

A MODIFICATION OF ROE'S SCHEME FOR ENTROPY  
SATISFYING SOLUTIONS OF SCALAR NON-LINEAR  
CONSERVATION LAWS

P. K. SWEBY

---

NUMERICAL ANALYSIS REPORT 6/82

A MODIFICATION OF ROE'S SCHEME FOR ENTROPY SATISFYING  
SOLUTIONS OF SCALAR NON-LINEAR CONSERVATION LAWS

P. K. Sweby

Introduction

In Sweby & Baines [2] convergence of Roe's scheme for a scalar non-linear conservation law is proved. However Roe's method suffers from an inability to deal satisfactorily with expansion fans.

In the present report a special procedure is devised for dealing with expansion fans and this is described in detail below. The resulting scheme is then put in a more general framework, which reduces to the scheme in Sweby & Baines (loc cit) away from expansions.

Consider the scalar non-linear equation

$$u_t + f(u)_x = 0 \quad -\infty < x < \infty, \quad 0 < t \leq T \quad (1a)$$

with initial data

$$u(x,0) = u_0(x), \quad (1b)$$

which is assumed to be (a) uniformly bounded in the  $L^\infty$  norm, and (b) of uniformly bounded variation, namely, for all real  $\delta$  and  $R > 0$

$$\int_{|x| < R} |u_0(x + \delta) - u_0(x)| dx \leq C(R) |\delta|. \quad (2)$$

We divide the  $x$  axis uniformly into cells  $(x_{k-1}, x_k)$  with

$$x_k = x_{k-1} + h \quad (3)$$

and take a time step  $\Delta t$ , defining the constant  $q$  by

$$q = \frac{\Delta t}{h} \quad (4)$$

We now define the piecewise constant projection  $u_h(x,t)$  of  $u(x,t)$  onto the computational grid by

$$u_h(x,t) = u_k^n \quad (5)$$

for  $(x,t) \in (x_{k-1}, x_{k+\frac{1}{2}}) \times (n-\frac{1}{2}qh, n+\frac{1}{2}qh)$ , where  $u_k^n$  are nodal values generated by a finite difference scheme. The initial data is projected by the restriction

$$u_k^0 = \frac{1}{h} \int_{x_{k-\frac{1}{2}}}^{x_{k+\frac{1}{2}}} u_0(x) dx \quad (6)$$

We now develop the scheme for generating the nodal values  $u_k^n$ . For brevity we adopt the notation

$$\left. \begin{aligned} u^k &= u_k^{n+1} \\ u_k &= u_k^n \end{aligned} \right\} \quad (7)$$

whenever there is no danger of confusion.

Due to the piecewise constant nature of  $u_h(x,t)$  we can, for sufficiently small  $\Delta t$ , treat the advancement of the solution in each cell from  $t = \tau$  to  $t = \tau + \Delta t$  as the solution of the local scalar Riemann problem, i.e.

$$\left. \begin{aligned} u_t + f(u)_x &= 0 \quad (x,t) \in (x_{k-1}, x_k) \times (\tau, \tau + \Delta t) \\ u(x,\tau) &= \begin{cases} u_{k-1} & x < x_{k-\frac{1}{2}} \\ u_k & x > x_{k-\frac{1}{2}} \end{cases} \end{aligned} \right\} \quad (8)$$

with  $\Delta t$  chosen such that the characteristics do not intersect.

Let  $a(x,t) = a(u(x,t)) = \frac{\partial f}{\partial u}$  be the wave speed and, to ease notation, make the following change of variables:

$$\left. \begin{aligned} X &= x - x_{k-\frac{1}{2}} \\ T &= t - \tau \end{aligned} \right\} \quad (9)$$

The problem (8) has solution (in terms of the new variables) [1]

$$u(X,T) = \begin{cases} u_{k-1} & X/T < a^L \\ u_k & X/T > a^R \end{cases} \quad (10)$$

where  $a^L, a^R$  are the left and right wave speeds respectively, i.e.

$$\begin{aligned} a^L &= a(x_{k-1}, \tau) \\ a^R &= a(x_k, \tau) \\ \text{and} \quad & -\infty < a^L \leq a^R < \infty. \end{aligned}$$

If the discontinuity is a shock then  $a^L = a^R = a(x_{k-\frac{1}{2}}, \tau)$  with

$$a(x_{k-\frac{1}{2}}, \tau) = \frac{f(u_k) - f(u_{k-1})}{u_k - u_{k-1}} \quad (11)$$

being the shock speed given by the Rankine-Hugoniot relation (see Fig. 1).

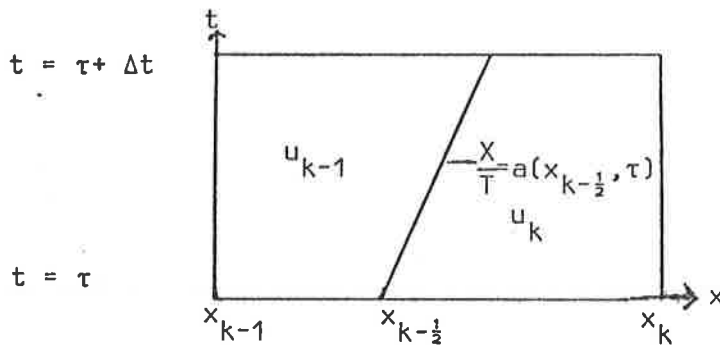


Fig. 1

If, however,  $a^L < a(x_{k-\frac{1}{2}}, \tau) < a^R$  then the correct physical solution is an expansion wave [1], i.e.

$$u(X, T) = \begin{cases} u_{k-1} & X/T \leq a^L \\ u_{k-1} + \theta(u_k - u_{k-1}) & a^L < X/T < a^R \\ u_k & X/T \geq a^R \end{cases} \quad (12)$$

where

$$\theta = \frac{(X/T - a^L)}{(a^R - a^L)} \quad (13)$$

(see Fig. 2).

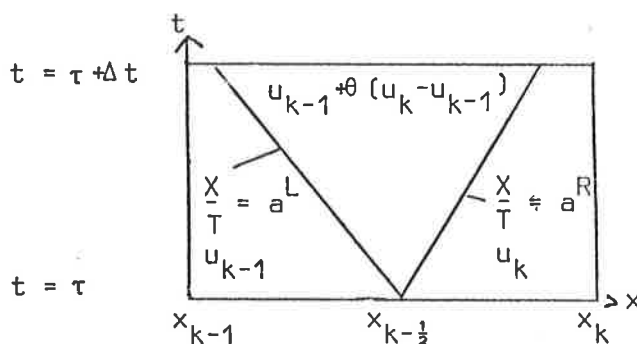


Fig. 2

We now project the piecewise linear solution (12) onto the space of piecewise constant functions.

## 2. Approximation of the Expansion

The linear solution (12) is approximated by the piecewise constant solution

$$u(X,T) = \begin{cases} u_{k-1} & X/T \leq a^L \\ u_{k-\frac{1}{2}} & a^L < X/T < a^R \\ u_k & X/T \geq a^R \end{cases} \quad (14)$$

To obtain  $u_{k-\frac{1}{2}}$  we integrate (1a) between  $x_{k-1}$  and  $x_k$  at time  $t$ ;

$$\text{i.e. } \frac{d}{dt} \int_{x_{k-1}}^{x_k} u(x,t) dx = -\{f(u(x_k,t)) - f(u(x_{k-1},t))\}. \quad (15)$$

Changing variables to  $X,T$  gives

$$\frac{d}{dT} \int_{-\frac{h}{2}}^{\frac{h}{2}} u(X,T) dX = -\{f(u_k) - f(u_{k-1})\} \quad (16)$$

Since  $u$  is piecewise constant we can evaluate the integral in (16) as

$$\begin{aligned} \int_{-\frac{h}{2}}^{\frac{h}{2}} u(X,T) dX &= \int_{-\frac{h}{2}}^{a^L T} u_{k-1} dX + \int_{a^L T}^{a^R T} u_{k-\frac{1}{2}} dX + \int_{a^R T}^{\frac{h}{2}} u_k dX \\ &= (a^L T + \frac{h}{2}) u_{k-1} + (a^R T - a^L T) u_{k-\frac{1}{2}} + (\frac{h}{2} - a^R T) u_k. \end{aligned} \quad (17)$$

So carrying out the differentiation in (16) we obtain

$$\begin{aligned} a^L (u_{k-\frac{1}{2}} - u_{k-1}) + a^R (u_k - u_{k-\frac{1}{2}}) &= f(u_k) - f(u_{k-1}) \\ &= a \{u_k - u_{k-1}\}, \end{aligned} \quad (18)$$

with  $a$  as defined in (11). Thus

$$(a^R - a^L) u_{k-\frac{1}{2}} = (a - a^L) u_{k-1} + (a^R - a) u_k. \quad (19)$$

Next we define numerical approximations  $a_{k-\frac{1}{2}}^L, a_{k-\frac{1}{2}}^R, a_{k-\frac{1}{2}}$  for  $a^L, a^R$  and  $a$ , respectively, by

$$a_{k-\frac{1}{2}} = \left. \begin{cases} \frac{f_k - f_{k-1}}{u_k - u_{k-1}} & u_k \neq u_{k-1} \\ \text{approximation for } \frac{\partial f}{\partial u} \Big|_{x_{k-\frac{1}{2}}} & u_k = u_{k-1} \end{cases} \right\} \quad (20)$$

$$a_{k-\frac{1}{2}}^L = \min \{a_{k-1}, a_{k-\frac{1}{2}}\} \quad (21)$$

$$a_{k-\frac{1}{2}}^R = \max \{a_{k-\frac{1}{2}}, a_k\},$$

where  $a_k, a_{k-1}$  are approximations for  $\frac{\partial f}{\partial u}(u_k), \frac{\partial f}{\partial u}(u_{k-1})$ , respectively. Note that the minimum and maximum in (21) ensure that  $a_{k-\frac{1}{2}}^L \leq a_{k-\frac{1}{2}} \leq a_{k-\frac{1}{2}}^R$ . We now have the situation as shown in Figure 3, and can devise a difference scheme as follows.

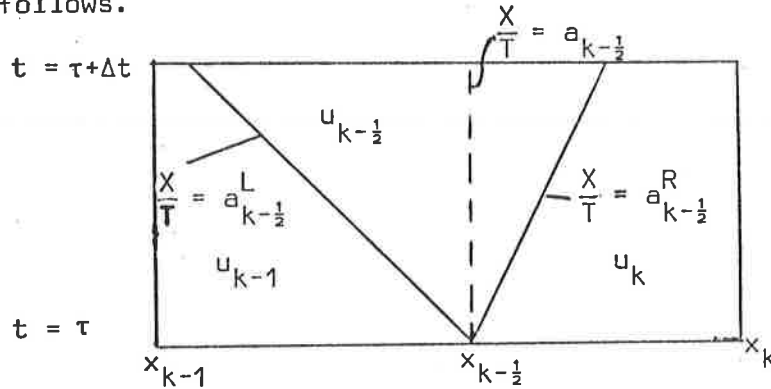


Fig. 3

### 3. The Finite Difference Scheme

We define left- and right-moving increments  $\phi_{k-\frac{1}{2}}^L$  and  $\phi_{k-\frac{1}{2}}^R$  as follows;

$$\phi_{k-\frac{1}{2}}^L = \frac{\Delta}{h} \int_{x_{k-1}}^{x_{k-\frac{1}{2}}} u \, dx, \quad \phi_{k-\frac{1}{2}}^R = \frac{\Delta}{h} \int_{x_{k-\frac{1}{2}}}^{x_k} u \, dx, \quad (22)$$

with a total fluctuation for the cell of

$$\phi_{k-\frac{1}{2}} = \frac{\Delta}{h} \int_{x_{k-1}}^{x_k} u \, dx = -\frac{1}{h} \int_{\tau}^{\tau+\Delta t} \int_{x_{k-1}}^{x_k} f_x \, dx \, dt$$

(arising from integrating (1) over the cell between  $t = \tau$  on  $t = \tau + \Delta t$ )

$$= -\frac{1}{h} \int_{\tau}^{\tau+\Delta t} (f(u_k) - f(u_{k-1})) \, dt$$

$$= -\frac{\Delta t}{h} a_{k-\frac{1}{2}} (u_k - u_{k-1}) \quad (23)$$

The assumption  $\frac{\Delta t}{h} |f'(u)| \leq \frac{1}{2} \quad (24a)$

and therefore  $\frac{\Delta t}{h} |a| \leq \frac{1}{2} \quad (24b)$

for  $a = a_{k-\frac{1}{2}}^L, a_{k-\frac{1}{2}}^R, a_{k-\frac{1}{2}}$  is made throughout, ensuring that discontinuities from adjacent cells do not interact.

We now perform the integration for  $\phi_{k-\frac{1}{2}}^L, \phi_{k-\frac{1}{2}}^R$  in (22) in a similar manner as for (16).

$$\begin{aligned} \phi_{k-\frac{1}{2}}^L &= \frac{\Delta}{h} \int_{x_{k-1}}^{x_{k-\frac{1}{2}}} u \, dx = \frac{1}{h} \left\{ \int_{x_{k-1}}^{x_{k-\frac{1}{2}}} u(x, t+\Delta t) \, dx - \int_{x_{k-1}}^{x_{k-\frac{1}{2}}} u(x, t) \, dx \right\} \\ &= \frac{1}{h} \left\{ u_k (-a_{k-\frac{1}{2}}^R)^+ \Delta t + u_{k-\frac{1}{2}} [(-a_{k-\frac{1}{2}}^L)^+ (-a_{k-\frac{1}{2}}^R)^+] \Delta t + u_{k-1} \left( \frac{h}{2} - (-a_{k-\frac{1}{2}}^L)^+ \Delta t \right) \right\} \frac{1}{h} \left\{ u_{k-1} \frac{h}{2} \right\} \\ &= q \left\{ (-a_{k-\frac{1}{2}}^R)^+ (u_k - u_{k-\frac{1}{2}}) + (-a_{k-\frac{1}{2}}^L)^+ (u_{k-\frac{1}{2}} - u_{k-1}) \right\}, \end{aligned} \quad (25)$$

where  $b^+ = \frac{1}{2}(b + |b|)$  is the positive part of  $b$ .

Defining

$$v_{k-\frac{1}{2}} = q a_{k-\frac{1}{2}} \quad (26)$$

we may now write (25) as

$$\phi_{k-\frac{1}{2}}^L = (-v_{k-\frac{1}{2}}^R)^+ (u_k - u_{k-\frac{1}{2}}) + (-v_{k-\frac{1}{2}}^L)^+ (u_{k-\frac{1}{2}} - u_{k-1}). \quad (27a)$$

Similarly

$$\phi_{k-\frac{1}{2}}^R = (v_{k-\frac{1}{2}}^R)^+ (u_k - u_{k-\frac{1}{2}}) - (v_{k-\frac{1}{2}}^L)^+ (u_{k-\frac{1}{2}} - u_{k-1}). \quad (27b)$$

We observe that

$$\begin{aligned} \phi_{k-\frac{1}{2}}^L + \phi_{k-\frac{1}{2}}^R &= -v_{k-\frac{1}{2}}^R (u_k - u_{k-\frac{1}{2}}) - v_{k-\frac{1}{2}}^L (u_{k-\frac{1}{2}} - u_{k-1}) \\ &= -v_{k-\frac{1}{2}} (u_k - u_{k-1}) \quad \text{from (18)} \\ &= \phi_{k-\frac{1}{2}}, \end{aligned} \quad (28)$$

as would be expected due to conservation.

Note that if  $v_{k-\frac{1}{2}}^R = v_{k-\frac{1}{2}}^L = v_{k-\frac{1}{2}}$  then

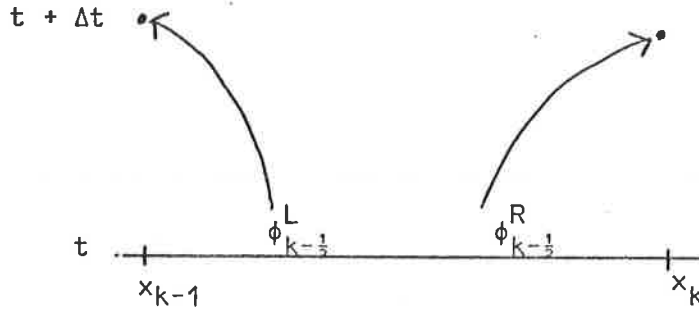
$$\left. \begin{aligned} \phi_{k-\frac{1}{2}}^L &= \phi_{k-\frac{1}{2}}, \quad \phi_{k-\frac{1}{2}}^R = 0 & \text{if } v_{k-\frac{1}{2}} < 0 \\ \phi_{k-\frac{1}{2}}^L &= 0, \quad \phi_{k-\frac{1}{2}}^R = \phi_{k-\frac{1}{2}} & \text{if } v_{k-\frac{1}{2}} > 0 \end{aligned} \right\} \quad (29)$$

(i) First Order

Incrementing  $u_{k-1}$  by  $\phi_{k-\frac{1}{2}}^L$  and  $u_k$  by  $\phi_{k-\frac{1}{2}}^R$  over a time step  $\Delta t$  yields the first order scheme

$$u^k = u_k + \phi_{k+\frac{1}{2}}^L + \phi_{k-\frac{1}{2}}^R, \quad (30)$$

which is shown graphically in Figure 4.



The First Order Scheme

Fig. 4

We note from (29) that this scheme is equivalent to Roe's first order scheme [3] everywhere except at expansions.

(ii) Second Order

We now introduce an antidiffusion stage by transferring some of the allocated fluctuation [2]. The result is a second order scheme almost everywhere (see below).

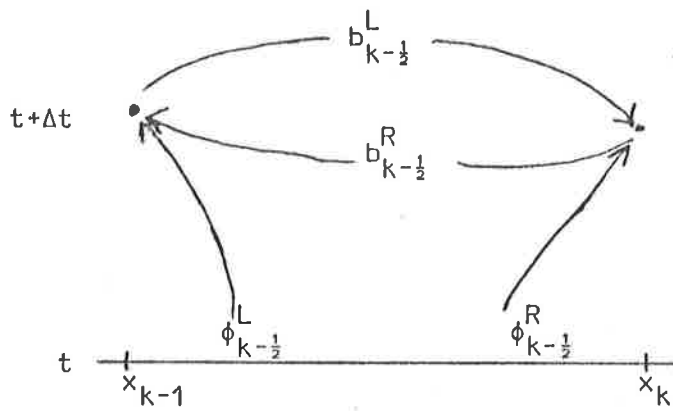
Let  $\alpha_{k-\frac{1}{2}} = \frac{1}{2}(1 - |v_{k-\frac{1}{2}}|)$  (31)

$$\left. \begin{aligned} b_{k-\frac{1}{2}}^L &= \text{minmod} \{ \alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^L, \alpha_{k+\frac{1}{2}} \phi_{k+\frac{1}{2}}^L \} \\ b_{k-\frac{1}{2}}^R &= \text{minmod} \{ \alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^L, \alpha_{k-\frac{3}{2}} \phi_{k-\frac{3}{2}}^R \} \end{aligned} \right\} \quad (32)$$

where  $\text{minmod} \{ \cdot, \cdot \}$  selects the argument with minimum absolute value. (In the case of arguments with equal modulus but opposite signs the first argument is chosen).

At the antidiffusion stage the quantity  $b_{k-\frac{1}{2}}^L$  is transferred left to right in the cell whilst  $b_{k-\frac{1}{2}}^R$  is transferred right to left (see Fig. 5). Again except at expansions either  $b_{k-\frac{1}{2}}^L$  or  $b_{k-\frac{1}{2}}^R$  will be zero as a consequence of (29) and (32).





The Second Order Scheme

Fig. 5

This scheme can be written

$$u^k = u_k + \xi_{k+1/2} \phi_{k+1/2}^L + \zeta_{k-1/2} \phi_{k-1/2}^R, \quad (33)$$

where

$$\xi_{k+1/2} = \begin{cases} 1 + \frac{(b_{k-1/2}^L - b_{k+1/2}^L)}{\phi_{k+1/2}^L} & \phi_{k+1/2}^L \neq 0 \\ 1 & \phi_{k+1/2}^L = 0 \end{cases} \quad \zeta_{k-1/2} = \begin{cases} 1 - \frac{(b_{k-1/2}^R - b_{k+1/2}^R)}{\phi_{k-1/2}^R} & \phi_{k-1/2}^R \neq 0 \\ 1 & \phi_{k-1/2}^R = 0 \end{cases} \quad (34)$$

and is second order accurate except at discontinuities of  $u$  (see below).

The scheme (33), (34) is equivalent to Roe's second order scheme [3], except at expansions, as studied by Sweby & Baines [2]. If the minmod operators in (32) always choose the first argument then the scheme is the Lax-Wendroff scheme [4]. If the second argument is always chosen it is the Warming and Beam fully upwind scheme [5].

### 3. Proof of Second Order Accuracy

We now demonstrate that the scheme (33), (34) is second order accurate away from discontinuities of the solution by comparing it with the Lax-Wendroff scheme.

A necessary and sufficient condition for a scheme formulated as in (33) to be second order accurate is that

$$(\xi_{k+\frac{1}{2}} - \xi_{k+\frac{1}{2}}^{LW}) \phi_{k+\frac{1}{2}}^L + (\zeta_{k-\frac{1}{2}} - \zeta_{k-\frac{1}{2}}^{LW}) \phi_{k-\frac{1}{2}}^R = O(h^2) \quad (35)$$

where  $\xi_{k+\frac{1}{2}}^{LW}$ ,  $\zeta_{k-\frac{1}{2}}^{LW}$  are the appropriate coefficients for the Lax-Wendroff scheme.

Except for the case when  $v_{k-\frac{1}{2}}^L < 0 < v_{k-\frac{1}{2}}^R$ , the Lax-Wendroff scheme is achieved by choosing the first argument in the minmod operators in (32) so it is only necessary to show that the arguments of the minmod operators differ by  $O(h^2)$ . That is we seek to show that

$$\alpha_{k+\frac{1}{2}} \phi_{k+\frac{1}{2}}^L - \alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^L = O(h^2) \quad (36)$$

and

$$\alpha_{k-3/2} \phi_{k-3/2}^R - \alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^R = O(h^2) \quad (37)$$

We shall verify only (36) here, but (37) may be verified in a similar manner. Assuming that  $u$  is smooth enough to be expanded in a Taylor series in  $x$  (and hence that  $f$  can be expanded as a function of  $x$ ),

$$\begin{aligned} \alpha_{k+\frac{1}{2}} \phi_{k+\frac{1}{2}}^L &= -\frac{1}{2}(1 - |v_{k+\frac{1}{2}}|) q (f_{k+1} - f_k) \\ &= -\frac{1}{2}(1 - |v_{k-\frac{1}{2}}| + O(h)) q (h \frac{\partial f}{\partial x}|_k + O(h^2)) \\ &= -(\alpha_{k-\frac{1}{2}} + O(h)) q (h \frac{\partial f}{\partial x}|_k + O(h^2)) \\ &= -\alpha_{k-\frac{1}{2}} q h \frac{\partial f}{\partial x}|_k + O(h^2) \end{aligned}$$

But also  $\alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^L = -\alpha_{k-\frac{1}{2}} q h \frac{\partial f}{\partial x}|_k + O(h^2)$

Hence  $\alpha_{k+\frac{1}{2}} \phi_{k+\frac{1}{2}}^L - \alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^L = O(h^2)$  verifying (36).

At discontinuities of  $u$  the Taylor expansion is invalid and hence we cannot claim second order accuracy at such points. Note, however, that the scheme is still second order at points where the minmod selection in (32) changes.

In cells where  $v_{k-\frac{1}{2}}^L < 0 < v_{k-\frac{1}{2}}^R$  it is well known that the Lax-Wendroff scheme may produce a non-physical solution, and thus we cannot use this comparison

to claim second order accuracy at such points.

#### 4. Convergence of the Scheme

We now consider convergence of the scheme after first giving some useful definitions.

$$\mu_{k-\frac{1}{2}}^s = \left\{ \begin{array}{ll} -\phi_{k-\frac{1}{2}}^s / (u_k - u_{k-1}) & u_k \neq u_{k-1} \\ \left\{ \begin{array}{ll} (-v_{k-\frac{1}{2}})^+ & s = L \\ (v_{k-\frac{1}{2}})^+ & s = R \end{array} \right. & u_k = u_{k-1} \end{array} \right\} \quad s = L, R \quad (38)$$

$$\beta_{k-\frac{1}{2}}^s = \left\{ \begin{array}{ll} b_{k-\frac{1}{2}}^s / \alpha_{k-\frac{1}{2}} \phi_{k-\frac{1}{2}}^s & \phi_{k-\frac{1}{2}}^s \neq 0 \\ 1 & \phi_{k-\frac{1}{2}}^s = 0 \end{array} \right\} \quad s = L, R \quad (39)$$

$$\gamma_{k-\frac{1}{2}}^s = \left\{ \begin{array}{ll} b_{k-\frac{1}{2}}^s / \alpha_{k+\frac{1}{2}} \phi_{k+\frac{1}{2}}^s & \phi_{k+\frac{1}{2}}^s \neq 0 \\ 1 & \phi_{k+\frac{1}{2}}^s = 0 \end{array} \right\} \quad s = L, R \quad (40)$$

We observe, from (24) and (31), (32), (39), (40), the inequalities

$$0 \leq \alpha_{k-\frac{1}{2}} \leq \frac{1}{2}$$

$$|\beta_{k-\frac{1}{2}}^s| \leq 1, \quad |\gamma_{k-\frac{1}{2}}^s| \leq 1 \quad s = L, R \quad (41)$$

and from (24) and (27), (38),

$$0 \leq \mu_{k-\frac{1}{2}}^R \leq \frac{1}{2}$$

$$0 \leq -\mu_{k-\frac{1}{2}}^L \leq \frac{1}{2}.$$

For convergence we require a uniform bound on the variation of the solution and also a uniform bound on the solution itself [2]. Sufficient conditions for the total variation to be non-increasing, and hence the variation to be bounded due to (2) are:

$$0 \leq -\xi_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^L \quad (43a)$$

$$0 \leq \zeta_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^R$$

$$0 \leq \zeta_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^R - \xi_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^L \leq 1, \quad (43b)$$

see [2].

The criterion we use for bounding the solution is

$$\inf\{u_{k-1}, u_k, u_{k+1}\} \leq u^k \leq \sup\{u_{k-1}, u_k, u_{k+1}\} \quad (44)$$

which ensures that, as well as preserving monotonicity of the data (also implied by a non-increasing total variation), maxima will not increase nor minima decrease. This implies in turn that the solution is bounded by the initial data in the  $L^\infty$  norm [2]. An additional sufficient condition to (43a) for (44) to hold is

$$0 \leq \zeta_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^R - \xi_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^L \leq 1. \quad (45)$$

[If both (45) and (43a) hold the solution will also be bounded in the  $l_1$  norm by the initial data (see [7])]

We note that, using (39), (40), equation (34) may be rewritten as

$$\begin{aligned} \xi_{k+\frac{1}{2}} &= 1 + (\gamma_{k-\frac{1}{2}}^L - \beta_{k+\frac{1}{2}}^L) \alpha_{k+\frac{1}{2}} \\ \zeta_{k-\frac{1}{2}} &= 1 - (\beta_{k-\frac{1}{2}}^R - \gamma_{k+\frac{1}{2}}^L) \alpha_{k+\frac{1}{2}} \end{aligned} \quad (46)$$

The inequalities (43a) then follow from (41).

Away from sign changes of  $f'(u)$  either  $\phi^L$  or  $\phi^R$  is zero. Consequently either  $\mu^L$  or  $\mu^R$  is zero, and the inequalities (43b) and (45) are easily verified [2]. It only remains to verify (43b) and (45) for expansions and shocks. This is done in (i) and (ii) below.

(i) Expansions

Consider first an expansion, the scheme for which is shown graphically in

Fig. 6

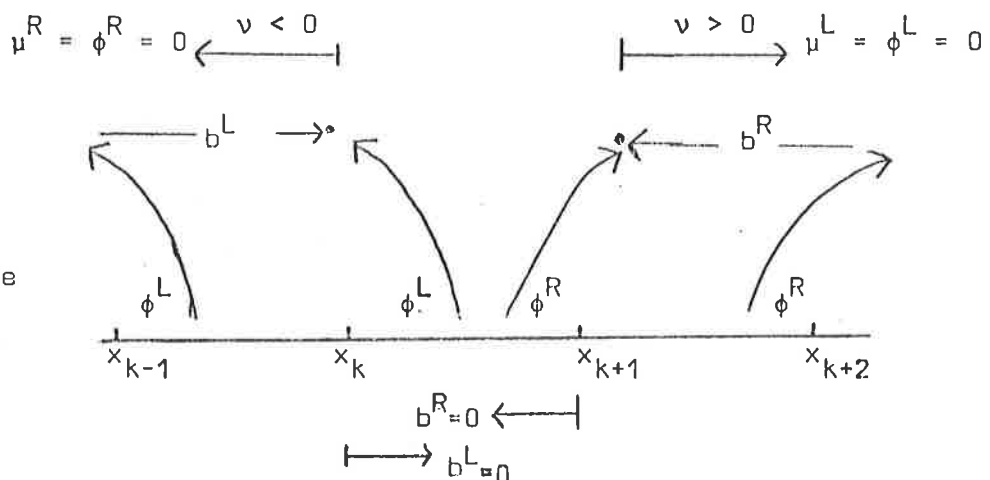


Fig. 6  
The Scheme  
for an  
expansion

It is easily seen that (45) holds. Now consider

$$\zeta_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^R - \xi_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^L = g(v_{k+\frac{1}{2}}), \text{ say.} \quad (47)$$

Clearly, from (43a),  $g(v_{k+\frac{1}{2}}) \geq 0$ .

Since  $b_{k+\frac{1}{2}}^L = b_{k+\frac{1}{2}}^R = 0$  and  $u_k \neq u_{k+1}$  for an expansion we may write

$$\begin{aligned} g(v_{k+\frac{1}{2}}) &= (1 + \gamma_{k-\frac{1}{2}}^L \alpha_{k+\frac{1}{2}}) \frac{\phi_{k+\frac{1}{2}}^L}{u_{k+1} - u_k} - (1 + \gamma_{k+\frac{3}{2}}^R \alpha_{k+\frac{1}{2}}) \frac{\phi_{k+\frac{1}{2}}^R}{u_{k+1} - u_k} \\ &= -v_k (1 + \gamma_{k-\frac{1}{2}}^L \alpha_{k+\frac{1}{2}}) \frac{(u_{k+\frac{1}{2}} - u_k)}{(u_{k+1} - u_k)} + v_{k+1} (1 + \gamma_{k+\frac{3}{2}}^R \alpha_{k+\frac{1}{2}}) \frac{(u_{k+1} - u_{k+\frac{1}{2}})}{(u_{k+1} - u_k)}. \end{aligned} \quad (48)$$

Using (19) we obtain

$$\frac{u_{k+\frac{1}{2}} - u_k}{u_{k+1} - u_k} = \frac{v_{k+1} - v_{k+\frac{1}{2}}}{v_{k+1} - v_k} \quad (49)$$

and

$$\frac{u_{k+1} - u_{k+\frac{1}{2}}}{u_{k+1} - u_k} = \frac{v_{k+\frac{1}{2}} - v_k}{v_{k+1} - v_k}.$$

Again for an expansion  $v_{k+1} \neq v_k$ , so the denominators are non-zero.

Substituting in (48) yields

$$\begin{aligned} g(v_{k+\frac{1}{2}}) &= -v_k \frac{(v_{k+1} - v_{k+\frac{1}{2}})}{(v_{k+1} - v_k)} (1 + \gamma_{k-\frac{1}{2}}^L \alpha_{k+\frac{1}{2}}) + v_{k+1} \frac{(v_{k+\frac{1}{2}} - v_k)}{(v_{k+1} - v_k)} (1 + \gamma_{k+\frac{3}{2}}^R \alpha_{k+\frac{1}{2}}) \\ &\leq \frac{3}{2} \left\{ \frac{-v_k (v_{k+1} - v_{k+\frac{1}{2}}) + v_{k+1} (v_{k+\frac{1}{2}} - v_k)}{v_{k+1} - v_k} \right\}, \end{aligned} \quad (50)$$

(using (41)), and therefore

$$\leq \max \left( -\frac{3}{2}v_k, \frac{3}{2}v_{k+1} \right) \leq 1, \quad (51)$$

since (50) is linear in  $v_{k+\frac{1}{2}}$ .

Thus (43b) holds, which completes the conditions for expansions.

(ii) Shocks

Next consider a shock, the scheme for which is illustrated in Fig. 7.

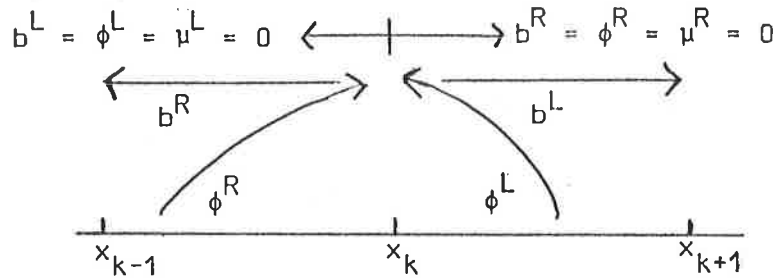


Fig. 7

The Scheme for a Shock

We have immediately that

$$0 \leq \zeta_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^R - \xi_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^L \leq 1,$$

verifying (43b), and also that

$$0 \leq \zeta_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^R - \xi_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^L.$$

Moreover

$$\begin{aligned} \zeta_{k-\frac{1}{2}} \mu_{k-\frac{1}{2}}^R - \xi_{k+\frac{1}{2}} \mu_{k+\frac{1}{2}}^L &= (1 - \beta_{k-\frac{1}{2}}^R \alpha_{k-\frac{1}{2}}) v_{k-\frac{1}{2}} - (1 - \beta_{k-\frac{1}{2}}^L \alpha_{k-\frac{1}{2}}) v_{k+\frac{1}{2}} \\ &\leq (1 + \alpha_{k-\frac{1}{2}}) v_{k-\frac{1}{2}} - (1 + \alpha_{k+\frac{1}{2}}) v_{k+\frac{1}{2}} \\ &\leq 1 \quad \text{using (31), if } |v_{k+\frac{1}{2}}| < 0.35 \text{ approximately} \end{aligned}$$

This result indicates that the inequality (45) may be violated if  $0.35 < |v_{k+\frac{1}{2}}| < 0.5$ . On close inspection we see that this only happens in the extreme case of oscillatory data with  $|v| \sim \frac{1}{2}$  adjacent to a shock. Although this is unlikely to occur for most convex  $f$ , we may achieve (45) for certain by selecting the sign of the first argument of the minmod operator adjacent to a shock.

Finally from the inequalities in (i), (ii) above we may use compactness arguments [2], [7] to show that we may select from  $\{u_h\}$  a convergent subsequence, and have therefore shown the convergence of the scheme (33), (34).

## 5. Entropy for the Semi-Discrete First Order Scheme

The first order semi-discrete scheme may be written as

$$\frac{\partial u_k^n}{\partial t} + \frac{1}{\Delta x} \Delta_- [f(u_{k+1}^n)(1-\theta_{k+\frac{1}{2}}^n) + \theta_{k+\frac{1}{2}}^n f(u_k^n)] = 0 \quad (52)$$

$$\text{where } \theta_{k+\frac{1}{2}}^n = \begin{cases} 1 & v > 0 \text{ in } (u_k, u_{k+1}) \\ 0 & v < 0 \text{ in } (u_k, u_{k+1}) \\ \frac{v_{k+\frac{1}{2}} - v_k}{v_{k+1} - v_k} \frac{v_{k-1}}{v_{k+\frac{1}{2}}} & v_k < 0 < v_{k+1} \\ \frac{1}{2}(1 + \text{sgn}(v_{k-\frac{1}{2}})) & v_{k+1} < 0 < v_k \end{cases} \quad (53)$$

We follow Osher [8] and multiply by  $u_k^n \psi_k^n \Delta x$  where  $\psi_k^n$  is a positive test function, and sum over  $k, n$ , giving

$$\sum_k \left\{ \psi_k^n u_k^n \frac{\partial u_k^n}{\partial t} + \frac{1}{\Delta x} \psi_k^n u_k^n \Delta_- [f(u_{k+1}^n)(1-\theta_{k+\frac{1}{2}}^n) + \theta_{k+\frac{1}{2}}^n f(u_k^n)] \right\} \Delta x = 0, \quad (54)$$

or

$$\sum_k \left\{ \psi_k^n \frac{\partial}{\partial t} (\frac{1}{2} u_k^n{}^2) + \frac{1}{\Delta x} \psi_k^n u_k^n [\Delta_+ f_k^n - \Delta_- \{\theta_{k+\frac{1}{2}}^n \Delta_+ f_k^n\}] \right\} \Delta x = 0 \quad (55)$$

As in Osher's paper we now add and subtract a term

$$\frac{1}{\Delta x} \psi_k^n \int_{u_k^n}^{u_{k+1}^n} s f'(s) ds$$

and sum the  $\Delta_-$  term by parts over  $k$  to give

$$\begin{aligned} & \sum_k \left\{ \psi_k^n \frac{\partial}{\partial t} (\frac{1}{2} u_k^n{}^2) - \frac{1}{\Delta x} (\Delta_+ \psi_k^n) \int_{\bar{u}}^{u_{k+1}^n} s f'(s) ds \right. \\ & \left. + \frac{1}{\Delta x} (\Delta_+ \psi_k^n) u_{k+1}^n \theta_{k+\frac{1}{2}}^n \Delta_+ f_k^n - \frac{1}{\Delta x} \psi_k^n \int_{u_k^n}^{u_{k+1}^n} [s - u_k^n - (\Delta_+ u_k^n) \theta_{k+\frac{1}{2}}^n] f'(s) ds \right\} \Delta x = 0, \quad (56) \end{aligned}$$

where  $f'(\bar{u}) = 0$ .

Note that the third term goes to zero with  $\Delta x$ . After integration by parts we obtain

$$\begin{aligned} & \sum_k \left\{ \psi_k^n \frac{\partial}{\partial t} (\frac{1}{2} u_k^n{}^2) - \frac{1}{\Delta x} (\Delta_+ \psi_k^n) \int_{\bar{u}}^{u_{k+1}^n} s f'(s) ds \right\} \Delta x \\ & = + \sum_k \frac{1}{\Delta x} \psi_k^n \int_{u_k^n}^{u_{k+1}^n} [s - u_k^n - (\Delta_+ u_k^n) \theta_{k+\frac{1}{2}}^n] f'(s) ds \Delta x. \quad (57) \end{aligned}$$

As  $\Delta x, \Delta t \rightarrow 0$  the left hand side tends to

$$\int \left\{ \psi \frac{\partial}{\partial t} \left( \frac{1}{2} u^2 \right) - \psi_x \int_{\underline{u}}^{\underline{u}} s f'(s) ds \right\} dx \quad (58)$$

We now consider the sign of the remaining term by examining the integral (dropping  $n$  superscripts) in (57), namely,

$$I_{k+\frac{1}{2}} = \int_{u_k}^{u_{k+1}} [s - u_k - \theta_{k+\frac{1}{2}} \Delta_+ u_k] f'(s) ds, \quad (59)$$

for the possible values of  $\theta_{k+\frac{1}{2}}$ .

If  $f'(s)$  is of constant sign in  $(u_k, u_{k+1})$  then  $\theta_{k+\frac{1}{2}} = 0$  or  $1$ .

Take the case  $f' > 0$ ,  $\theta_{k+\frac{1}{2}} = 1$ : then

$$I_{k+\frac{1}{2}} = \int_{u_k}^{u_{k+1}} [s - u_{k+1}] f'(s) ds \leq 0,$$

and similarly, for  $f' < 0$ ,  $\theta_{k+\frac{1}{2}} = 0$ ,

$$I_{k+\frac{1}{2}} = \int_{u_k}^{u_{k+1}} [s - u_k] f'(s) ds \leq 0.$$

If, however,  $f'(u_k) > 0 > f'(u_{k+1})$  we have  $u_k > \bar{u} > u_{k+1}$  (since  $f$  is convex) and

$$\theta_{k+\frac{1}{2}} = \begin{cases} 1 & \text{if } \frac{\Delta_+ f_k}{\Delta_+ u_k} = \frac{\int_{u_k}^{u_{k+1}} f'(s) ds}{\Delta_+ u_k} \geq 0 \\ 0 & \text{if } \frac{\Delta_+ f_k}{\Delta_+ u_k} = \frac{\int_{u_k}^{u_{k+1}} f'(s) ds}{\Delta_+ u_k} < 0 \end{cases} \quad (60)$$

i.e.

$$\theta_{k+\frac{1}{2}} = \begin{cases} 1 & \text{if } \int_{u_{k+1}}^{u_k} f'(s) ds \geq 0 \\ 0 & \text{if } \int_{u_{k+1}}^{u_k} f'(s) ds < 0 \end{cases} \quad (61)$$



Taking the first of these as an example:

$$\begin{aligned}
 I_{k+\frac{1}{2}} &= \int_{u_{k+1}}^{u_k} [u_{k+1} - s] f'(s) ds \\
 &= [(u_{k+1} - s)f(s)]_{u_{k+1}}^{u_k} + \int_{u_{k+1}}^{u_k} f(s) ds \\
 &= (u_{k+1} - u_k)f(u_k) - (u_{k+1} - u_k)f(\xi) \quad \xi \in (u_{k+1}, u_k) \\
 &= (u_{k+1} - u_k)(f(u_k) - f(\xi)) \\
 &\leq 0 \quad \text{since} \quad \int_{\xi}^{u_k} f'(s) ds \geq \int_{u_{k+1}}^{u_k} f'(s) ds \geq 0
 \end{aligned}$$

If  $\int_{u_{k+1}}^{u_k} f'(s) ds < 0$  a similar argument gives  $I_{k+\frac{1}{2}} \leq 0$ .

The remaining case is  $f'(u_{k+1}) > 0 > f'(u_k)$  with  $u_{k+1} > \bar{u} > u_{k+1}$ .

Here  $\theta_{k+\frac{1}{2}} = \frac{v_{k+\frac{1}{2}} - v_k}{v_{k+1} - v_k} \frac{v_{k+1}}{v_{k+\frac{1}{2}}}$  with  $v_k \leq v_{k+\frac{1}{2}} \leq v_{k+1}$

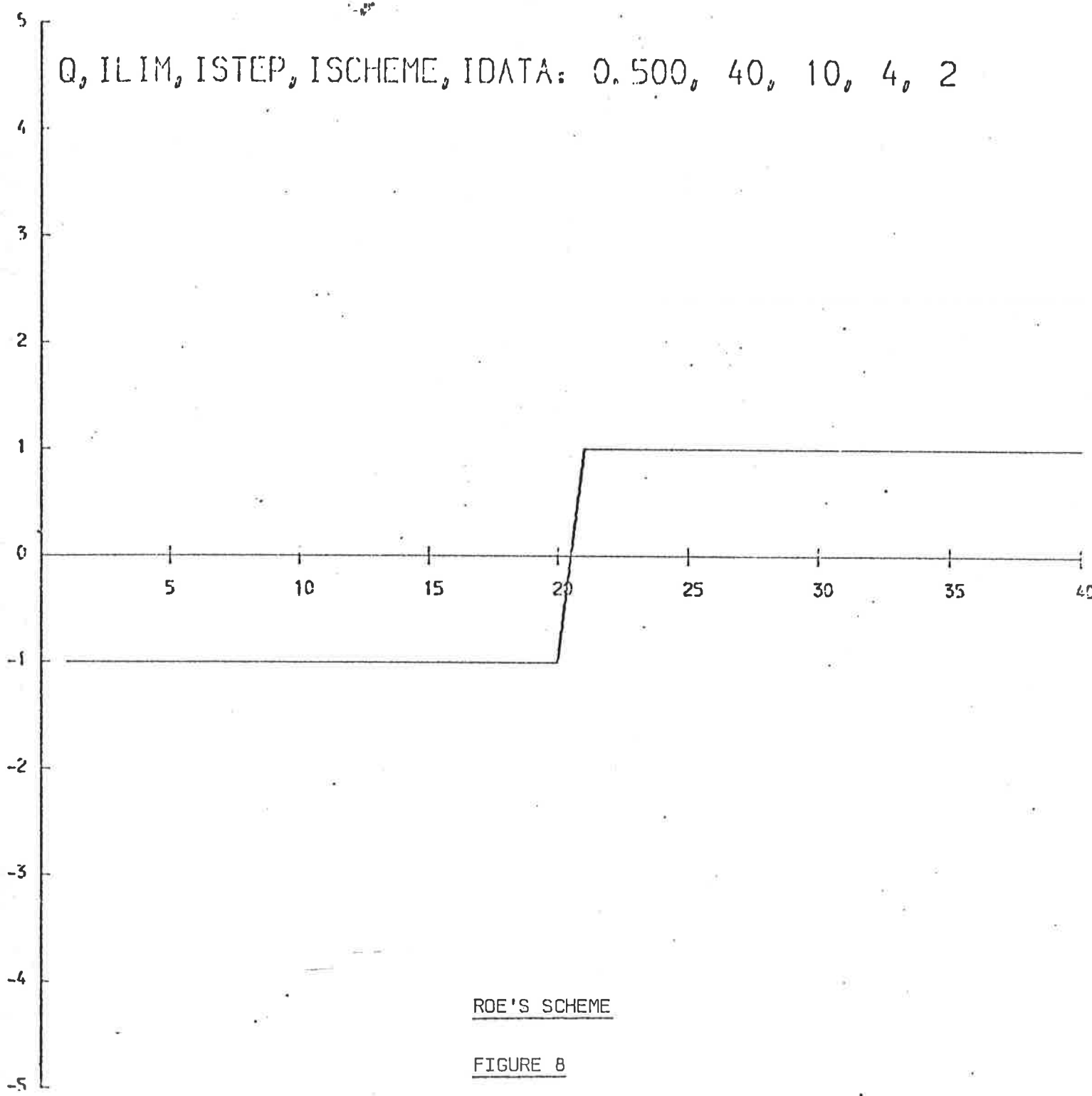
so, if  $\int_{u_k}^{u_{k+1}} f'(s) ds \geq 0$ ,  $\theta_{k+\frac{1}{2}} \geq 1$ :

we have  $\int_{u_k}^{u_{k-1}} f'(s) ds < 0$ ,  $\theta_{k+\frac{1}{2}} \leq 0$ , which gives

$$\begin{aligned}
 I_{k+\frac{1}{2}} &= \int_{u_k}^{u_{k+1}} [s-c] f'(s) ds, \quad \text{where } c \geq u_{k+1} \quad \text{if } \int_{u_k}^{u_{k+1}} f'(s) ds \geq 0 \\
 & \quad \quad \quad c \leq u_k \quad \text{if } \int_{u_k}^{u_{k+1}} f'(s) ds < 0.
 \end{aligned} \tag{62}$$

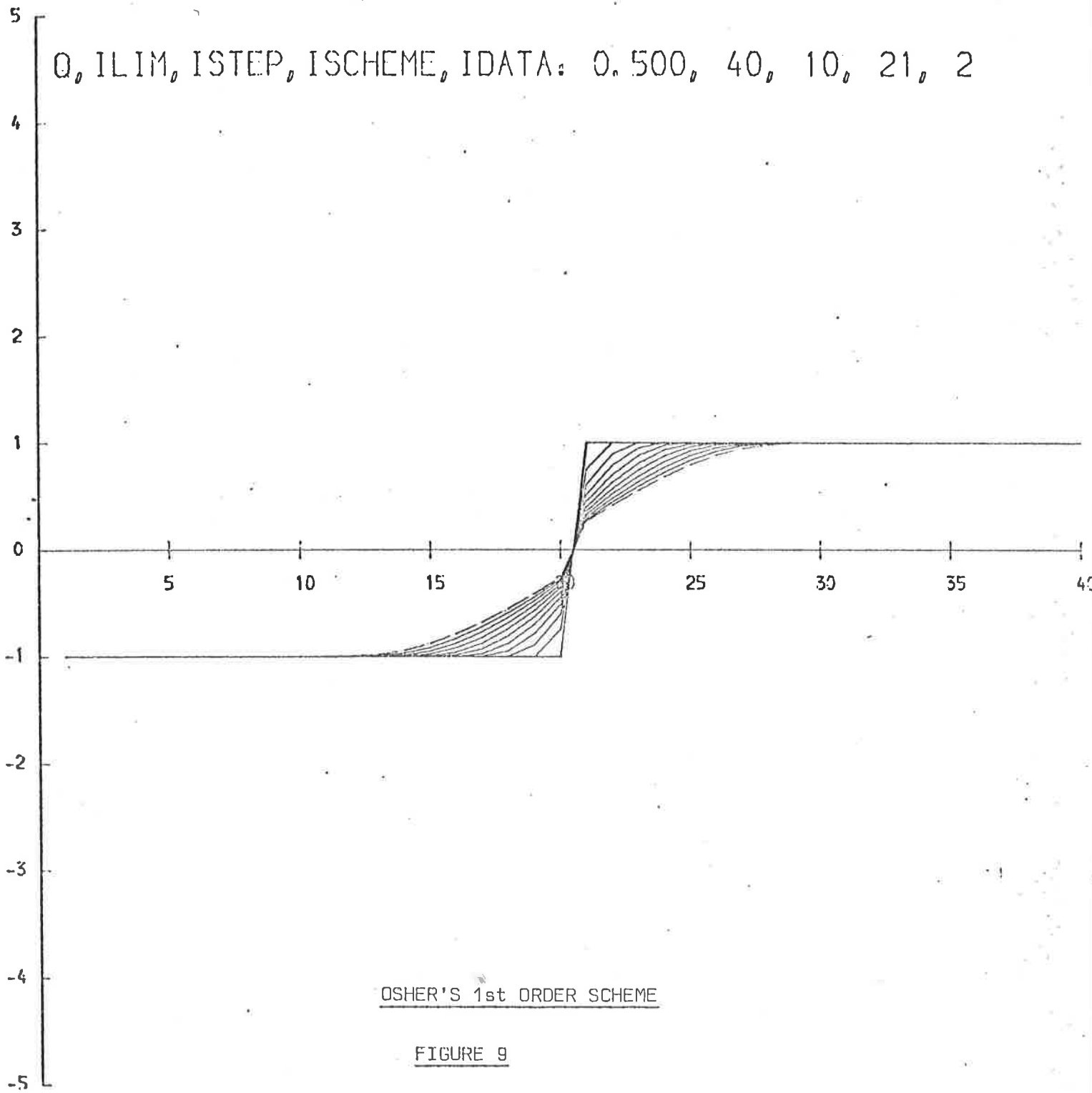
Taking the case  $\int_{u_k}^{u_{k+1}} f'(s) ds \geq 0$ ,  $c \geq u_{k+1}$ , we have

$$\begin{aligned}
 I_{k+\frac{1}{2}} &= \int_{u_k}^{u_{k+1}} [s-c] f'(s) ds \\
 &= [(s-c)f(s)]_{u_k}^{u_{k+1}} - \int_{u_k}^{u_{k+1}} f(s) ds \\
 &= (u_{k+1}-c)f(u_{k+1}) - (u_k-c)f(u_k) - (u_{k+1}-u_k)f(\xi) \quad \xi \in (u_k, u_{k+1}) \\
 &= (u_k-c)\{f(u_{k+1})-f(u_k)\} + (u_{k+1}-u_k)\{f(u_{k+1})-f(\xi)\}. \tag{63}
 \end{aligned}$$

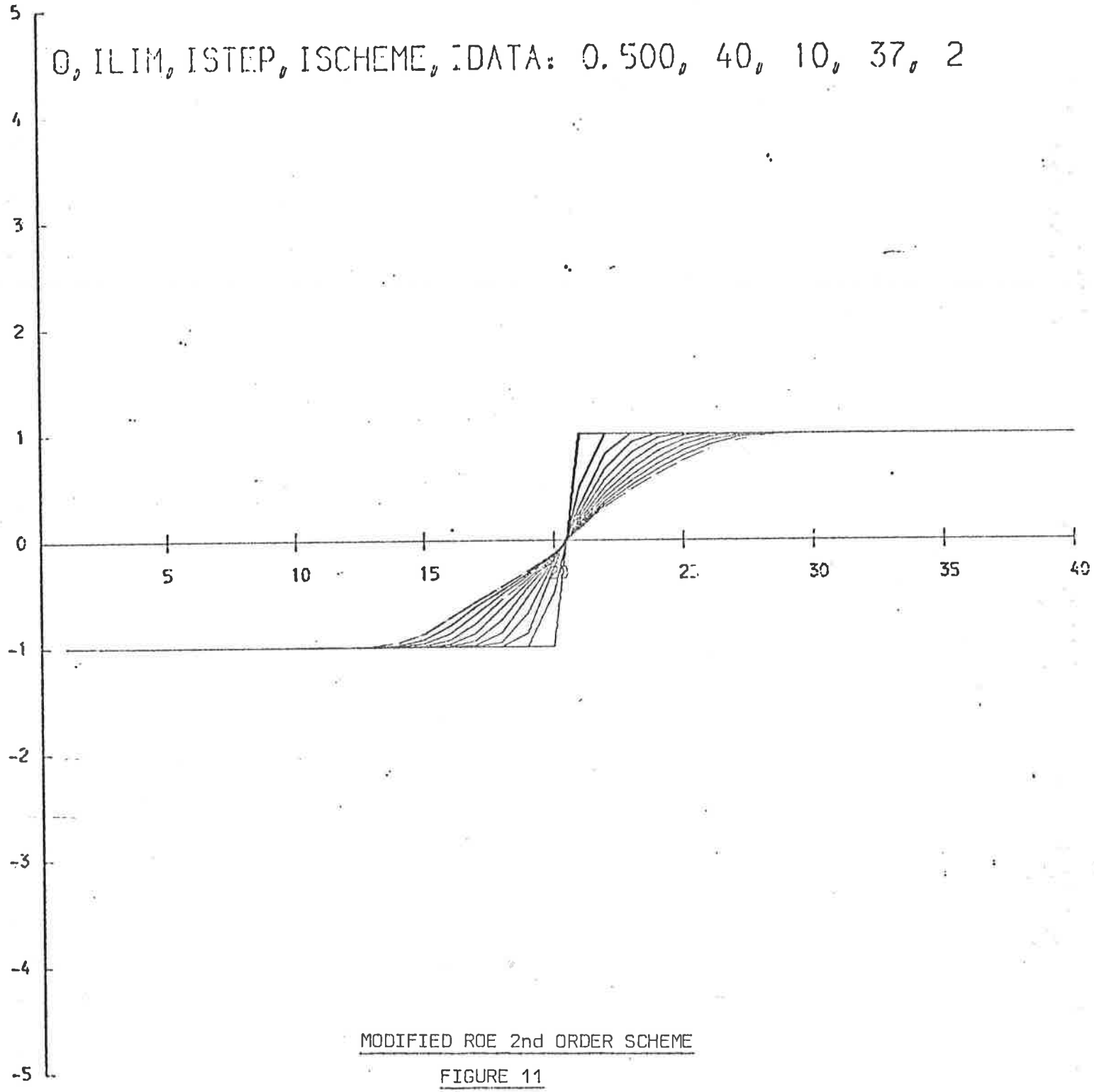


ROE'S SCHEME

FIGURE 8

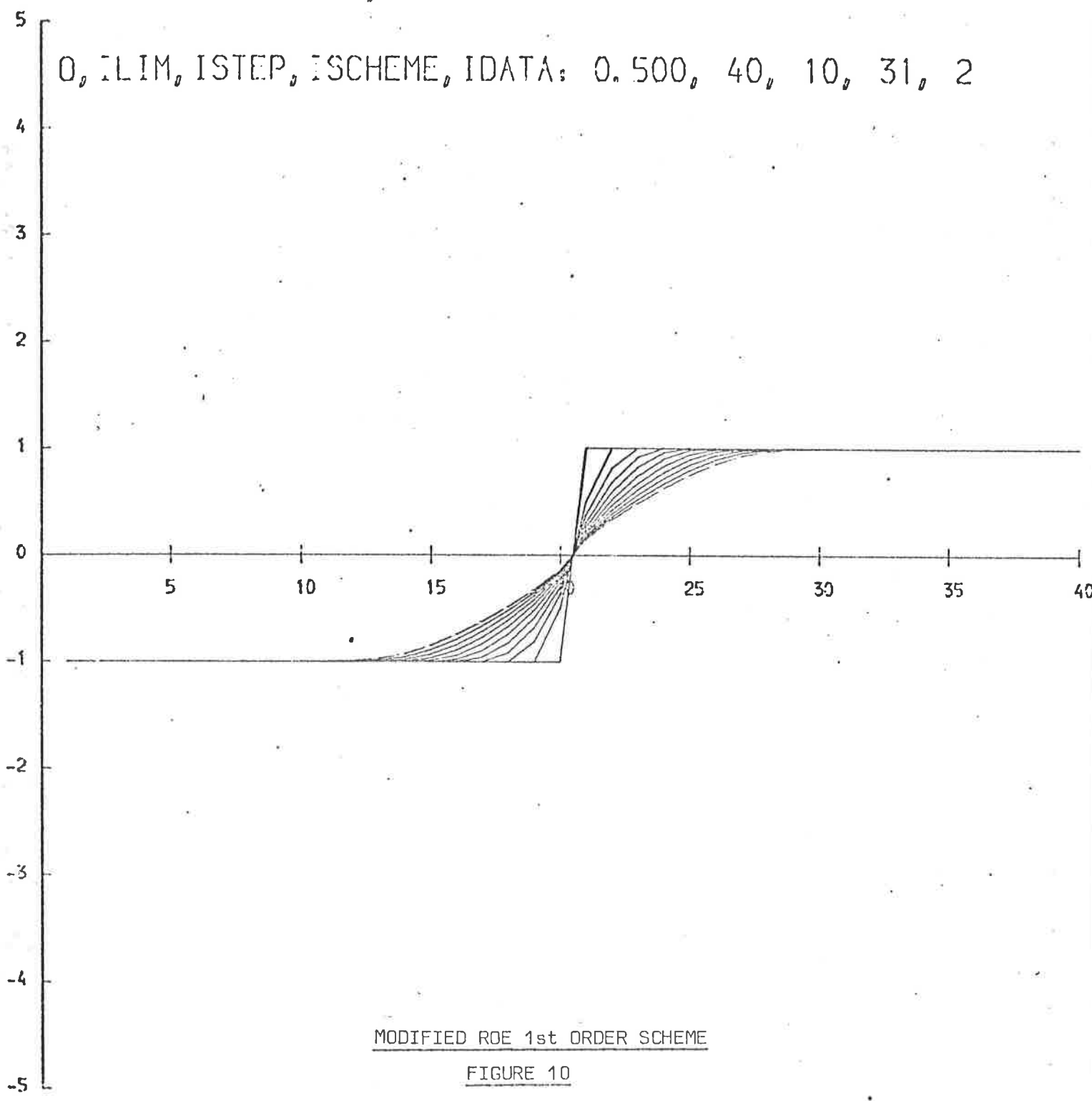


0, ILIM, ISTEP, ISCHEME, I DATA: 0.500, 40, 10, 37, 2



MODIFIED ROE 2nd ORDER SCHEME

FIGURE 11



The first term is negative, whilst the second is positive since  $\int_{\xi}^{u_{k+1}} f'(s) ds \geq \int_{u_k}^{u_{k+1}} f'(s) ds \geq 0$ . Comparing absolute values of the terms,

$$\left| \frac{(u_{k+1} - u_k)(f(u_{k+1}) - f(\xi))}{(c - u_k)(f(u_{k-1}) - f(u_k))} \right| \leq 1, \quad (64)$$

since  $c \geq u_{k+1}$ : hence the first term is dominant giving  $I_{k+\frac{1}{2}} \leq 0$ . Similarly if  $\int_{u_k}^{u_{k-1}} f'(s) ds \leq 0$ ,  $c \leq u_k$  we again have  $I_{k+\frac{1}{2}} \leq 0$ .

Note that if  $v_{k+\frac{1}{2}} = 0$  we may consider  $\theta_{k+\frac{1}{2}}$  indefinitely large, i.e.  $c \gg u_{k+1}$ , for the purposes of this proof: in the scheme itself  $\theta_{k+\frac{1}{2}} v_{k+\frac{1}{2}}$  is finite.

Collecting these results together, we have from (59)

$$\int \left\{ \psi \frac{\partial}{\partial t} \left( \frac{1}{2} u^2 \right) - \psi_x \int_u^U s f'(s) ds \right\} dx \leq 0,$$

that is, in the sense of distributions,

$$\frac{\partial V(u)}{\partial t} + \frac{\partial G(u)}{\partial x} \leq 0.$$

where  $V(u) = \frac{1}{2} u^2$  is the entropy function, and  $G(u) = \int_u^U V'(s) f'(s) ds$  is the entropy flux.

It follows that the solutions of the first order scheme converge to the correct entropy satisfying solution of the equation

$$u_t + f(u)_x = 0$$

for convex  $f(u)$ . (See [6]).

## 5. Results and Conclusions

Comparisons of Roe's Original Scheme [3], Osher's first order Scheme [8], and the modified Roe 1st and 2nd order schemes described here, are shown for a centred rarefaction in Figs. 8-11, respectively. The test equation used is

$$u_t + \left( \frac{1}{2} u^2 \right)_x = 0.$$

It can be seen that Roe's original scheme produces an entropy violating solution, whilst Osher's, although giving the correct physical solution,

contains a large "dog-leg" compared to the two modified Roe schemes.

### Acknowledgements

I would like to thank Dr. M.J. Baines of the University of Reading and Prof. S. Osher of the University of California, Los Angeles for invaluable discussions. I would also like to acknowledge the financial support of the Science & Engineering Research Council and the Royal Aircraft Establishment, Farnborough.

### References

- [1] Lax, P.D. Hyperbolic Conservation Laws and Mathematical Theory of Shock Waves, 1972, SIAM Regional Conference Series, Lectures in Applied Mathematics, II.
- [2] Sweby, P.K. & Baines, M.J. Convergence of Roe's Scheme for the general non-linear scalar wave equation (1981). Numerical Analysis Report 8/81, University of Reading, also submitted to Math. Comp.
- [3] Roe, P.L. Numerical algorithms for the linear wave equation, 1981, RAE Technical Report 81047. Also Roe, P.L. & Baines, M.J., 1981. Proceedings of the 4th GAMM Conference: Notes on Numerical Fluid Mechanics, Vol. 5 (Vieweg).
- [4] Lax & Wendroff. Comm. Pure Applied Maths, 13, p. 212, 1960.
- [5] Warming & Beam. AIAA Jnl. 14, p. 1241, 1976.
- [6] Kruzkov, S.N. First order quasilinear equations in several independent variables. Math. USSR Sbornik, 1970, Vol. 10, p. 217-243. Also, Oleinik, O.A. Discontinuous solutions of non-linear differential equations. Amer. Math. Soc. Trans. Ser. 2, No. 26, p. 95-172.
- [7] Sanders. On convergence of Finite Difference Schemes with Variable Spatial Differencing. Ph.D. abstract, Mathematics Dept. University of California, Los Angeles, (1980).
- [8] Engquist, B., & Osher, S. Stable and entropy satisfying approximations for transonic flow calculations. Math. Comp. 34 No. 149, Jan. 1980, p. 45-75.