k^2 q_{xx}$  will then be of  $O(k^2 h) = O(k^3) = O(h^3)$ . We can therefore truncate the series (1.40) to just the first term and approximate

$$(1.41) \quad \frac{\partial^2 q}{\partial x^2}(x_j, t^n) \cong \frac{1}{h^2} \Delta_{x-}^2 q(x_j, t^n) = \frac{1}{h^2} [q(x_j, t^n) - 2q(x_{j-1}, t^n) + q(x_{j-2}, t^n)]$$

$$(1.42) \quad \cong \frac{1}{h^2} (Q_j^n - 2Q_{j-1}^n + Q_{j-2}^n)$$

Note that this is different from the procedure used in deriving the Lax-Wendroff scheme where a second order accurate expression of  $q_{xx}$  was used. The reason is that in the Lax-Wendroff scheme the computational stencil already included  $Q_{j-1}^n$ ,  $Q_j^n$ ,  $Q_{j+1}^n$  from the approximation of the first derivative  $q_x$ . Since the second order accurate approximation of  $q_{xx}$  does not widen the stencil there is no penalty in using the more accurate, second order approximation of  $q_{xx}$ . In the one-sided scheme we are deriving here however, using a second order accurate approximation of  $q_{xx}$  would involve widening the computational stencil to include  $Q_{j-3}^n$ . This increases the arithmetic cost of applying the formula without noticeable gain so we choose to use an  $O(h)$  approximation of  $q_{xx}$ . Combining the above results we obtain

$$(1.43) \quad Q_j^{n+1} = Q_j^n - \frac{uk}{2h} (3Q_j^n - 4Q_{j-1}^n + Q_{j-2}^n) + \frac{1}{2} \left( \frac{uk}{h} \right)^2 (Q_j^n - 2Q_{j-1}^n + Q_{j-2}^n)$$

for  $u > 0$ , which is known as the *Beam-Warming scheme*.

1.1.3. *Stability analysis.* We now turn to the analysis of the stability of the various schemes introduced above. The analysis can be done using the techniques for systems of ODE's or using Von Neumann analysis. We shall carry out both procedures.

Semi-discretized system. The matrix  $B$  arising in the semi-discretized approach is skew-symmetric and will have purely imaginary eigenvalues. We can check this by explicitly calculating the eigenvalues. As usual, we guess that

$$(1.44) \quad \mathbf{W}_p = \begin{bmatrix} e^{iph} & e^{ip2h} & \dots & e^{ipjh} & \dots & e^{ipMh} \end{bmatrix}$$

will be an eigenvector since  $B$  discretizes a derivation operator. Computing the  $j^{th}$  component of  $B\mathbf{W}_p$  we get

$$(1.45) \quad (B\mathbf{W}_p)_j = e^{ip(j+1)h} - e^{ip(j-1)h} = 2i \sin ph \, e^{ipjh} = 2i \sin ph \, (\mathbf{W}_p)_j$$

so the eigenvalue associated with  $\mathbf{W}_p$  is

$$(1.46) \quad \lambda_p = 2i \sin ph$$

and is indeed purely imaginary.

To establish the stability region for the FTCS method (??) we use the eigenvalues  $\lambda$  of  $B$  in the criterion

$$(1.47) \quad |1 + z| \leq 1$$

with  $z = k\lambda$ . It is immediately apparent that the scheme will be unconditionally unstable because  $\lambda$  is purely imaginary  $\lambda = ai$  so

$$(1.48) \quad |1 + z| = \sqrt{1 + (ka)^2} > 1$$

for all  $a > 0$ .

The interval of stability for the midpoint scheme is  $\text{Re } z = 0$ ,  $|\text{Im}(z)| \leq 1$ . Here we would have

$$(1.49) \quad z = -\frac{uk}{h} i \sin ph$$

and the method is stable for

$$(1.50) \quad \left| \frac{uk}{h} \right| \leq 1 .$$

Von Neumann analysis. Obtaining analytical expression for the matrices arising from the semi-discretized approach becomes increasingly difficult as we use more accurate approximations of the derivatives in the PDE or study PDE's more complex than the advection equation. Von Neumann analysis is typically simpler to apply. We start by determining the stability region for the FTCS scheme (1.13). Substituting a typical wavemode  $Q^n = \hat{Q}^n e^{i\xi jh}$  we obtain

$$(1.51) \quad \hat{Q}^{n+1} e^{i\xi jh} = \hat{Q}^n e^{i\xi jh} - \frac{uk}{2h} \left( \hat{Q}^n e^{i\xi(j+1)h} - \hat{Q}^n e^{i\xi(j-1)h} \right) .$$

The amplification ratio is

$$(1.52) \quad G = \frac{\hat{Q}^{n+1}}{\hat{Q}^n} = 1 - \frac{uk}{h} i \sin \xi h$$

which is always greater than 1

$$(1.53) \quad |G| \geq 1 .$$

Thus the scheme is unconditionally unstable as we expected from the semi-discretized stability analysis done above.

For the Lax-Friedrichs scheme we obtain

$$(1.54) \quad \hat{Q}^{n+1} e^{i\xi j h} = \frac{1}{2} \left( \hat{Q}^n e^{i\xi(j+1)h} + \hat{Q}^n e^{i\xi(j-1)h} \right) - \frac{uk}{2h} \left( \hat{Q}^n e^{i\xi(j+1)h} - \hat{Q}^n e^{i\xi(j-1)h} \right)$$

and the amplification factor is

$$(1.55) \quad G = \cos \xi h - \frac{uk}{h} i \sin \xi h$$

Let us introduce the notation

$$(1.56) \quad \nu = \frac{uk}{h}, \quad \theta = \xi h.$$

The stability condition is that

$$(1.57) \quad |G| = \cos^2 \theta + \nu^2 \sin^2 \theta \leq 1 = \cos^2 \theta + \sin^2 \theta$$

from where we obtain

$$(1.58) \quad (\nu^2 - 1) \sin^2 \theta \leq 0.$$

The inequality is satisfied for

$$(1.59) \quad |\nu| \leq 1.$$

The quantity  $\nu$  that appears repeatedly in analysis of numerical schemes for the advection equation is known as the *Courant-Friedrichs-Lewy* number or more concisely as the *CFL number*. We say that the Lax-Friedrichs scheme is stable for CFL numbers up to 1, it being implicitly understood that we're considering the absolute value of the velocity  $|u|$ . From the stability criterion we obtain a bound on the time step that we can use in the Lax-Friedrichs scheme

$$(1.60) \quad k \leq \frac{h}{|u|}.$$

For the Lax-Wendroff scheme the amplification ratio is

$$(1.61) \quad G = 1 - \nu i \sin \theta + \nu^2 (\cos \theta - 1)$$

We have

$$(1.62) \quad |G| = 1 + 2\nu^2 (\cos \theta - 1) + \nu^4 (\cos \theta - 1)^2 + \nu^2 \sin^2 \theta$$

$$(1.63) \quad = 1 - 4\nu^2 \sin^2 \frac{\theta}{2} \left( 1 - \cos^2 \frac{\theta}{2} \right) + 4\nu^4 \sin^4 \frac{\theta}{2}$$

The stability condition is  $|G| \leq 1$  leads to

$$(1.64) \quad (\nu^2 - 1) \sin^2 \frac{\theta}{2} \leq 0$$

so again the domain of stability is

$$(1.65) \quad |\nu| \leq 1.$$

Lax-Wendroff is a more efficient scheme than Lax-Friedrichs since we obtain  $O(h^2, k^2)$  precision as opposed to  $O(h, k^2)$  under the same time step restriction  $k \leq h/|u|$ .

Turning now to the one-sided schemes, for upwind when  $u > 0$  we have

$$(1.66) \quad G = 1 - \nu (1 - e^{-i\theta})$$

$$(1.67) \quad |G| = 1 - 2\nu (1 - \cos \theta) + \nu^2 (1 - \cos \theta)^2 + \nu^2 \sin^2 \theta$$

FIGURE 2. Amplification factor  $|G(\nu, \theta)|$  for the Beam-Warming scheme evaluated at  $\theta = m\pi/8$ ,  $m = 0, 1, \dots, 16$ .

The stability condition  $|G| \leq 1$  leads to

$$(1.68) \quad -2\nu(1 - \cos \theta) + \nu^2(1 - \cos \theta)^2 + \nu^2 \sin^2 \theta \leq 0$$

which can be rewritten in terms of the half-angle  $\theta/2$  to give

$$(1.69) \quad -4\nu \sin^2 \frac{\theta}{2} + 4\nu^2 \sin^4 \frac{\theta}{2} + 4\nu^2 \sin^2 \frac{\theta}{2} \cos^2 \frac{\theta}{2} \leq 0$$

and finally

$$(1.70) \quad \nu(\nu - 1) \leq 0$$

so the stability region is again  $\nu \leq 1$ .

For the Beam-Warming scheme we have

$$(1.71) \quad G = 1 - \frac{\nu}{2}(3 - 4e^{-i\theta} + e^{-2i\theta}) + \frac{1}{2}\nu^2(1 - 2e^{-i\theta} + e^{-2i\theta}) .$$

Notice that as we look at more complicated schemes the amplification factors become increasingly difficult to evaluate analytically. We can however use a numerical evaluation of  $G(\nu, \theta)$  to generate plots such as Fig. 2. From the plot we deduce that the stability region is  $\nu \leq 2$ .

1.1.4. *Lax equivalence theorem.* The importance of establishing consistency and stability for a finite difference scheme for the advection equation is that these two properties guarantee convergence by the Lax equivalence theorem.

THEOREM 4. *A finite difference scheme for a linear PDE is convergent if the scheme is consistent with the PDE and it is stable.*

Convergence means that

$$(1.72) \quad \lim_{k, h \rightarrow 0} Q_j^n = q(x_j, t^n)$$

where  $k, h$  go to zero in accordance with the stability criterion for the scheme. Convergence is obtained when the scheme is consistent, i.e. the truncation error goes to zero

$$(1.73) \quad \lim_{k, h \rightarrow 0} \tau_j^n = 0$$

and the step sizes satisfy the stability criterion.

1.1.5. *Modified equations.* We have established a number of methods for solving the advection equation (1.1). Up to now we have characterized the error of any one scheme by its truncation error. Though indicative of the overall quality of an approximation, the precise nature of the error in the scheme is not apparent. It has proved very fruitful in the development of better methods to more accurately describe how a numerical approximation differs from the exact solution. A question one can ask is whether a given numerical scheme is perhaps a more accurate discretization of another PDE than the one it was originally designed for. Let us exemplify using the upwind scheme for the advection equation with  $u > 0$

$$(1.74) \quad Q_j^{n+1} = Q_j^n - \nu(Q_j^n - Q_{j-1}^n) .$$



We know that this scheme is  $O(k, h)$  accurate for the equation  $q_t + uq_x = 0$ . Suppose that the scheme is an exact discretization of some unknown PDE  $Ls = 0$  with  $L$  an unknown differential operator and  $s = s(x, t)$ . Then we would have

$$(1.75) \quad s(x, t + k) = s(x, t) - \nu [s(x, t) - s(x - h, t)] ,$$

exactly. Let us carry out Taylor series expansion of  $s$  around  $(x, t)$

$$(1.76) \quad s + ks_t + \frac{k^2}{2}s_{tt} + \frac{k^3}{6}s_{ttt} + \dots = s - \frac{uk}{h} \left[ hs_x - \frac{h^2}{2}s_{xx} + \frac{h^3}{6}s_{xxx} - \dots \right] .$$

To obtain a more concise notation the function arguments have been dropped. We obtain

$$(1.77) \quad s_t + us_x = -\frac{k}{2}s_{tt} + \frac{uh}{2}s_{xx} - \frac{k^2}{6}s_{ttt} - \frac{uh^2}{6}s_{xxx} + \dots$$

This is of the form  $As = E_{(h,k)}s$  with  $A$  the advection operator  $A = \partial_t + u\partial_x$  and  $E_{(h,k)}$  an operator giving the deviation of the modified equation from the advection equation. Note that if  $k = h = 0$  we obtain the advection equation for which the scheme (1.74) is  $O(h, k)$  accurate. We can interpret (1.77) as stating that the scheme (1.74) is:

(1) first order accurate for

$$(1.78) \quad s_t + us_x = 0$$

(2) second order accurate for

$$(1.79) \quad s_t + us_x = -\frac{k}{2}s_{tt} + \frac{uh}{2}s_{xx}$$

(3) third order accurate for

$$(1.80) \quad s_t + us_x = -\frac{k}{2}s_{tt} + \frac{uh}{2}s_{xx} - \frac{k^2}{6}s_{ttt} - \frac{uh^2}{6}s_{xxx}$$

The equations obtained above are called *modified equations*. These statements can be verified by explicit computation of the truncation error. For example let us compute the truncation error in applying (1.74) to (??)

$$(1.81) \quad \tau_j^n = (\tilde{D} - D) s(x_j, t^n)$$

The finite difference approximation operator is

$$(1.82) \quad \tilde{D}s(x_j, t^n) = \frac{s_j^{n+1} - s_j^n}{k} + \frac{u}{h} (s_j^n - s_{j-1}^n)$$

The exact operator for the modified equation (??) is

$$(1.83) \quad D = \frac{\partial}{\partial t} + u \frac{\partial}{\partial x} + \frac{k}{2} \frac{\partial^2}{\partial t^2} - \frac{uh}{2} \frac{\partial^2}{\partial x^2}$$

and we have  $Ds = 0$ . We now expand  $s_j^{n+1}$ ,  $s_{j-1}^n$ ,  $s_{j+1}^n$  from (1.82) around  $(x_j, t^n)$  and obtain

$$(1.84) \quad \tau_j^n = \frac{1}{k} \left( s + ks_t + \frac{k^2}{2}s_{tt} + \dots \right) - \frac{1}{k} s + \frac{u}{h} \left( s - s + hs_x - \frac{h^2}{2}s_{xx} + \frac{h^3}{6}s_{xxx} + \dots \right)$$

$$(1.85) \quad \tau_j^n = s_t + \frac{k}{2}s_{tt} + us_x + \frac{uh}{2}s_{xx} + O(k^2, h^2) = Ds + O(k^2, h^2) = O(k^2, h^2)$$

so the truncation error is indeed of second order.

Now let us show the benefits of looking at the modified equation by using (??) for which the upwind scheme (1.77) is second order accurate. First we recast (??) so as to eliminate higher order derivatives in time. We can rewrite (??) as

$$(1.86) \quad s_t = -us_x - \frac{k}{2}s_{tt} + \frac{uh}{2}s_{xx}$$

and differentiate with respect to  $t$  to obtain

$$(1.87) \quad s_{tt} = -us_{xt} - \frac{k}{2}s_{ttt} + \frac{uh}{2}s_{xxt}$$

Replacing (1.87) in (??) gives

$$(1.88) \quad s_t + us_x = -\frac{k}{2} \left( -us_{xt} - \frac{k}{2}s_{ttt} + \frac{uh}{2}s_{xxt} \right) + \frac{uh}{2}s_{xx}$$

$$(1.89) \quad = \frac{uk}{2}s_{xt} + \frac{uh}{2}s_{xx} + O(k^2, h^2, kh)$$

We can neglect the higher order terms since this is consistent with the order of accuracy used in obtaining (??). Differentiating (??) with respect to  $x$  yields

$$(1.90) \quad s_{tx} = -us_{xx} - \frac{k}{2}s_{ttx} + \frac{uh}{2}s_{xxx}$$

and replacing this in (??) gives

$$(1.91) \quad s_t + us_x = -\frac{u^2k}{2}s_{xx} + \frac{uh}{2}s_{xx} = \frac{uh}{2}(1 - \nu)s_{xx}$$

Equation (1.91) is the usual way to express the modified equation for the upwind scheme applied to the advection equation to second order. It shows that the upwind scheme does indeed model the advection equation in the limit  $h \rightarrow 0$ . For finite step sizes however the upwind scheme more accurately models the equation (1.91). The difference between (1.91) and the advection equation is the term

$$(1.92) \quad \frac{uh}{2}(1 - \nu)s_{xx}$$

Note that this is a diffusive term whose effect is to smooth out any variations in  $s(x, t)$  as long as  $|1 - \nu| \geq 0$  as has been seen in the study of the heat equation. The condition  $|1 - \nu| \geq 0$  is exactly the stability criterion for the upwind scheme. Indeed if  $\nu > 1$  then we would obtain a negative diffusion coefficient for which the initial value problem is ill-defined. We can see that at exactly  $\nu = 1$  there is no diffusion indicating that for  $\nu = 1$  the upwind scheme achieves higher order accuracy for the advection equation. When  $\nu < 1$  the error in the upwind scheme with respect to the true solution  $q(x, t)$  of the advection equation will be diffusive: gradients will be smoothed out instead of being simply advected.

Now that we have seen the nature of the error introduced by the upwind scheme applied to the advection equation, we can also use this information to derive better schemes. Since the error is known to be given by (1.92) we can change the upwind scheme

$$(1.93) \quad Q_j^{n+1} = Q_j^n - \nu(Q_j^n - Q_{j-1}^n)$$

to counteract the known error by including a discretization of (1.92)

$$(1.94) \quad Q_j^{n+1} = Q_j^n - \nu(Q_j^n - Q_{j-1}^n) - k \frac{uh}{2}(1 - \nu) \frac{Q_{j+1}^n - 2Q_j^n + Q_{j-1}^n}{h^2}$$

Working this through leads to the scheme

$$(1.95) \quad Q_j^{n+1} = Q_j^n - \frac{\nu}{2} (Q_{j+1}^n - Q_{j-1}^n) + \frac{\nu^2}{2} (Q_{j+1}^n - 2Q_j^n + Q_{j-1}^n).$$

Thus we have obtained the Lax-Wendroff scheme (1.30) via another route.

The procedure can be continued to higher orders. We can now ask what is the modified equation more accurately described by the Lax-Wendroff scheme. Repeating the same procedures as above we first write

$$(1.96) \quad s(x, t+k) = s(x, t) - \frac{\nu}{2} [s(x+h, t) - s(x-h, t)] + \frac{\nu^2}{2} [s(x+h, t) - 2s(x, t) + s(x-h, t)]$$

and then carry out Taylor series expansions around  $(x, t)$  to obtain

$$(1.97) \quad s + ks_t + \frac{k^2}{2}s_{tt} + \frac{k^3}{6}s_{ttt} + \dots = s - \frac{uk}{2h} \left[ 2hs_x + \frac{h^3}{3}s_{xxx} + \dots \right] +$$

$$(1.98) \quad \frac{u^2k^2}{2h^2} \left[ h^2s_{xx} + \frac{h^4}{12}s_{xxx} \right]$$

from where

$$(1.99) \quad s_t + us_x = -\frac{k}{2}(s_{tt} - u^2s_{xx}) - \frac{k^2}{6}s_{ttt} - \frac{uh^2}{6}s_{xxx} + O(k^3, h^3)$$

Note the appearance of the  $O(k)$  term. Had we carried out the Taylor expansion for the advection equation this term would have been proportional to  $q_{tt} - u^2q_{xx}$  which is zero according to (1.26). Here we cannot assume that  $s_{tt} - u^2s_{xx}$  is zero a priori. We must carry the term in the ensuing computations, expecting that it will give a higher order correction. Let us neglect the  $O(k^3, h^3)$  contributions and proceed with our technique of replacing higher order time derivatives with spatial derivatives using

$$(1.100) \quad s_{tt} = -us_{xt} - \frac{k}{2}(s_{ttt} - u^2s_{xxt})$$

$$(1.101) \quad s_{ttt} = -us_{xtt}.$$

Higher order terms have been dropped since they would lead to  $O(k^3, h^3)$  contributions in (1.99). Our intermediate result is

$$(1.102) \quad s_t + us_x = -\frac{k}{2} \left[ -us_{xt} - \frac{k}{2}(s_{ttt} - u^2s_{xxt}) - u^2s_{xx} \right] + \frac{k^2}{6}us_{xtt} - \frac{uh^2}{6}s_{xxx}$$

and we continue by eliminating mixed derivatives. In the above formula we wish to express  $s_{xtt}$  in terms of  $x$  derivatives to  $O(1)$

$$(1.103) \quad s_{xtt} = s_{ttx} = (s_t)_{tx} = (-us_x)_{tx} = -u(s_t)_{xx} = u^2s_{xxx}.$$

We also need to express  $s_{xt}$  in terms of  $x$  derivatives to  $O(k, h)$

$$(1.104) \quad s_{xt} = s_{tx} = -us_{xx} - \frac{k}{2}(s_{ttx} - u^2s_{xxx})$$

and  $s_{ttt}, s_{xxt}$  to  $O(1)$

$$(1.105) \quad s_{ttt} = -u^3s_{xxx}, \quad s_{xxt} = -us_{xxx}.$$

Replacing in (1.102) leads to

(1.106)

$$s_t + us_x = -\frac{k}{2} \left\{ -u \left[ -us_{xx} - \frac{k}{2} (s_{ttx} - u^2 s_{xxx}) \right] - u^2 s_{xx} - \frac{k}{2} (-u^3 s_{xxx} + u^3 s_{xxx}) \right\} + \frac{u}{6} (k^2 u^2 - h^2) s_{xxx}$$

which simplifies to

$$(1.107) \quad s_t + us_x = \frac{k^2}{4} (s_{ttx} - u^2 s_{xxx}) + \frac{u}{6} (k^2 u^2 - h^2) s_{xxx}$$

Since  $s_{ttx} = u^2 s_{xxx}$  to  $O(1)$  we obtain in final

$$(1.108) \quad s_t + us_x = -\frac{uh^2}{6} (1 - \nu^2) s_{xxx} .$$

The third order derivative now obtained shows that the Lax-Wendroff scheme introduces a *dispersive error* with different wave numbers traveling at different speeds. As expected, the dispersive error is proportional to  $h^2$  since the Lax-Wendroff scheme is second order. A scheme more accurate than Lax-Wendroff could be obtained by adding a correction term modeling the dispersive error. Since this involves a third-order derivative the stencil of the scheme would become wider by at least one unit thereby entailing more computational work.

**1.2. Non-linear scalar equations.** We have introduced a number of finite difference methods for the simple constant-velocity advection equation (1.1). Of course, there is hardly much need for a numerical method in order to solve (1.1). Rather we have used (1.1) as a model problem to study the properties of numerical schemes on a simple case. We now proceed to consider more complicated problems and investigate how the methods already derived apply to these problems.

A general first order, hyperbolic scalar equation is given by

$$(1.109) \quad q_t + u(x, t, q)q_x = \sigma(x, t, q)$$

where  $u$  may be interpreted as local advection velocity that depends in general upon  $x, t$  and  $q$ . In a wide range of problems equations of the form

$$(1.110) \quad q_t + f(q)_x = \sigma(x, t, q)$$

arise where  $f$  is known as the *flux function*. If  $f$  is differentiable we can write

$$(1.111) \quad q_t + f_q q_x = \sigma(x, t, q)$$

so  $f_q$  plays the role of the local advection velocity. Generally  $u, f$  depend on  $q$  so that the equations become non-linear in  $q$ . Equation (1.110) is said to be in *conservative form* as opposed to (1.109) which is said to be in *non-conservative form*. Generally we say that a PDE is in conservative form when it can be expressed as the space-time divergence of a vector field. For equation (1.110) the vector field would be  $(q, f(q))$  and the space-time divergence is  $\nabla_{(t,x)} \cdot = (\partial_t, \partial_x) \cdot$  so another way of writing (1.110) is

$$(1.112) \quad \nabla_{(t,x)} \cdot (q, f(q)) = 0 .$$

An initial value problem is defined by specifying a solution domain along  $x$  and an initial condition  $q_0(x)$ .

FIGURE 3. Characteristic curves for  $q_t + e^{x+t}q_x = -\beta q$ .

1.2.1. *Solution by characteristics.* We can solve (1.109) using the method of characteristics. We again ask whether there are any special curves  $\Gamma$  within the  $(x, t)$  plane on which (1.109) reduces to a simpler form. Along the curves specified by the differential system

$$(1.113) \quad \frac{dt}{ds} = 1, \quad \frac{dx}{ds} = u(x, t, q)$$

we do indeed obtain the simpler form

$$(1.114) \quad \frac{dq}{ds} = \sigma.$$

The essential difference with respect to the constant-velocity case is that the curves are no longer simple straight lines but depend on  $x, t$  and  $q$ . Let us consider some examples in order to see the complications involved.

Variable-velocity advection. Consider the equation

$$(1.115) \quad q_t + u(x, t)q_x = \sigma$$

which describes the advection of the unknown field variable  $q$  by an imposed velocity field  $u(x, t)$ . The velocity field is not influenced by  $q$  itself;  $q$  is said to be a *passive tracer*. The characteristic curves are given by the ODE

$$(1.116) \quad \frac{dx}{dt} = u(x, t).$$

Note that we are no longer guaranteed that the characteristics exist for all times as they did for the constant-velocity advection equation. This is the case only if  $u$  is uniformly Lipschitz.

EXAMPLE 7. *Consider the velocity field*

$$(1.117) \quad u(x, t) = x + t,$$

*the initial condition*

$$(1.118) \quad q_0(x) = \sin x,$$

*and the source term*

$$(1.119) \quad \sigma = -\beta q.$$

*The characteristic curves are*

$$(1.120) \quad x(t) = Ce^t - t - 1$$

*which are shown in Fig. (3). At  $t = 0$  the characteristic labeled by  $C$  passes through the  $x$  coordinate  $x_0 = C - 1$ . Along each characteristic the variable-velocity advection equation reduces to the ODE*

$$(1.121) \quad \frac{dq}{dt} = -\beta q$$

*which has the solution  $q(x, t) = Ae^{-\beta t}$ . We have to determine the constant  $A$  from the initial conditions. Through any given point  $(x, t)$  there passes the characteristic curve labeled by  $C = e^{-t}(x+t+1)$ . This particular characteristic curve will intersect*

FIGURE 4. Solution of  $q_t + (x+t)q_x = -\beta q$ ,  $q(x, t=0) = \sin x$  for  $\beta = 0.1$  at  $t = 0, 0.2, \dots, 1$ .

FIGURE 5. Crossing characteristics for inviscid Burgers equation with initial condition  $q_0(x) = \sin x$  (shown in thick line).

the  $x$ -axis at  $x_0 = C - 1$  and this is the position from which we must take the initial value for  $q$

$$(1.122) \quad q(x, t) = q_0(e^{-t}(x + t + 1) - 1)e^{-\beta t} = \sin(e^{-t}(x + t + 1) - 1)e^{-\beta t} .$$

We have found the solution to the PDE using the simpler expression of the PDE along the characteristics. The solution can be verified by direct substitution in (1.115) and is depicted in Fig. (4). The initial condition is spread out due to the spreading out of the characteristic curves and attenuated due to the source term  $\sigma$ .

Burgers equation. A model equation used extensively in the study of non-linear equations is

$$(1.123) \quad q_t + qq_x = 0$$

known as the *inviscid Burgers equation*. It is given in non-conservative form above. In conservative form it becomes

$$(1.124) \quad q_t + \left( \frac{q^2}{2} \right)_x = 0$$

so the flux function is

$$(1.125) \quad f(q) = q^2/2 .$$

The characteristic curves are given by

$$(1.126) \quad \frac{dx}{dt} = q(x, t)$$

and along a characteristic curve  $\Gamma$  equation (1.123) reduces to

$$(1.127) \quad \left( \frac{dq}{ds} \right)_\Gamma = 0 ,$$

i.e. there is no variation in  $q$  along the characteristic. This implies that the slope of each characteristic curve is constant and specified by the initial condition  $q(x, t=0) = q_0(x)$ .

The type of difficulties that arise for non-linear equations is immediately apparent from the consideration of simple initial conditions. Consider  $q_0(x) = \sin x$ . The characteristics are sketched in Fig. 5. The problem is that the characteristic curves cross one another. At such a crossing point it is not apparent what the correct value of  $q$  should be since different values are being transported along each of the crossing characteristics.

To get a better idea of what is happening it is useful to simplify the initial condition as much as possible. This leads to the so-called *Riemann problem*

$$(1.128) \quad q_0(x) = \begin{cases} q_l & x < 0 \\ q_r & x > 0 \end{cases}$$

Let us try to solve Burgers equation for this initial condition.

If  $q_l > q_r$  characteristics from  $x < 0$  will overtake those from  $x > 0$ . This will occur on some ray from the origin of equation  $x = st$ . To the left of this separating ray we will observe the value  $q_l$  while to the right we will observe the value  $q_r$ . The solution is therefore

$$(1.129) \quad q(x, t) = \begin{cases} q_l & x < st \\ q_r & x > st \end{cases}$$

The initial discontinuity propagates at a velocity  $s$ . The discontinuity is called a *shock* using the language of compressible gas dynamics and  $s$  is the *shock velocity*. The shock velocity can be determined by using the integrating Burgers equation over a domain having the shock as its diagonal  $[st_1, st_2] \times [t_1, t_2]$

$$(1.130) \quad \int_{st_1}^{st_2} \int_{t_1}^{t_2} [q_t + f(q)_x] dt dx = 0$$

from where

$$(1.131) \quad s = \frac{f(q_r) - f(q_l)}{q_r - q_l}.$$

If  $q_l < q_r$  two solutions are possible. We can again have the shock solution (1.129) but also the solution

$$(1.132) \quad q(x, t) = \begin{cases} q_l & x < q_l t \\ x/t & q_l t \leq x \leq q_r t \\ q_r & x > q_r t \end{cases}$$

called a *rarefaction solution*, again using terms from gas dynamics. This an even worse conundrum, not only can discontinuities arise which invalidate the differentiation operations but multiple solutions seem to be possible. Clearly something is wrong and a way to correct the model that led to equation (??) must be found. From the physical point of view certain effects have been neglected, namely the viscosity of the fluid and we might be led to studying the viscous Burgers equation

$$(1.133) \quad q_t + qq_x = \nu q_{xx}$$

as a remedy to the difficulties encountered. This can be done and leads to smooth solutions with very large gradients in the regions where shocks would have formed for the inviscid Burgers equation. These large gradients are difficult to resolve properly requiring very fine grids, much finer than needed elsewhere in the solution domain. So a way that enables us to still work with the inviscid equation is quite useful.

**1.2.2. Weak solutions.** The possibility of crossing characteristic curves is indicative with a breakdown of the modeling assumptions that led to a certain hyperbolic PDE. In this situation one must revisit the method by which a certain PDE is derived and consider the validity of all intermediate hypotheses used in the derivation. Burgers equation serves as a useful example. The PDE

$$(1.134) \quad q_t + f(q)_x = 0$$

with  $f = q^2/2$  was proposed as a model for fluid flow in which the quantity  $q$  is conserved but being advected by itself. The correct formulation of a conservation principle is through the integral statement

$$(1.135) \quad \int_{x_1}^{x_2} [q(x, t_2) - q(x, t_1)] dx = - \int_{t_1}^{t_2} [f(q(x_2, t)) - f(q(x_1, t))] dt$$

the one-dimensional expression of (1.10). In this form one can replace

$$(1.136) \quad q(x, t_2) - q(x, t_1) = \int_{t_1}^{t_2} \frac{\partial q}{\partial t}(x, t) dt$$

$$(1.137) \quad f(q(x_2, t)) - f(q(x_1, t)) = \int_{x_1}^{x_2} \frac{\partial f}{\partial x} dx$$

and obtain Burgers equation by going to the limits  $t_2 \rightarrow t_1$ ,  $x_2 \rightarrow x_1$  if the derivatives  $\partial q/\partial t$ ,  $\partial f/\partial x$  exist. However one cannot do this if  $q$  is discontinuous. In this case only the integral form (1.135) is valid.

Nonetheless it is typically much more convenient to work with differential equations instead of integral equations. Therefore it is useful to extend the meaning we associate to “ $q$  is a solution of a PDE” to cover the case where  $q$  might be discontinuous at a few points. This is done through the techniques of the theory of distributions by requiring that  $q$  satisfy a certain integral condition. Namely we consider the integral

$$(1.138) \quad I = \int_0^\infty \int_{-\infty}^{+\infty} \phi [q_t + f(q)_x] dx dt$$

with  $\phi$  a smooth function of finite support and impose  $I = 0$ . Typically we require that  $\phi$  be at least differentiable. We can integrate by parts to obtain

$$(1.139) \quad \int_0^\infty \int_{-\infty}^{+\infty} [\phi_t q + \phi_x f(q)] dx dt = - \int_{-\infty}^{+\infty} \phi(x, 0) q(x, 0) dx .$$

By this technique all differentiation operations on  $q$  have been removed. We say that  $q$  is a *weak solution* of (1.134) if (2.28) is satisfied for all  $\phi$  from some space of test functions such as  $\phi \in C^1(\mathbb{R} \times \mathbb{R})$ .

**1.2.3. Difficulties of finite difference methods for non-linear hyperbolic equations.** The possibility of shocks for non-linear hyperbolic equations should alert us to possible difficulties with the finite difference methods we have introduced for the linear advection equation. Since these are based upon Taylor series expansions of  $q(x, t)$  and  $q$  can be discontinuous, the expansions will break down and not be valid near the discontinuities. Nevertheless, we would expect the methods to be adequate in regions where  $q$  is smooth.

Let us see how we would apply the methods to a non-linear equation, taking Burgers equation as an example. One possibility is to interpret  $q$  as the local advection velocity  $u$ . The upwind method for

$$(1.140) \quad q_t + qq_x = 0$$

then becomes

$$(1.141) \quad Q_j^{n+1} = Q_j^n - \begin{cases} Q_j^n (Q_j^n - Q_{j-1}^n) & \text{if } Q_j^n \geq 0 \\ Q_j^n (Q_{j+1}^n - Q_j^n) & \text{if } Q_j^n < 0 \end{cases}$$



and the Lax-Wendroff methods reads

$$(1.142) \quad Q_j^{n+1} = Q_j^n - \frac{Q_j^n k}{2h} (Q_{j+1}^n - Q_{j-1}^n) + \frac{(Q_j^n)^2 k^2}{2h^2} (Q_{j+1}^n - 2Q_j^n + Q_{j-1}^n) .$$

Applying this for a Riemann problem leads to a numerical solution similar to the exact shock solution but with oscillations near the shock (Fig. 1.2.3). There is also a smearing of the shock, instead of sharp discontinuity we have a smoothing of  $q$  in the vicinity of the shock. Far from the shock the numerical solution is quite good however. This therefore leads to the search for so-called *high-resolution algorithms* that are able to preserve a high order of accuracy away from discontinuities and also sharply capture discontinuities.

## 2. Systems of hyperbolic equations

### 2.1. Linear systems.

2.1.1. *Classification of linear systems.* Consider now that we are interested in the simultaneous time evolution of a number of quantities

$$(2.1) \quad q = \begin{bmatrix} q_1 & q_2 & \dots & q_m \end{bmatrix}^T$$

that satisfy

$$(2.2) \quad q_t + \mathbf{A}q_x = 0$$

with  $\mathbf{A}$  a constant  $m \times m$  matrix of real numbers. Such a system is said to be hyperbolic if the eigenvectors of  $\mathbf{A}$  form a basis for real  $m$ -vectors.

EXAMPLE 8. *The second order wave equation is given in canonical form as*

$$(2.3) \quad \phi_{tt} - c^2 \phi_{xx} = 0 .$$

*It can be reduced to a system of two first-order equations by introducing*

$$(2.4) \quad u = \phi_t, \quad v = \phi_x$$

*We have*

$$(2.5) \quad u_t - c^2 v_x = 0$$

*and since  $\phi_{xt} = \phi_{tx}$*

$$(2.6) \quad v_t - u_x = 0$$

*In vector form we obtain*

$$(2.7) \quad \frac{\partial}{\partial t} \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} 0 & -c^2 \\ -1 & 0 \end{bmatrix} \frac{\partial}{\partial x} \begin{bmatrix} u \\ v \end{bmatrix} = 0$$

*or*

$$(2.8) \quad q_t + \mathbf{A}q_x = 0$$

*with*

$$(2.9) \quad q = \begin{bmatrix} u \\ v \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 0 & -c^2 \\ -1 & 0 \end{bmatrix}$$

*The eigenvalues of  $\mathbf{A}$  are  $\lambda_{1,2} = \pm c$  and the eigenvectors are*

$$(2.10) \quad r_1 = \begin{bmatrix} c \\ 1 \end{bmatrix}, \quad r_2 = \begin{bmatrix} -c \\ 1 \end{bmatrix}$$

The eigenvectors are independent for  $c \neq 0$  and therefore they form a basis for the space of real 2-vectors. The system (2.8) is hyperbolic.

EXAMPLE 9. Applying the same procedure to the Laplace equation

$$(2.11) \quad \phi_{tt} + \phi_{xx} = 0$$

leads to the matrix

$$(2.12) \quad \mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

whose eigenvalues are  $\lambda_{1,2} = \pm i$  and eigenvectors are

$$(2.13) \quad \begin{bmatrix} i \\ 1 \end{bmatrix}, \begin{bmatrix} -i \\ 1 \end{bmatrix}$$

These have complex values and are not a basis for real 2-vectors. The system (2.12) is not hyperbolic, it is elliptic.

### 2.1.2. Solution by method of characteristics and reduction to diagonal form.

For hyperbolic systems we can apply a procedure similar to that used for systems of ODE's. We can write

$$(2.14) \quad \mathbf{A} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1}$$

with

$$(2.15) \quad \mathbf{\Lambda} = \text{diag} \{ \lambda_1, \lambda_2, \dots, \lambda_m \}$$

$$(2.16) \quad \mathbf{T} = [r_1 \ r_2 \ \cdots \ r_m]$$

$$(2.17) \quad \mathbf{A}r_j = \lambda_j r_j, \ j = 1, 2, \dots, m.$$

and write

$$(2.18) \quad q_t + \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1}q_x = 0.$$

Introducing the notation

$$(2.19) \quad w = \mathbf{T}^{-1}q$$

we obtain

$$(2.20) \quad w_t + \mathbf{\Lambda}w_x = 0.$$

Since  $\mathbf{\Lambda}$  is a diagonal matrix, the equations of the original system have been decoupled and we can write the scalar  $j^{\text{th}}$  component equation

$$(2.21) \quad w_t^{(j)} + \lambda_j w_x^{(j)} = 0,$$

for  $j = 1, 2, \dots, m$ . These are now simple constant-velocity advection equations for which we know the solution

$$(2.22) \quad w^{(j)}(x, t) = w_0^{(j)}(x - \lambda_j t)$$

with  $w_0^{(j)}$  given by the initial conditions on  $q$

$$(2.23) \quad w_0 = \mathbf{T}^{-1}q_0.$$

The value of each individual component of  $w^{(j)}$  is constant along the family of characteristics  $x - \lambda_j t = C_j$ . Therefore  $w$  are known as the *conservative variables*.

From a knowledge of the conservative variable solution we can recover the solution for the original variables

$$(2.24) \quad q = \mathbf{T}w .$$

2.1.3. *Finite difference methods.* The finite difference methods derived for the constant-velocity advection equation can be applied formally to hyperbolic systems also. For example, the Lax-Wendroff scheme is

$$(2.25) \quad Q_j^{n+1} = Q_j^n - \frac{k}{2h} \mathbf{A} (Q_{j+1}^n - Q_{j-1}^n) + \frac{k^2}{2h^2} \mathbf{A}^2 (Q_{j+1}^n - 2Q_j^n + Q_{j-1}^n)$$

There are some new features though due to the fact that there is no longer just a single “advection” or characteristic velocity. Let us try to apply the upwind method to the system (2.8). It is not apparent what the upwind direction should be for  $q$ . We can ascertain this for the conservative variables though. We have  $\lambda_1 = -c$ ,  $\lambda_2 = c$ ,

$$(2.26) \quad r_1 = \begin{bmatrix} c \\ 1 \end{bmatrix}, \quad r_2 = \begin{bmatrix} -c \\ 1 \end{bmatrix}$$

$$(2.27) \quad T = \begin{bmatrix} c & -c \\ 1 & 1 \end{bmatrix}, \quad T^{-1} = \begin{bmatrix} \frac{1}{2c} & \frac{1}{2} \\ -\frac{1}{2c} & \frac{1}{2} \end{bmatrix}, \quad \Lambda = \begin{bmatrix} -c & 0 \\ 0 & c \end{bmatrix}$$

and the conservative variable system is

$$(2.28) \quad \frac{\partial}{\partial t} \begin{bmatrix} w^{(1)} \\ w^{(2)} \end{bmatrix} + \begin{bmatrix} -c & 0 \\ 0 & c \end{bmatrix} \frac{\partial}{\partial x} \begin{bmatrix} w^{(1)} \\ w^{(2)} \end{bmatrix} = 0 .$$

This system can be discretized in an upwind manner and we obtain the scheme

$$(2.29) \quad \left(W^{(1)}\right)_j^{n+1} = \left(W^{(1)}\right)_j^n + \frac{ck}{h} \left[ \left(W^{(1)}\right)_{j+1}^n - \left(W^{(1)}\right)_j^n \right]$$

$$(2.30) \quad \left(W^{(2)}\right)_j^{n+1} = \left(W^{(2)}\right)_j^n - \frac{ck}{h} \left[ \left(W^{(2)}\right)_j^n - \left(W^{(2)}\right)_{j-1}^n \right]$$

In matrix form this reads

$$(2.31) \quad W_j^{n+1} = (1 - \nu) W_j^n + \mathbf{C} W_{j-1}^n + \mathbf{D} W_{j+1}^n$$

with

$$(2.32) \quad \nu = \frac{ck}{h}, \quad \mathbf{C} = \begin{bmatrix} 0 & 0 \\ 0 & \nu \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} \nu & 0 \\ 0 & 0 \end{bmatrix}$$

Multiplying by  $T$  leads to

$$(2.33) \quad Q_j^{n+1} = (1 - \nu) Q_j^n + \mathbf{TCT}^{-1} Q_{j-1}^n + \mathbf{TDT}^{-1} Q_{j+1}^n$$

$$(2.34) \quad \mathbf{TCT}^{-1} = \frac{\nu}{2} \begin{bmatrix} 1 & -c \\ -\frac{1}{c} & 1 \end{bmatrix}, \quad \mathbf{TDT}^{-1} = \frac{\nu}{2} \begin{bmatrix} 1 & c \\ \frac{1}{c} & 1 \end{bmatrix}$$

This is the upwind scheme for the system (2.8).

## 2.2. Non-linear systems.

2.2.1. *Classification.* Non-linear systems are written as

$$(2.35) \quad q_t + f(q)_x = 0$$

with

$$(2.36) \quad q = [q_1 \quad q_2 \quad \dots \quad q_m]^T, \quad f = [f_1 \quad f_2 \quad \dots \quad f_m]^T$$

The classification of non-linear systems is made in accordance with the properties of the Jacobian of  $f$  with respect to  $q$

$$(2.37) \quad f_q = \begin{bmatrix} \frac{\partial f_1}{\partial q_1} & \frac{\partial f_1}{\partial q_2} & \dots & \frac{\partial f_1}{\partial q_m} \\ \frac{\partial f_2}{\partial q_1} & \frac{\partial f_2}{\partial q_2} & \dots & \frac{\partial f_2}{\partial q_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial q_1} & \frac{\partial f_m}{\partial q_2} & \dots & \frac{\partial f_m}{\partial q_m} \end{bmatrix}$$

If the eigenvectors of  $f_q$  form a basis for  $q$ -vectors the system is said to be hyperbolic, otherwise it is elliptic or parabolic. Note that in this case the eigenvectors typically depend on the variables  $q$  themselves so that the same system of equations may be hyperbolic in some regions and elliptic in others. The classification of PDE's as hyperbolic, parabolic and elliptic may be more familiar from the classification of second order equations. Let us show the equivalence of the two usages.

The canonical elliptic second order PDE is the Laplace equation

$$(2.38) \quad \phi_{tt} + \phi_{xx} = 0.$$

We reduce it to a system of first-order PDE's by introducing  $u = \phi_t$ ,  $v = \phi_x$ . The Laplace equation states  $u_t + v_x = 0$  and we also have  $u_x = v_t$  by the equality of mixed derivatives. These two relations can be written in matrix form as

$$(2.39) \quad q_t + \mathbf{A}q_x = 0$$

$$(2.40) \quad q = \begin{bmatrix} u \\ v \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} = 0$$

The matrix  $\mathbf{A}$  has the eigenvalues  $\lambda_{1,2} = \pm i$  and eigenvectors  $r_{1,2} = [\pm i \quad 1]^T$ . The eigenvectors  $r_{1,2}$  do not form a basis for two-component real vectors such as  $q$  so the system is classified as elliptic in accord with the second-order Laplace equation's classification.

The canonic hyperbolic second order PDE is the wave equation

$$(2.41) \quad \phi_{tt} - \phi_{xx} = 0.$$

Following the same procedure we arrive at the study of the eigensystem of

$$(2.42) \quad \mathbf{B} = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix}$$

which is given by  $\lambda_{1,2} = \pm 1$ ,  $r_{1,2} = [\pm 1 \quad 1]^T$ . The eigenvectors now do form a basis for two-component real vectors and the system is classified as hyperbolic as expected from the wave equation.

Finally, the typical parabolic equation is

$$(2.43) \quad \phi_x = \phi_{tt}$$

for which we denote  $u = \phi_t$  to obtain

$$(2.44) \quad q_t + Cq_x = \sigma$$

with

$$(2.45) \quad q = \begin{bmatrix} \phi \\ u \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \quad \sigma = \begin{bmatrix} u \\ 0 \end{bmatrix}.$$

The eigensystem of  $\mathbf{C}$  is  $\lambda_{1,2} = 0$ ,  $r_1 = \begin{bmatrix} 0 & 1 \end{bmatrix}$ ,  $r_2 = \begin{bmatrix} 0 & 0 \end{bmatrix}$  which does not form a basis for two-component real vectors. Note that in this case the rank of  $\mathbf{C}$  is less than the dimension of the system; this is characteristic of parabolic equations.

2.2.2. *Solution by characteristics.* Let  $\mathbf{A}$  be the Jacobian matrix for a non-linear hyperbolic system

$$(2.46) \quad q_t + \mathbf{A}(q)q_x = 0,$$

with  $q$  a vector with  $m$  components. By the definition of a hyperbolic system we know that  $\mathbf{A}$  can be represented as

$$(2.47) \quad \mathbf{A} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1}.$$

The difference with respect to the linear system case is that the matrices  $\mathbf{T}$ ,  $\mathbf{\Lambda}$  are no longer constant but depend on  $q$  and hence on  $(x, t)$ . Nevertheless, we can follow the same procedure of reduction to characteristic form locally for some neighborhood of a point  $(x_0, t_0)$  where  $q(x, t) = q_0$ . We can write

$$(2.48) \quad q(x, t) = q_0 + \tilde{q}(x, t)$$

where  $\tilde{q}$  is the perturbation from the value  $q_0$ . System (2.46) can now be written

$$(2.49) \quad \tilde{q}_t + \mathbf{A}_0\tilde{q}_x = 0,$$

from where we obtain

$$(2.50) \quad \tilde{w}_t + \mathbf{\Lambda}_0\tilde{w}_x = 0$$

with the perturbation characteristic variables given by

$$(2.51) \quad \tilde{w} = \mathbf{T}^{-1}\tilde{q}_0.$$

The characteristic system (2.50) leads to the ODE's

$$(2.52) \quad \frac{d\tilde{w}^{(i)}}{ds_i} = 0, i = 1, \dots, m.$$

where  $d/ds_i$  indicates the derivative along the  $i^{th}$  characteristic direction whose slope is given by the  $\lambda_{0i}$  eigenvalue of  $\mathbf{A}_0$

$$(2.53) \quad \frac{d}{ds_i} = \frac{\partial}{\partial t} + \lambda_{0i} \frac{\partial}{\partial x}.$$

A solution to (2.46) can be found by locally solving the ODE's (2.52). This is the *method of characteristics* for non-linear hyperbolic systems.