

# Numerical methods for one-dimensional hyperbolic conservation laws

***Citation for published version (APA):***

Berkenbosch, A. C., Kaasschieter, E. F., & Thijs Boonkamp, ten, J. H. M. (1992). *Numerical methods for one-dimensional hyperbolic conservation laws*. (RANA : reports on applied and numerical analysis; Vol. 9215). Eindhoven University of Technology.

***Document status and date:***

Published: 01/01/1992

***Document Version:***

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

***Please check the document version of this publication:***

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

***General rights***

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

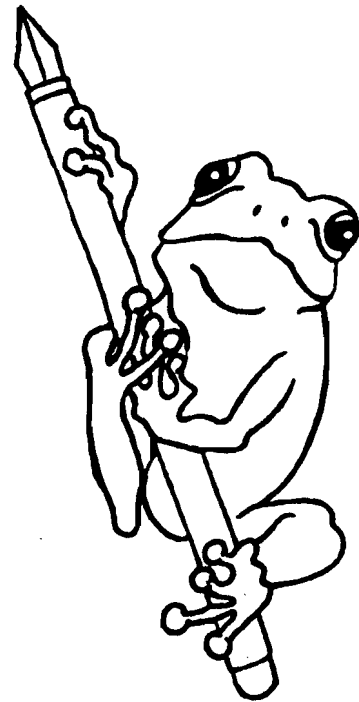
***Take down policy***

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

RANA 92-15  
October 1992  
Numerical Methods for  
One-Dimensional  
Hyperbolic Conservation Laws  
by  
A.C. Berkenbosch  
E.F. Kaasschieter  
J.H.M. ten Thije Boonkkamp



Reports on Applied and Numerical Analysis  
Department of Mathematics and Computing Science  
Eindhoven University of Technology  
P.O. Box 513  
5600 MB Eindhoven  
The Netherlands  
ISSN: 0926-4507

# Numerical Methods for One-Dimensional Hyperbolic Conservation Laws

A.C. Berkenbosch, E.F. Kaasschieter and J.H.M. ten Thije Boonkkamp

*Eindhoven University of Technology,  
Department of Mathematics and Computing Science,  
P.O. Box 513, 5600 MB Eindhoven, The Netherlands.*

## Abstract

This paper contains a survey of some important numerical methods for one-dimensional hyperbolic conservation laws. Weak solutions of hyperbolic conservation laws are introduced and the concept of entropy stability is discussed. Furthermore, the Riemann problem for hyperbolic conservation laws is solved. An introduction to numerical methods is given for which important concepts such as e.g. conservativity, stability and consistency are introduced. Godunov-type methods are elaborated for general systems of hyperbolic conservation laws. Finally, flux limiter methods are developed for the scalar non-linear conservation law.

A.M.S. Classifications: 35A40, 35L65, 65M06, 65M99, 76M99

Keywords : Hyperbolic conservation laws, Euler equations, entropy, Riemann problem, conservative schemes, Godunov-type schemes, high resolution schemes, flux limiters.

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Introduction to Hyperbolic Conservation Laws</b>	<b>3</b>
2.1	Definition of hyperbolic conservation laws . . . . .	3
2.2	Solutions of hyperbolic conservation laws . . . . .	5
<b>3</b>	<b>The Riemann Problem for Hyperbolic Conservation Laws</b>	<b>8</b>
3.1	Preliminaries . . . . .	8
3.2	Solution of the Riemann problem . . . . .	8
3.3	Riemann invariants . . . . .	12
<b>4</b>	<b>Introduction to Numerical Methods</b>	<b>16</b>
4.1	Some basic concepts . . . . .	16
4.2	Examples of conservative methods . . . . .	19
4.3	Modified equations . . . . .	22
4.4	Numerical entropy stability . . . . .	24
<b>5</b>	<b>Godunov-type methods</b>	<b>26</b>
5.1	Introduction . . . . .	26
5.2	The basic Godunov method . . . . .	26
5.3	Osher's method . . . . .	28
5.4	Roe's method . . . . .	33
<b>6</b>	<b>High Resolution Methods</b>	<b>38</b>
6.1	Some convergence results . . . . .	38
6.2	Flux limiter methods . . . . .	41
	<b>References</b>	<b>46</b>

# 1 Introduction

Many (practical) problems in science and engineering involve conservation laws. A special class are the so-called hyperbolic conservation laws, which can be formulated as a system of first order partial differential equations. An important example are the Euler equations of gas dynamics (cf. [18]). Other examples arise in meteorology and astrophysics. In general it is not possible to derive exact solutions of these equations, and therefore, we have to devise and study numerical methods to approximate solutions (cf. [14]).

Apart from the practical applications, there are two other reasons for studying numerical methods for hyperbolic conservation laws. Firstly, there are special difficulties associated with solving hyperbolic conservation laws (e.g. shock formation) that must be dealt with carefully in developing numerical methods. Methods based on naive finite difference approximations may behave well for smooth solutions but can give disastrous results when discontinuities are present. Secondly, a great deal is known about the mathematical structure of these equations and their solutions (cf. [14]). This theory can be exploited to develop special methods that overcome some of the numerical difficulties arising from a more naive approach.

Usually practical problems are in two or three space dimensions. However, most of the methods currently in use are heavily based on one-dimensional methods, generalized by 'dimensional splitting' or similar techniques. For this reason we will consider only one-dimensional conservation laws.

Many methods have been derived for the one-dimensional hyperbolic conservation law. Methods developed using straightforward finite difference discretizations are inappropriate near discontinuities, since they are based on truncated Taylor series expansions. A survey of these methods (with applications to the Euler equations) is given in [9] and [10]. Another important class of numerical methods are the Godunov-type methods (cf. [3], [8], [16]). These methods use, in some way, the exact solution of the Riemann problem and do not produce oscillations around discontinuities. Unfortunately, these methods are only first order; hence the solutions are smoothed around discontinuities. Therefore, other methods have been developed. A very popular class of methods are the high resolution methods, which are second order accurate in smooth regions and give good results (no oscillations) around shocks (cf. [5], [6], [21], [28], [30]).

This paper is organized as follows. In the next section one-dimensional hyperbolic conservation laws are introduced. Furthermore, weak solutions of these equations are defined. These weak solutions turn out to be non-unique and therefore an extra condition (i.e. entropy stability) is introduced to characterize the physically relevant solution. In Section 3, the Riemann problem is introduced and solved. This Riemann problem is important, because it forms the underlying physical model for the Godunov-type methods. Section 4 is of preliminary nature. Some basic numerical concepts are introduced, which are important to study the behaviour of numerical methods. Furthermore, some well-known methods are discussed for the scalar, linear convection equation. In Section 5, Godunov-type methods are discussed. Two examples of these methods are given, namely a method developed by Osher (cf. [22]) and a method developed by Roe (cf. [24]). Finally in the last section, high resolution methods are introduced. As an important example the flux limiter methods are considered in more detail for non-linear scalar conservation laws.

## 2 Introduction to Hyperbolic Conservation Laws

### 2.1 Definition of hyperbolic conservation laws

In this report we consider *nonlinear conservation laws* in one space dimension. The general form of such conservation laws is (cf. e.g. [17], [32])

$$\frac{d}{dt} \int_{x_1}^{x_2} u(x, t) dx = f(u(x_1, t)) - f(u(x_2, t)), \quad (2.1)$$

stating that the rate of change of the variable  $\int_{x_1}^{x_2} u(x, t) dx$  is equal to the difference in the fluxes  $f(u(x, t))$  at  $x_1$  and at  $x_2$ . Another integral form is obtained by integrating (2.1) in time, giving

$$\int_{t_1}^{t_2} \int_{x_1}^{x_2} \left\{ \frac{\partial}{\partial t} u(x, t) + \frac{\partial}{\partial x} f(u(x, t)) \right\} dx dt = 0. \quad (2.2)$$

Since (2.2) should hold for arbitrary  $x_1, x_2$  and  $t_1, t_2$ , the integrand in (2.2) must be equal to 0, i.e.

$$\frac{\partial}{\partial t} u(x, t) + \frac{\partial}{\partial x} f(u(x, t)) = 0. \quad (2.3)$$

This is the differential form of the conservation law, which only holds if the solution  $u : \mathbb{R} \times [0, \infty) \rightarrow \mathbb{R}^m$  and the *flux-function*  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$  are continuously differentiable. Finally, the quasi-linear form of the conservation law reads

$$\frac{\partial}{\partial t} u(x, t) + A(u(x, t)) \frac{\partial}{\partial x} u(x, t) = 0, \quad (2.4)$$

where  $A(u)$  is the *Jacobian matrix*, defined by

$$A(u) = \frac{\partial}{\partial u} f(u). \quad (2.5)$$

A *hyperbolic conservation law* is defined as follows (cf. e.g. [14], [17]).

**Definition 2.1** *The system (2.3) is called a hyperbolic conservation law if there exists a real diagonal matrix  $\Lambda(u)$  and a non-singular real matrix  $R(u)$  such that*

$$A(u)R(u) = R(u)\Lambda(u), \quad \forall u \in \mathbb{R}^m. \quad (2.6)$$

Here  $\Lambda(u) = \text{diag}(\lambda_1(u), \lambda_2(u), \dots, \lambda_m(u))$  is the diagonal matrix of the eigenvalues of  $A(u)$  and  $R(u) = (r^{(1)}(u), r^{(2)}(u), \dots, r^{(m)}(u))$  is the matrix of the corresponding right eigenvectors of  $A(u)$ . We assume that the eigenvalues are labeled in increasing order, i.e.  $\lambda_1(u) \leq \lambda_2(u) \leq \dots \leq \lambda_m(u)$ .

A very important example of a system of hyperbolic conservation laws are the *Euler equations* (for a more complete description, cf. e.g. [10], [18]).

**Example 2.2 (The Euler equations)** The Euler equations of gas dynamics describe the flow of an inviscid, non-heat-conducting compressible fluid (a gas). They represent the conservation of *mass*, *momentum* and *energy*. With *density*  $\rho(x, t)$ , *velocity*  $u(x, t)$ ,

*total energy*  $E(x, t)$ , *stagnation enthalpy*  $H(x, t)$  and the *pressure* of the gas  $p(x, t)$ , these conservation laws read respectively

$$\begin{aligned}\frac{\partial}{\partial t}(\rho) + \frac{\partial}{\partial x}(\rho u) &= 0, \\ \frac{\partial}{\partial t}(\rho u) + \frac{\partial}{\partial x}(\rho u^2 + p) &= 0, \\ \frac{\partial}{\partial t}(\rho E) + \frac{\partial}{\partial x}(\rho u H) &= 0,\end{aligned}\tag{2.7}$$

where the enthalpy  $H$  is defined by

$$H = E + \frac{p}{\rho}.$$

The system (2.7) has to be completed with an equation of state, which relates the pressure  $p$  with  $\rho$ ,  $E$  and  $u$ . In the remainder of the paper a *perfect gas* is considered for which the equation of state can be written as

$$p = (\gamma - 1)\rho(E - \frac{1}{2}u^2),\tag{2.8}$$

where  $\gamma = c_p/c_v$  is the *specific heat ratio*. Here  $c_p$  and  $c_v$  are the *specific heats* at constant pressure and at constant volume, respectively. An important quantity which we use later, is the (so-called) *entropy*  $s$ , which is given by

$$s = c_v \ln \frac{p}{\rho^\gamma}.$$

If the vector of *conservative variables*  $\mathbf{u}$  and the flux vector  $\mathbf{f}(\mathbf{u})$  are defined by

$$\mathbf{u} = (\rho, \rho u, \rho E)^T\tag{2.9}$$

and

$$\mathbf{f}(\mathbf{u}) = (\rho u, \rho u^2 + p, \rho u H)^T,\tag{2.10}$$

then the Euler equations can be written in the general form (2.3). Let the Jacobian matrix  $A(\mathbf{u})$  be defined as in (2.5). For the Euler equations the eigenvalues and right eigenvectors of  $A(\mathbf{u})$  are given by (cf. [10])

$$\lambda_1(\mathbf{u}) = u - c, \quad \lambda_2(\mathbf{u}) = u, \quad \lambda_3(\mathbf{u}) = u + c,\tag{2.11}$$

and

$$\begin{aligned}\mathbf{r}^{(1)}(\mathbf{u}) &= -\frac{\rho}{2c}(1, u - c, H - uc)^T, \\ \mathbf{r}^{(2)}(\mathbf{u}) &= (1, u, \frac{u^2}{2})^T, \\ \mathbf{r}^{(3)}(\mathbf{u}) &= \frac{\rho}{2c}(1, u + c, H + uc)^T.\end{aligned}\tag{2.12}$$

In (2.11) and (2.12),  $c$  is the *speed of sound*, which is given by

$$c = \left(\frac{\gamma p}{\rho}\right)^{\frac{1}{2}}.$$

Obviously, the Euler equations are hyperbolic (cf. e.g. [10], [26]).



## 2.2 Solutions of hyperbolic conservation laws

A function  $\mathbf{u}$  satisfying (2.3) has to be continuously differentiable. As was mentioned above, the original form of the conservation law is an integral equation (see (2.1)). In practice also discontinuous solutions occur (cf. e.g. [14], [25], [26]), so the restriction of the solution to be continuously differentiable is too strong. This is the reason why *weak solutions* of the system (2.3) are interesting. These weak solutions are obtained by multiplying (2.3) with an arbitrary test function  $\varphi \in C_0^1(\mathbb{R} \times [0, \infty))$  (thus,  $\varphi(x, 0) = 0$  for all  $x$ ) and, subsequently, partially integrating this equation in space and time. This leads to the following definition.

**Definition 2.3** *The function  $\mathbf{u} \in L_2(\mathbb{R} \times [0, \infty))$  is called a weak solution of the conservation law (2.3) if*

$$\int_0^\infty \int_{-\infty}^{+\infty} \left\{ \mathbf{u}(x, t) \frac{\partial}{\partial t} \varphi(x, t) + \mathbf{f}(\mathbf{u}(x, t)) \frac{\partial}{\partial x} \varphi(x, t) \right\} dx dt = 0, \quad (2.13)$$

for all functions  $\varphi \in C_0^1(\mathbb{R} \times [0, \infty))$ .

From now on by a solution of (2.3) is meant a weak solution of (2.3) in the sense of Definition 2.3. Thus also discontinuous solutions of (2.3) are allowed. Let  $\mathbf{u}$  have a jump discontinuity along a smooth curve  $\Gamma$ , i.e.  $\mathbf{u}$  has well defined limits on both sides of  $\Gamma$ . Let  $\Gamma$  be given by  $x = x(t)$ , then the values  $\mathbf{u}_L = \mathbf{u}(x(t) - 0, t)$  and  $\mathbf{u}_R = \mathbf{u}(x(t) + 0, t)$  are well defined. Not every discontinuity is permissible: in fact the condition (2.13) places severe restrictions on the curve of discontinuity  $\Gamma$ . It is possible to show (cf. [14], [25]) that

$$\bar{s}(\mathbf{u}_L - \mathbf{u}_R) = \mathbf{f}(\mathbf{u}_L) - \mathbf{f}(\mathbf{u}_R) \quad (2.14)$$

must hold at each point on  $\Gamma$ , where  $\bar{s} = x'(t)$  is the speed of the discontinuity. Relation (2.14) is called the *jump condition*; in gas dynamics it is also known as the *Rankine-Hugoniot condition*.

A difficulty is that the weak solutions of (2.3) turn out to be non-unique for a given set of initial data, and it remains to characterise the "physically relevant" weak solution. We therefore remark that (2.3) can be obtained, in the limit for  $\mu \downarrow 0$ , from the equation

$$\frac{\partial}{\partial t} \mathbf{u}_\mu(x, t) + \frac{\partial}{\partial x} \mathbf{f}(\mathbf{u}_\mu(x, t)) = \mu \frac{\partial^2}{\partial x^2} \mathbf{u}_\mu(x, t), \quad (2.15)$$

with  $\mu$  the *viscosity coefficient* ( $\mu > 0$ ). Hence the unique, *physically relevant weak solution* is defined as, roughly speaking, the stable limit of a vanishing viscosity mechanism (cf. [29]). The usual criterion to identify such vanishing viscosity solution is *entropy stability* (cf. e.g. [29]). The idea of entropy stability is to add an extra condition to the solution, such that a physically relevant solution is obtained.

Therefore, consider a twice continuously differentiable function  $\eta : \mathbb{R}^m \rightarrow \mathbb{R}$ . The function  $\eta$  is called *convex* if its Hessian (denoted by  $\eta_{\mathbf{u}\mathbf{u}}$ ) is symmetric positive definite. Thus, for a convex function  $\eta$  the following inequality holds

$$(\eta_{\mathbf{u}\mathbf{u}} \frac{\partial}{\partial x} \mathbf{u})^T \frac{\partial}{\partial x} \mathbf{u} \geq 0, \quad \forall \mathbf{u} \in \mathbb{R}^m \setminus \{0\}. \quad (2.16)$$

**Definition 2.4** A twice continuously differentiable, convex function  $\eta : \mathbb{R}^m \rightarrow \mathbb{R}$  is called an entropy function for the conservation law (2.3), if there exists a continuously differentiable function  $\psi : \mathbb{R}^m \rightarrow \mathbb{R}$ , such that

$$\nabla \psi(\mathbf{u})^T = \nabla \eta(\mathbf{u})^T \frac{\partial}{\partial \mathbf{u}} \mathbf{f}(\mathbf{u}), \quad \forall \mathbf{u} \in \mathbb{R}^m. \quad (2.17)$$

The function  $\psi$  is called an entropy flux.

Here  $\nabla \psi(\mathbf{u})^T = (\frac{\partial}{\partial u_1} \psi(\mathbf{u}), \dots, \frac{\partial}{\partial u_m} \psi(\mathbf{u}))$  and  $\nabla \eta(\mathbf{u})^T$  is defined analogously. A straightforward computation shows that

$$\frac{\partial}{\partial t} \eta(\mathbf{u}(x, t)) + \frac{\partial}{\partial x} \psi(\mathbf{u}(x, t)) = 0 \quad (2.18)$$

holds, if  $\mathbf{u}$  is a continuously differentiable solution of (2.3).

The system of equations (2.17) has two unknowns,  $\eta$  and  $\psi$ . If the system has too many equations, then it may have no solution. If the Jacobian matrix  $A(\mathbf{u})$  is symmetric for all  $\mathbf{u} \in \mathbb{R}^m$ , i.e.  $\partial f_i / \partial u_j = \partial f_j / \partial u_i$ , then there exists a function  $g : \mathbb{R}^m \rightarrow \mathbb{R}$ , such that  $\partial g / \partial u_i = f_i$ . Now it is clear that  $\eta$  and  $\psi$  can be chosen as  $\eta(\mathbf{u}) = \frac{1}{2} \sum_j u_j^2$  and  $\psi(\mathbf{u}) = \sum_j u_j f_j - g(\mathbf{u})$ . In [25] other examples are given for which nontrivial solutions of (2.17) exist.

For the solution of the viscous equation (2.15), we can associate the (small but) positive viscosity term with an entropy inequality. If it is assumed that the solution of the parabolic equation (2.15) is twice continuously differentiable, then equation (2.17) leads to

$$\frac{\partial}{\partial t} \eta(\mathbf{u}_\mu(x, t)) + \frac{\partial}{\partial x} \psi(\mathbf{u}_\mu(x, t)) = \mu \nabla \eta(\mathbf{u}_\mu(x, t))^T \frac{\partial^2}{\partial x^2} \mathbf{u}_\mu(x, t). \quad (2.19)$$

Let  $\varphi \in C_0^1(\mathbb{R} \times [0, \infty))$  be an arbitrary test function such that  $\varphi(x, t) \geq 0$  for all  $x \in \mathbb{R}$  and  $t \in [0, \infty)$ , and assume that the solution  $\mathbf{u}_\mu$  of (2.19) is bounded. Using (2.16) and (2.19) it is easy to see that

$$\begin{aligned} & \int_0^\infty \int_{-\infty}^{+\infty} \left\{ \frac{\partial}{\partial t} \eta(\mathbf{u}_\mu(x, t)) + \frac{\partial}{\partial x} \psi(\mathbf{u}_\mu(x, t)) \right\} \varphi(x, t) dx dt = \\ & \int_0^\infty \int_{-\infty}^{+\infty} \left\{ \mu \nabla \eta(\mathbf{u}_\mu(x, t))^T \frac{\partial^2}{\partial x^2} \mathbf{u}_\mu(x, t) \right\} \varphi(x, t) dx dt = \\ & \mu \int_0^\infty \int_{-\infty}^{+\infty} \left\{ \frac{\partial^2}{\partial x^2} \eta(\mathbf{u}_\mu(x, t)) - (\eta_{\mathbf{u}\mathbf{u}}(\mathbf{u}_\mu) \frac{\partial}{\partial x} \mathbf{u}_\mu(x, t))^T \frac{\partial}{\partial x} \mathbf{u}_\mu(x, t) \right\} \varphi(x, t) dx dt \leq \\ & \mu \int_0^\infty \int_{-\infty}^{+\infty} \frac{\partial^2}{\partial x^2} \eta(\mathbf{u}_\mu(x, t)) \varphi(x, t) dx dt \rightarrow 0, \text{ for } \mu \downarrow 0, \end{aligned}$$

since the latter integral is bounded (cf. [14]). Next an *entropy stable solution* is defined (cf. [29]).

**Definition 2.5** A bounded solution  $\mathbf{u}$  of (2.3) is called an entropy stable solution if, for all convex entropy functions  $\eta$  and corresponding entropy fluxes  $\psi$ , the inequality

$$\frac{\partial}{\partial t} \eta(\mathbf{u}(x, t)) + \frac{\partial}{\partial x} \psi(\mathbf{u}(x, t)) \leq 0 \quad (2.20)$$

is satisfied in the weak sense (for all nonnegative test functions).

In [25] it is proved that (2.20) is equivalent to the condition

$$\bar{s}(\eta(\mathbf{u}_L) - \eta(\mathbf{u}_R)) \leq \psi(\mathbf{u}_L) - \psi(\mathbf{u}_R), \quad (2.21)$$

which holds at a discontinuity of a piecewise continuous solution  $\mathbf{u}$ . Hence this criterion is also often used as the definition of entropy stable solutions. For more details, cf. [14], [25].

Consider the scalar nonlinear conservation law, i.e. (2.3) with  $m = 1$ , be considered. Suppose that an entropy function  $\eta$  and a corresponding entropy flux  $\psi$  are given by

$$\begin{aligned} \eta(u(x, t)) &= |u(x, t) - z|, \\ \psi(u(x, t)) &= \{f(u(x, t)) - f(z)\} \operatorname{sgn}(u(x, t) - z), \end{aligned} \quad (2.22)$$

where  $z$  is an arbitrary real number and  $\operatorname{sgn}(x) = 1$  for  $x > 0$  and  $\operatorname{sgn}(x) = -1$  for  $x < 0$ . It has been shown by Krushkov that the entropy stable solution is unique and is characterized by these choices for the entropy functions and the entropy fluxes. Furthermore, this unique entropy solution is equal to the physically relevant solution (cf. [13]).

Therefore, by analogy with the scalar case, condition (2.20) or (2.21) is often imposed in order to identify the physically relevant solution.

### 3 The Riemann Problem for Hyperbolic Conservation Laws

#### 3.1 Preliminaries

In this section the *Riemann problem* is introduced. The Riemann problem is very important because it forms the underlying physical model of many upwind schemes for hyperbolic conservation laws. For instance, the well-known *Godunov upwind schemes* use the exact solution of the Riemann problem for the numerical solution of hyperbolic conservation laws (cf. [8], [10], [17]).

**Definition 3.1** *The Riemann problem for a general hyperbolic system is the following initial value problem*

$$\frac{\partial}{\partial t} \mathbf{u}(x, t) + \frac{\partial}{\partial x} \mathbf{f}(\mathbf{u}(x, t)) = 0 \quad (3.1)$$

with

$$\mathbf{u}(x, 0) = \begin{cases} \mathbf{u}_L, & x < 0, \\ \mathbf{u}_R, & x > 0, \end{cases} \quad (3.2)$$

where  $\mathbf{u}_L \in \mathbb{R}^m$  and  $\mathbf{u}_R \in \mathbb{R}^m$  are constant states.

Since (3.1) is assumed to be a hyperbolic system, the Jacobian matrix  $A(\mathbf{u})$  has  $m$  real eigenvalues  $\lambda_1(\mathbf{u}), \dots, \lambda_m(\mathbf{u})$  and  $m$  linearly independent right eigenvectors  $\mathbf{r}^{(1)}(\mathbf{u}), \dots, \mathbf{r}^{(m)}(\mathbf{u})$ . For calculating the solution of the Riemann problem the following concepts are introduced (cf. [14], [26]).

**Definition 3.2** *Consider the hyperbolic system (3.1). Let  $\mathbf{r}^{(k)}(\mathbf{u})$  be a right eigenvector of  $A(\mathbf{u})$  and  $\lambda_k(\mathbf{u})$  the corresponding eigenvalue. The eigenvector  $\mathbf{r}^{(k)}(\mathbf{u})$  is called genuinely nonlinear if*

$$(\nabla \lambda_k(\mathbf{u}), \mathbf{r}^{(k)}(\mathbf{u})) \neq 0, \quad \forall \mathbf{u} \in \mathbb{R}^m. \quad (3.3)$$

*The eigenvector  $\mathbf{r}^{(k)}(\mathbf{u})$  is called linearly degenerate if*

$$(\nabla \lambda_k(\mathbf{u}), \mathbf{r}^{(k)}(\mathbf{u})) = 0, \quad \forall \mathbf{u} \in \mathbb{R}^m. \quad (3.4)$$

Here  $\nabla \lambda_k(\mathbf{u}) = (\frac{\partial}{\partial u_1} \lambda_k(\mathbf{u}), \dots, \frac{\partial}{\partial u_m} \lambda_k(\mathbf{u}))^T$  and  $(\cdot, \cdot)$  denotes the usual inner product in  $\mathbb{R}^m$ . It is assumed that  $\mathbf{f}$  is twice continuously differentiable, so that  $\nabla \lambda_k(\mathbf{u})$  exists for all  $k$ . In the following subsection a procedure is given to calculate the solution of the Riemann problem (3.1)-(3.2).

**Example 3.3** In this example a theorem is given, which shows whether the eigenvectors belonging to the Euler equations, given in (2.12), are genuinely nonlinear or linearly degenerate. Its proof is a straightforward calculation (cf. [26]).

**Theorem 3.4** *Let  $\mathbf{r}^{(k)}(\mathbf{u})$ ,  $k = 1, 2, 3$ , be the eigenvectors given in (2.12). Then  $\mathbf{r}^{(k)}(\mathbf{u})$  is linearly degenerate for  $k = 2$  and genuinely nonlinear for  $k = 1$  and  $k = 3$ .*

#### 3.2 Solution of the Riemann problem

The solution of the Riemann problem for a nonlinear hyperbolic system is hard to give in general. But for certain pairs  $(\mathbf{u}_L, \mathbf{u}_R)$  the solution of the Riemann problem can be

derived easily. In [14] it is proved, that if  $\mathbf{u}_L$  and  $\mathbf{u}_R$  are sufficiently close, the initial value problem (3.1)-(3.2) has a unique solution for which an explicit expression can be given. In this subsection it is shown how to derive this solution.

*In the remainder of this section it is assumed that each eigenvector is either linearly degenerate or genuinely nonlinear and that all eigenvalues are distinct.*

Note that this assumption is fulfilled in the case of the Euler equations (see Theorem 3.4). The following theorem describes the general form of solutions of the Riemann problem (cf. [26]).

**Theorem 3.5** *Suppose that there exists a unique solution  $\mathbf{u}$  of the Riemann problem (3.1)-(3.2). Then this solution can be written in the similarity form  $\mathbf{u}(x, t) = \tilde{\mathbf{u}}(x/t)$ .*

**Proof** Define  $\mathbf{u}_\alpha(x, t) = \mathbf{u}(\alpha x, \alpha t)$  with  $\alpha > 0$ . Then it is easily verified that  $\mathbf{u}_\alpha(x, t)$  is also a solution of the Riemann problem. Hence,  $\mathbf{u}(x, t) = \mathbf{u}(\alpha x, \alpha t)$  for all  $\alpha > 0$ , so  $\mathbf{u}(x, t) = \tilde{\mathbf{u}}(x/t)$ .  $\square$

In order to analyse the Riemann problem, the part of a solution associated with a single eigenvector is considered. Here different possibilities exist. If the eigenvector is linearly degenerate, then a *contact discontinuity* appears. In the genuinely nonlinear case there are two possibilities: firstly  $\lambda(\mathbf{u}_L) < \lambda(\mathbf{u}_R)$  in which case a *simple wave* (or *rarefaction wave*) is found, and secondly,  $\lambda(\mathbf{u}_L) > \lambda(\mathbf{u}_R)$  in which case a *shock wave* is found.

To calculate the contact discontinuity or the simple wave solution, the *phase space* is considered. This phase space is simply the  $m$ -dimensional space that contains all values of  $\mathbf{u} = (u_1, u_2, \dots, u_m)^T$ . For each right eigenvector we define a curve in the phase space, such that it starts in some arbitrary given state  $\mathbf{u}_L$  and is tangent to the right eigenvector. It is shown that these curves describe the simple wave and contact discontinuity solution of the Riemann problem.

*Simple wave solution of the Riemann problem.*

Suppose that  $\mathbf{r}^{(k)}(\mathbf{u})$  is a genuinely nonlinear eigenvector and  $\lambda_k(\mathbf{u}_L) < \lambda_k(\mathbf{u}_R)$ . Then  $\mathbf{r}^{(k)}(\mathbf{u})$  can be normalized such that

$$(\nabla \lambda_k(\mathbf{u}), \mathbf{r}^{(k)}(\mathbf{u})) = 1, \quad \forall \mathbf{u} \in \mathbb{R}^m.$$

For an arbitrary state  $\mathbf{u}_L$  and  $k \in \{1, \dots, m\}$ , consider the following initial value problem:

$$\begin{cases} \frac{d}{d\xi} \tilde{\mathbf{u}}(\xi) = \mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi)), \\ \tilde{\mathbf{u}}(0) = \mathbf{u}_L, \end{cases} \quad (3.5)$$

for all  $\xi \in [0, \xi_R]$ . Let  $\xi_R$  be chosen such that  $\mathbf{u}_R = \tilde{\mathbf{u}}(\xi_R)$  is well defined. So  $\tilde{\mathbf{u}}$  describes a curve in the phase space from  $\mathbf{u}_L$  to  $\mathbf{u}_R$ . Since

$$\frac{d}{d\xi} \lambda_k(\tilde{\mathbf{u}}(\xi)) = (\nabla \lambda_k(\tilde{\mathbf{u}}(\xi)), \mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi))) = 1, \quad (3.6)$$

it is obvious that  $\lambda_k(\tilde{\mathbf{u}}(\xi)) = \xi + \lambda_k(\mathbf{u}_L)$  for all  $\xi \in [0, \xi_R]$ . Define the function  $\mathbf{u}$  by

$$\mathbf{u}(x, t) = \begin{cases} \mathbf{u}_L & \text{if } x/t < \lambda_k(\mathbf{u}_L), \\ \tilde{\mathbf{u}}(x/t - \lambda_k(\mathbf{u}_L)) & \text{if } \lambda_k(\mathbf{u}_L) < x/t < \lambda_k(\mathbf{u}_R), \\ \mathbf{u}_R & \text{if } \lambda_k(\mathbf{u}_R) < x/t. \end{cases} \quad (3.7)$$

It will be verified that  $\mathbf{u}(x, t)$  defined in (3.7) is the solution of the Riemann problem. From now on we restrict ourselves to the case  $\lambda_k(\mathbf{u}_L) < x/t < \lambda_k(\mathbf{u}_R)$ . The cases  $x/t < \lambda_k(\mathbf{u}_L)$  or  $x/t > \lambda_k(\mathbf{u}_R)$  are trivial. It is easy to see that

$$\begin{aligned}\lambda_k(\mathbf{u}(x, t)) &= \lambda_k(\tilde{\mathbf{u}}(x/t - \lambda_k(\mathbf{u}_L))) \\ &= x/t - \lambda_k(\mathbf{u}_L) + \lambda_k(\mathbf{u}_L) = x/t.\end{aligned}\tag{3.8}$$

Therefore, using (3.5) with  $\xi = x/t - \lambda_k(\mathbf{u}_L)$  and (3.8), we have

$$\begin{aligned}\frac{\partial}{\partial t}\mathbf{u}(x, t) + \frac{\partial}{\partial x}\mathbf{f}(\mathbf{u}(x, t)) &= \frac{\partial}{\partial t}\mathbf{u}(x, t) + A(\mathbf{u}(x, t))\frac{\partial}{\partial x}\mathbf{u}(x, t) \\ &= -\frac{x}{t^2}\mathbf{r}^{(k)}(\mathbf{u}(x, t)) + \frac{1}{t}A(\mathbf{u}(x, t))\mathbf{r}^{(k)}(\mathbf{u}(x, t)) \\ &= -\frac{x}{t^2}\mathbf{r}^{(k)}(\mathbf{u}(x, t)) + \frac{1}{t}\lambda_k(\mathbf{u}(x, t))\mathbf{r}^{(k)}(\mathbf{u}(x, t)) = 0.\end{aligned}$$

Thus indeed  $\mathbf{u}(x, t)$  defined in (3.7) is the solution of the Riemann problem with initial states  $(\mathbf{u}_L, \mathbf{u}_R)$ . This solution is called a  $k$ th simple wave (or rarefaction wave).

*Shock wave solution of the Riemann problem.*

Another elementary type of solutions of the Riemann problem is given by shock wave solutions. Here a short description of shock wave solutions is given, since a detailed description does not contribute very much to the understanding of the numerical methods that will be discussed. Suppose that  $\mathbf{r}^{(k)}(\mathbf{u})$  is a genuinely nonlinear eigenvector and  $\lambda_k(\mathbf{u}_L) > \lambda_k(\mathbf{u}_R)$ . Recall that for every discontinuity in a solution the jump condition (2.14) must be satisfied. A discontinuity satisfying (2.14) is called a  $k$ th shock wave if

$$\begin{aligned}\lambda_{k-1}(\mathbf{u}_L) &< \bar{s} < \lambda_k(\mathbf{u}_L), \\ \lambda_k(\mathbf{u}_R) &< \bar{s} < \lambda_{k+1}(\mathbf{u}_R),\end{aligned}\tag{3.9}$$

where  $\bar{s}$  is the *speed of the discontinuity* (in e.g. [25] it is proved that shock waves can occur for the Riemann problem). Condition (3.9) (which is sufficient for the solution to be entropy stable, cf. [14]) ensures that nonphysical solutions, such as expansion shocks, do not appear. For a detailed description of shock wave solutions, cf. [14], [25].

*Contact discontinuity solution of the Riemann problem.*

Suppose that  $\mathbf{r}^{(k)}(\mathbf{u})$  is a linearly degenerate eigenvector. Let  $\tilde{\mathbf{u}}(\xi)$  be the solution of (3.5) and suppose that the value  $\xi_R$  is chosen such that  $\mathbf{u}_R = \tilde{\mathbf{u}}(\xi_R)$  is well defined. Since

$$\frac{d}{d\xi}\lambda_k(\tilde{\mathbf{u}}(\xi)) = (\nabla\lambda_k(\tilde{\mathbf{u}}(\xi)), \mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi))) = 0,\tag{3.10}$$

it is obvious that  $\lambda_k(\tilde{\mathbf{u}}(\xi)) = \lambda_k(\mathbf{u}_L) = \lambda_k(\mathbf{u}_R)$  for all  $\xi \in [0, \xi_R]$ . Define the discontinuous function  $\mathbf{u}$  by

$$\mathbf{u}(x, t) = \begin{cases} \mathbf{u}_L & \text{if } x/t < \lambda_k(\mathbf{u}_L) = \lambda_k(\mathbf{u}_R), \\ \mathbf{u}_R & \text{if } x/t > \lambda_k(\mathbf{u}_L) = \lambda_k(\mathbf{u}_R). \end{cases}\tag{3.11}$$

It will be shown that the function  $\mathbf{u}$  defined in (3.11) is the solution of the Riemann problem. For this it suffices to show that the jump condition (2.14) is satisfied. Since

$$\begin{aligned} \frac{d}{d\xi} [\mathbf{f}(\tilde{\mathbf{u}}(\xi)) - \lambda_k(\tilde{\mathbf{u}}(\xi))\tilde{\mathbf{u}}(\xi)] &= A(\tilde{\mathbf{u}}(\xi))\frac{d}{d\xi}\tilde{\mathbf{u}}(\xi) - \lambda_k(\tilde{\mathbf{u}}(\xi))\frac{d}{d\xi}\tilde{\mathbf{u}}(\xi) \\ &= A(\tilde{\mathbf{u}}(\xi))\mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi)) - \lambda_k(\tilde{\mathbf{u}}(\xi))\mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi)) \\ &= \lambda_k(\tilde{\mathbf{u}}(\xi))\mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi)) - \lambda_k(\tilde{\mathbf{u}}(\xi))\mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi)) = 0, \end{aligned}$$

it is easy to see that (2.14) holds with  $\bar{s} = x'(t) = \lambda_k(\mathbf{u}_L) = \lambda_k(\mathbf{u}_R)$ .

Next a theorem is given which describes the total solution of the Riemann problem (for a proof cf. [25]).

**Theorem 3.6** *Let  $\mathbf{u}_L \in \mathbb{R}^m$  be given and suppose that the system (3.1) is hyperbolic with distinct eigenvalues. Further assume that each eigenvector of the Jacobian matrix of  $\mathbf{f}$  is either genuinely nonlinear or linearly degenerate. Then there exists a neighbourhood  $\Omega \subset \mathbb{R}^m$  of  $\mathbf{u}_L$  such that, if  $\mathbf{u}_R \in \Omega$ , the Riemann problem (3.1)-(3.2) has a unique solution. This solution consists of at most  $(m+1)$ -constant states separated by shocks, simple waves or contact discontinuities.*

**Example 3.7** In this example the solution of the Riemann problem for linear systems is considered (cf. [8], [17]) (which will be used in the derivation of Roe's numerical scheme for nonlinear conservation laws, see Subsection 5.4). Therefore, let  $\mathbf{f}(\mathbf{u}) = A\mathbf{u}$ , where  $A$  is a constant  $m \times m$ -matrix. Hence (3.1)-(3.2) simplifies to

$$\frac{\partial}{\partial t}\mathbf{u}(x, t) + A\frac{\partial}{\partial x}\mathbf{u}(x, t) = 0 \quad (3.12)$$

with

$$\mathbf{u}(x, 0) = \begin{cases} \mathbf{u}_L, & x < 0, \\ \mathbf{u}_R, & x > 0, \end{cases} \quad (3.13)$$

where  $\mathbf{u}_L \in \mathbb{R}^m$  and  $\mathbf{u}_R \in \mathbb{R}^m$  are constant states. Since in the linear case all eigenvalues of  $A$  are constant, all eigenvectors are linearly degenerate. Thus, only contact discontinuities appear in the solution. The eigenvectors  $\mathbf{r}^{(k)}$ ,  $k = 1, \dots, m$ , are linearly independent and therefore,  $\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(m)}\}$  can be used as a basis for  $\mathbb{R}^m$ . A solution of (3.12)-(3.13) can be expressed with respect to this basis, i.e.

$$\mathbf{u}(x, t) = \mathbf{u}_L + \sum_{k=1}^m \beta_k(x, t)\mathbf{r}^{(k)},$$

where  $\beta_k : \mathbb{R} \times [0, \infty) \rightarrow \mathbb{R}$ , for all  $k$  with  $1 \leq k \leq m$ . Substitution of this expression in (3.12) leads to

$$\frac{\partial}{\partial t}\mathbf{u}(x, t) + A\frac{\partial}{\partial x}\mathbf{u}(x, t) = \sum_{k=1}^m \left\{ \frac{\partial}{\partial t}\beta_k(x, t) + \lambda_k\frac{\partial}{\partial x}\beta_k(x, t) \right\} \mathbf{r}^{(k)} = 0.$$

Since the eigenvectors are linearly independent, the following equalities should hold

$$\frac{\partial}{\partial t}\beta_k(x, t) + \lambda_k\frac{\partial}{\partial x}\beta_k(x, t) = 0, \quad k = 1, \dots, m.$$

The general solutions of these equations are given by

$$\beta_k(x, t) = \beta_k^0(x - \lambda_k t), \quad k = 1, \dots, m,$$

where  $\beta_k^0 : \mathbb{R} \rightarrow \mathbb{R}$ . Hence, the general solution of (3.12)-(3.13) reads

$$\mathbf{u}(x, t) = \mathbf{u}_L + \sum_{k=1}^m \beta_k^0(x - \lambda_k t) \mathbf{r}^{(k)}. \quad (3.14)$$

Let the initial states be decomposed as  $\mathbf{u}_R - \mathbf{u}_L = \sum_{k=1}^m \alpha_k \mathbf{r}^{(k)}$  and recall the definition of the *Heavyside function*  $H$  by  $H(x) = 1$  if  $x > 0$  and  $H(x) = 0$  if  $x < 0$ . Then,

$$\mathbf{u}(x, 0) = \mathbf{u}_L + H(x)(\mathbf{u}_R - \mathbf{u}_L) = \mathbf{u}_L + \sum_{k=1}^m \alpha_k H(x) \mathbf{r}^{(k)}.$$

Comparing this equation and (3.14), the solution of the Riemann problem (3.12)-(3.13) is obviously,

$$\mathbf{u}(x, t) = \mathbf{u}_L + \sum_{k=1}^m \alpha_k H(x - \lambda_k t) \mathbf{r}^{(k)}. \quad (3.15)$$

### 3.3 Riemann invariants

For the construction of the solution of Riemann problems the so-called *Riemann invariants* are useful, which are defined as follows.

**Definition 3.8** Consider the hyperbolic system (2.3). Let  $\mathbf{r}^{(k)}(\mathbf{u})$  be the  $k$ th right eigenvector of the Jacobian matrix  $A(\mathbf{u})$ . A  $k$ -Riemann invariant is a continuously differentiable function  $w_k : \mathbb{R}^m \rightarrow \mathbb{R}$  such that

$$(\nabla w_k(\mathbf{u}), \mathbf{r}^{(k)}(\mathbf{u})) = 0, \quad \forall \mathbf{u} \in \mathbb{R}^m. \quad (3.16)$$

Note that if  $\mathbf{r}^{(k)}(\mathbf{u})$  is linearly degenerate, then the corresponding eigenvalue  $\lambda_k(\mathbf{u})$  is a Riemann invariant (see (3.4)). In general there are  $m - 1$   $k$ -Riemann invariants whose gradients are linearly independent in  $\mathbb{R}^m$ . For the construction of simple waves or contact discontinuities (3.5) has to be solved. Let  $\mathbf{u}(\xi) = \tilde{\mathbf{u}}(\xi)$ ,  $0 \leq \xi \leq \xi_R$ , be the solution of (3.5). Then

$$\frac{d}{d\xi} w_k(\tilde{\mathbf{u}}(\xi)) = (\nabla w_k(\tilde{\mathbf{u}}(\xi)), \mathbf{r}^{(k)}(\tilde{\mathbf{u}}(\xi))) = 0.$$

Therefore, a  $k$ -Riemann invariant is constant along the curve described by (3.5). If there are  $m - 1$   $k$ -Riemann invariants  $w_k^1, w_k^2, \dots, w_k^{m-1}$ , with linearly independent gradients, then it is easily seen that the curve described by (3.5) is part of the curve described by

$$\left\{ \mathbf{u} \in \mathbb{R}^m \mid w_k^1(\mathbf{u}) = w_k^1(\mathbf{u}_L), w_k^2(\mathbf{u}) = w_k^2(\mathbf{u}_L), \dots, w_k^{m-1}(\mathbf{u}) = w_k^{m-1}(\mathbf{u}_L) \right\}.$$

**Example 3.9** A simple example illustrating the interesting behaviour of the solution of a Riemann problem is the *shock tube problem* of gas dynamics. The physical set-up is a tube filled with gas, initially divided by a membrane in two sections. The density and the pressure of the gas in one part of the tube are larger than in the other part, and the velocity is zero everywhere. At time  $t = 0$ , the membrane is suddenly removed or broken, and the gas flows. It is expected that the gas moves in the direction of lower



pressure. Assuming that the flow is uniform across the tube, there is variation in only one direction and the Riemann problem (3.1) corresponding to the one-dimensional Euler equations (2.7) is relevant. The eigenvectors belonging to these equations are given in (2.12). In the following theorem the Riemann invariants corresponding with the eigenvectors  $\mathbf{r}^{(k)}(\mathbf{u})$  are given.

**Theorem 3.10** *The Riemann invariants  $w_k^1$  and  $w_k^2$  corresponding to the eigenvectors  $\mathbf{r}^{(k)}(\mathbf{u})$  are given by, respectively,*

$$\begin{aligned} w_1^1(\mathbf{u}) &= u + \frac{2}{\gamma-1}c, & w_1^2(\mathbf{u}) &= s, \\ w_2^1(\mathbf{u}) &= u, & w_2^2(\mathbf{u}) &= p, \\ w_3^1(\mathbf{u}) &= u - \frac{2}{\gamma-1}c, & w_3^2(\mathbf{u}) &= s. \end{aligned}$$

The proof is just simple computation (cf. e.g. [25], [26]). Using this theorem the simple wave solution, shock wave solution and the contact discontinuity solution of the Riemann problem for the Euler equations can be derived.

If the pair  $(\mathbf{u}_L, \mathbf{u}_R)$  is such that

$$u_L + \frac{2}{\gamma-1}c_L = u_R + \frac{2}{\gamma-1}c_R; \quad s_L = s_R \quad (3.17)$$

and

$$u_L - c_L < u_R - c_R, \quad (3.18)$$

then a simple wave solution, corresponding to  $\mathbf{r}^{(1)}(\mathbf{u})$  exists and is given by

$$\begin{aligned} \mathbf{u} &= \mathbf{u}_L & \text{if} & \quad x/t < u_L - c_L, \\ \left. \begin{aligned} u + \frac{2}{\gamma-1}c &= u_L + \frac{2}{\gamma-1}c_L \\ s &= s_L \\ u - c &= x/t \end{aligned} \right\} & \text{if} & \quad u_L - c_L < x/t < u_R - c_R, \\ \mathbf{u} &= \mathbf{u}_R & \text{if} & \quad u_R - c_R < x/t. \end{aligned} \quad (3.19)$$

Note that  $u - c = x/t$  follows from (3.8). If  $(\mathbf{u}_L, \mathbf{u}_R)$  is such that (3.17) holds, while  $u_L - c_L > u_R - c_R$ , the solution given by (3.19) corresponds to a multivalued solution, called a *compression wave*. Although such a compression wave has no physical meaning, it is shown in [22] that allowing compression waves, an approximate solution of the Riemann problem can be obtained, which leads to an excellent upwind scheme for the Euler equations.

Similarly, if  $\mathbf{u}_L$  and  $\mathbf{u}_R$  are such that

$$u_L - \frac{2}{\gamma-1}c_L = u_R - \frac{2}{\gamma-1}c_R; \quad s_L = s_R \quad (3.20)$$

and

$$u_L + c_L < u_R + c_R, \quad (3.21)$$

then a simple wave solution, corresponding to  $\mathbf{r}^{(3)}(\mathbf{u})$  exists and is given by

$$\begin{aligned}
 \mathbf{u} &= \mathbf{u}_L && \text{if} && x/t < u_L + c_L, \\
 \left. \begin{aligned} u - \frac{2}{\gamma-1}c &= u_L - \frac{2}{\gamma-1}c_L \\ s &= s_L \\ u + c &= x/t \end{aligned} \right\} && \text{if} && u_L + c_L < x/t < u_R + c_R, \\
 \mathbf{u} &= \mathbf{u}_R && \text{if} && u_R + c_R < x/t.
 \end{aligned} \tag{3.22}$$

If  $u_L + c_L > u_R + c_R$  and (3.20) still hold, then (3.22) generates a compression wave.

If the pair  $(\mathbf{u}_L, \mathbf{u}_R)$  is such that (3.17) holds and  $u_L - c_L > u_R - c_R$ , then a shock wave solution, corresponding to  $\mathbf{r}^{(1)}(\mathbf{u})$  exists and is given by

$$\begin{aligned}
 \mathbf{u} &= \mathbf{u}_L && \text{if} && x/t < \bar{s}, \\
 \mathbf{u} &= \mathbf{u}_R && \text{if} && x/t > \bar{s},
 \end{aligned} \tag{3.23}$$

where  $\bar{s}$  is defined by the Rankine-Hugoniot condition (2.14).

Similarly, if the pair  $(\mathbf{u}_L, \mathbf{u}_R)$  is such that (3.20) and  $u_L + c_L > u_R + c_R$  hold, then a shock wave solution, corresponding to  $\mathbf{r}^{(3)}(\mathbf{u})$  exists, which is again given by (3.23).

Finally, if  $\mathbf{u}_L$  and  $\mathbf{u}_R$  are such that  $u_L = u_R$  and  $p_L = p_R$ , then a contact discontinuity solution, corresponding to  $\mathbf{r}^{(2)}(\mathbf{u})$  exists and is given by

$$\begin{aligned}
 \mathbf{u} &= \mathbf{u}_L && \text{if} && x/t < u_L = u_R, \\
 \mathbf{u} &= \mathbf{u}_R && \text{if} && x/t > u_L = u_R.
 \end{aligned} \tag{3.24}$$

A combination of the equations (3.19)-(3.24) gives the total solution. An example of such a solution is given in Figure 1. The structure of the solution in the  $x$ - $t$  plane is also shown.

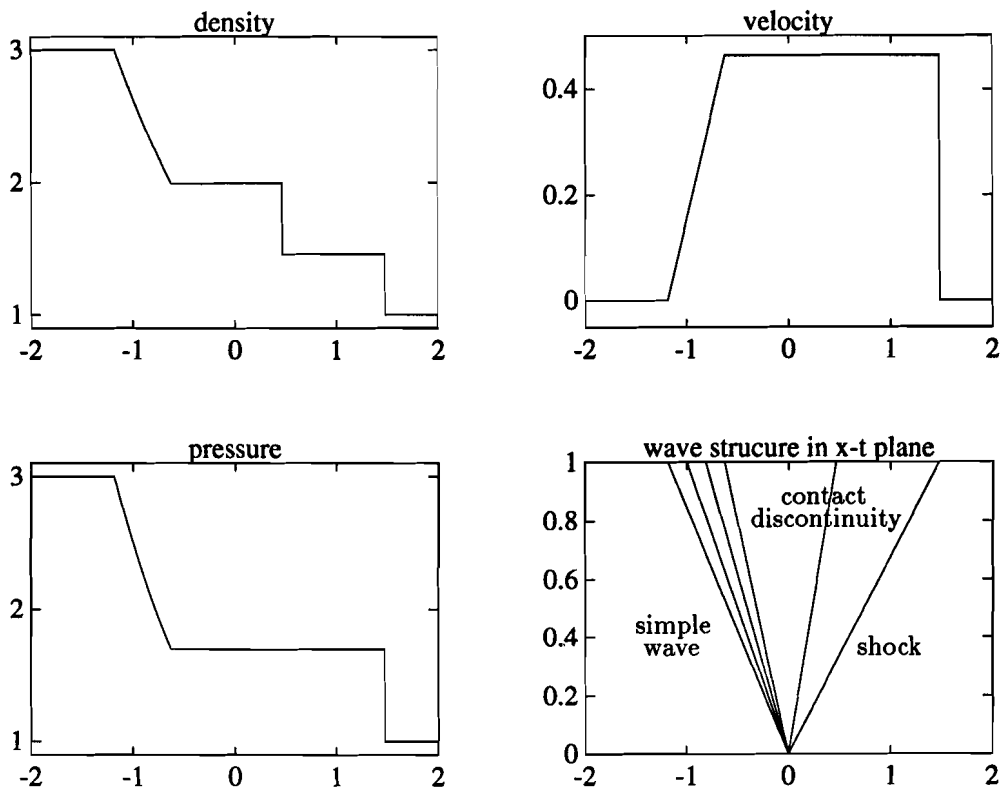


Figure 1: Solution at time  $t = 1$  of a shock tube problem for the one-dimensional Euler equations (2.7) with initial conditions  $p(x, 0) = 3$ ,  $\rho(x, 0) = 3$ ,  $u(x, 0) = 0$  if  $x < 0$  and  $p(x, 0) = 1$ ,  $\rho(x, 0) = 1$ ,  $u(x, 0) = 0$  if  $x > 0$ .

## 4 Introduction to Numerical Methods

### 4.1 Some basic concepts

When solutions of (2.3) are calculated numerically, new problems arise. A finite difference discretization of the conservation law (2.3) is expected to be inappropriate near discontinuities, since it is based on truncated Taylor series expansions. Indeed, if discontinuous solutions of conservation laws are computed using standard finite difference methods, very poor numerical results are obtained (cf. [9],[17]). Later in this section a short description of these standard methods is given, since they are the starting point for more sophisticated methods (cf. [17]).

A mesh is defined in the  $(x, t)$ -space by choosing a mesh width  $\Delta x$  and a time step  $\Delta t$ . For simplicity a uniform mesh is taken. The discrete mesh points  $(x_i, t_n)$  are defined by

$$\begin{aligned} x_i &= i\Delta x, \quad i = \dots, -2, -1, 0, 1, 2, \dots, \\ t_n &= n\Delta t, \quad n = 0, 1, 2, \dots \end{aligned}$$

It will also be useful to define intermediate points

$$x_{i+\frac{1}{2}} = (i + \frac{1}{2})\Delta x.$$

The finite difference methods we shall consider, produce approximations  $U_i^n \in \mathbb{R}^m$  to the true solution  $u(x_i, t_n)$  at the discrete mesh points. The average of  $u(\cdot, t_n)$  on the cell  $[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})$  is defined by

$$\bar{u}_i^n = \frac{1}{\Delta x} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} u(x, t_n) dx. \quad (4.1)$$

For conservation laws it is often convenient to view  $U_i^n$  as an approximation to this average, since the integral form (2.1) of the conservation law describes the evolution in time of integrals such as (4.1). For convenience sake we define a piecewise constant function  $U_{\Delta t}(x, t)$  for all  $x$  and  $t$  from the discrete values  $U_i^n$  by (cf. [17])

$$U_{\Delta t}(x, t) = U_i^n, \quad (x, t) \in [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}) \times [t_n, t_{n+1}). \quad (4.2)$$

In the following it is assumed that the mesh width  $\Delta x$  and time step  $\Delta t$  are related by

$$\frac{\Delta t}{\Delta x} = \tau, \quad (4.3)$$

with  $\tau > 0$  a given constant. Hence, the choice of  $\Delta t$  defines a unique mesh.

Many difficulties for a numerical method are caused by the fact that a discontinuous solution of (2.3) can occur. It is not surprising that the method might converge to the wrong solution, since in general a weak solution is not unique. Therefore the discrete solution of this problem is often required to satisfy a discrete form of entropy stability, as defined in Definition 2.5. More surprisingly a method may converge to a function that is not a weak solution at all. The latter problem is avoided by considering *conservative methods* only, which are consistent with the conservation law (2.3).

**Definition 4.1** *Let a  $(2k + 1)$ -point finite difference method, with 2 time levels, for the hyperbolic conservation law (2.3) be given. The numerical method is said to be conservative, if the corresponding scheme can be written as*

$$U_i^{n+1} = U_i^n - \tau(F_{i+\frac{1}{2}} - F_{i-\frac{1}{2}}), \quad (4.4)$$

where  $\mathbf{F}$  is a function of the values of  $\mathbf{U}$  at  $2k$  points, i.e.

$$\mathbf{F}_{i+\frac{1}{2}} = \mathbf{F}(\mathbf{U}_{i+k}^n, \dots, \mathbf{U}_{i-k+1}^n). \quad (4.5)$$

$\mathbf{F}$  is called the *numerical flux (function)*.

Another important concept is the *local discretization error*. The local discretization error  $\mathbf{D}_{\Delta t}(x, t)$  is a measure how well the difference equation approximates the differential equation locally (at the point  $(x, t)$ ). Let the conservative,  $(2k+1)$ -point scheme (4.4) be written as

$$\mathbf{U}_i^{n+1} = \mathcal{H}_{\Delta t}(\mathbf{U}_{i+k}^n, \dots, \mathbf{U}_{i-k}^n), \quad (4.6)$$

where  $\mathcal{H}_{\Delta t} : (\mathbb{R}^m)^{2k+1} \rightarrow \mathbb{R}^m$  is a finite difference operator. Equation (4.6) can be generalized for arbitrary  $x$  and  $t$  by replacing the numerical solution by the piecewise constant function  $\mathbf{U}_{\Delta t}$  defined in (4.2):

$$\mathbf{U}_{\Delta t}(x, t + \Delta t) = \mathcal{H}_{\Delta t}(\mathbf{U}_{\Delta t}(x + k\Delta x, t), \dots, \mathbf{U}_{\Delta t}(x - k\Delta x, t)). \quad (4.7)$$

If now each  $\mathbf{U}_{\Delta t}(x, t)$  in (4.7) is replaced by the exact solution of (2.3) at the corresponding point, then in general the equality will not hold exactly. This leads to the following definition.

**Definition 4.2** Consider a conservative,  $(2k+1)$ -point method written in the generic form (4.7) for arbitrary  $x$  and  $t$ . The local discretization error  $\mathbf{D}_{\Delta t}$  of this method at the point  $(x, t)$  is defined by

$$\mathbf{D}_{\Delta t}(x, t) = \frac{1}{\Delta t} \left\{ \mathbf{u}(x, t + \Delta t) - \mathcal{H}_{\Delta t}(\mathbf{u}(x + k\Delta x, t), \dots, \mathbf{u}(x - k\Delta x, t)) \right\}, \quad (4.8)$$

where  $\mathbf{u}$  is the solution of (2.3).

Often  $\mathbf{u}$  is assumed to be a smooth solution of (2.3), since Taylor series expansions are used to calculate the local discretization error.

Using the local discretization error we can define the concept of *consistency* (cf. [17]).

**Definition 4.3** Consider a conservative,  $(2k+1)$ -point method. The numerical method is called consistent of order  $p$  with the conservation law (2.3), if the local discretization error satisfies the equality

$$\mathbf{D}_{\Delta t}(x, t) = \mathcal{O}(\Delta t^p),$$

with  $p > 0$ . The method is called consistent with the conservation law (2.3), if it is consistent of order  $p$ .

It can be shown (cf. [17]) that for consistency of a conservative method it is sufficient to require the flux-function  $\mathbf{F}$  of the corresponding scheme (4.4) to be Lipschitz continuous and to satisfy

$$\mathbf{F}(\mathbf{u}, \dots, \mathbf{u}) = \mathbf{f}(\mathbf{u}).$$

The final concepts we introduce in this subsection are *convergence* and *stability*. For simplicity it is assumed that the numerical method is linear, i.e.  $\mathcal{H}_{\Delta t}$  is a linear difference operator. First the *global discretization error*  $\mathbf{E}_{\Delta t}(x, t)$  of a conservative,  $(2k+1)$ -point method is defined for arbitrary  $x$  and  $t$  as

$$\mathbf{E}_{\Delta t}(x, t) = \mathbf{U}_{\Delta t}(x, t) - \mathbf{u}(x, t). \quad (4.9)$$

With this definition, convergence of a numerical method is defined (cf. [17]).

**Definition 4.4** Consider a conservative,  $(2k + 1)$ -point method. The method is called convergent in some particular norm  $\|\cdot\|$ , if

$$\|\mathbf{E}_{\Delta t}(\cdot, t)\| \rightarrow 0, \text{ as } \Delta t \rightarrow 0,$$

for any fixed  $t \geq 0$  and for all initial data  $\mathbf{u}_0$  with  $\|\mathbf{u}_0\|$  finite.

Note that (4.8) can be rewritten in the form

$$\mathbf{u}(x, t + \Delta t) = \mathcal{H}_{\Delta t}(\mathbf{u}(x + k\Delta x, t), \dots, \mathbf{u}(x - k\Delta x, t)) + \Delta t \mathbf{D}_{\Delta t}(x, t).$$

Since the numerical solution satisfies (4.7), after subtracting these two equations a simple recurrence relation for the global discretization error  $\mathbf{E}_{\Delta t}$  is obtained,

$$\mathbf{E}_{\Delta t}(x, t + \Delta t) = \mathcal{H}_{\Delta t}(\mathbf{E}_{\Delta t}(x + k\Delta x, t), \dots, \mathbf{E}_{\Delta t}(x - k\Delta x, t)) - \Delta t \mathbf{D}_{\Delta t}(x, t).$$

Note that linearity is essential here. The latter equation can be rewritten in the functional form

$$\mathbf{E}_{\Delta t}(\cdot, t + \Delta t) = \mathcal{H}_{\Delta t} \mathbf{E}_{\Delta t}(\cdot, t) - \Delta t \mathbf{D}_{\Delta t}(\cdot, t).$$

The global error  $\mathbf{E}_{\Delta t}$  at time  $t + \Delta t$  consists of two parts. One is the new local error  $-\Delta t \mathbf{D}_{\Delta t}$  introduced in this time step. The other part is the cumulative error from previous time steps. By applying this relation recursively we obtain an expression for the global error at time  $t_n$ :

$$\mathbf{E}_{\Delta t}(\cdot, t_n) = \mathcal{H}_{\Delta t}^n \mathbf{E}_{\Delta t}(\cdot, 0) - \Delta t \sum_{j=1}^n \mathcal{H}_{\Delta t}^{n-j} \mathbf{D}_{\Delta t}(\cdot, t_{j-1}). \quad (4.10)$$

Here superscripts on  $\mathcal{H}_{\Delta t}$  represent powers of the linear operator obtained by repeated applications. In order to obtain a bound on the global error, we must insure that the local error  $\mathbf{D}_{\Delta t}(\cdot, t_{j-1})$  is not unduly amplified by applying  $n - j$  steps of the method. Note that a bound is always with respect to some given norm  $\|\cdot\|$ . Next the concept of stability is introduced (cf. [17]).

**Definition 4.5** Consider a conservative,  $(2k + 1)$ -point method written in the generic form (4.7) for arbitrary  $x$  and  $t$ . The numerical method is called stable in some particular norm  $\|\cdot\|$ , if for each time  $T > 0$  there exists a constant  $C$  and a value  $k_0$  such that

$$\|\mathcal{H}_{\Delta t}^n\| \leq C, \text{ for all } n\Delta t \leq T, \Delta t < k_0$$

holds.

In practice, instead of Definition 4.5, often the Von Neumann method for stability analysis is used (cf. [9]). This method gives necessary conditions for a numerical method to be stable. Unfortunately, these conditions are not always sufficient for stability. We return to these concepts later on.

In the remainder of this report only conservative numerical methods, which are consistent with the conservation law (2.3), are considered.

## 4.2 Examples of conservative methods

In this subsection short descriptions are given of four well-known standard finite difference methods. Two first order methods, namely Lax-Friedrichs and basic upwind, and two second order methods, namely Lax-Wendroff and Beam-Warming, are given. Numerical results for the methods are presented in Figure 2. All methods mentioned here, reproduce discontinuous solutions very poorly. The methods are considered for the scalar, linear convection equation

$$\frac{\partial}{\partial t}u(x, t) + a\frac{\partial}{\partial x}u(x, t) = 0. \quad (4.11)$$

Hence, the flux-function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is given by  $f(u) = au$ . It is assumed that the solution  $u$  is three times continuously differentiable. It is useful to introduce the *Courant number*, which is defined by

$$\sigma = a\tau = a\frac{\Delta t}{\Delta x}. \quad (4.12)$$

**Example 4.6 (Lax-Friedrichs)** If in (4.11)  $\partial u/\partial t$  is replaced by a forward difference approximation and  $\partial u/\partial x$  is replaced by a central difference approximation, then the following difference scheme is obtained,

$$U_i^{n+1} = U_i^n - \frac{1}{2}\sigma(U_{i+1}^n - U_{i-1}^n).$$

However, this scheme is unconditionally unstable (cf. [9]). If  $U_i^n$  is replaced by the average  $\frac{1}{2}(U_{i-1}^n + U_{i+1}^n)$ , then the *Lax-Friedrichs method* is obtained:

$$U_i^{n+1} = \frac{1}{2}(U_{i-1}^n + U_{i+1}^n) - \frac{1}{2}\sigma(U_{i+1}^n - U_{i-1}^n). \quad (4.13)$$

This method is stable in the  $L_1$ -norm under the CFL-condition (cf. [17])

$$|\sigma| \leq 1. \quad (4.14)$$

If the numerical flux-function  $F$  is defined by

$$F_{i+\frac{1}{2}} = \frac{1}{2\tau}(U_i^n - U_{i+1}^n) + \frac{1}{2}a(U_i^n + U_{i+1}^n),$$

then it is easy to see that the Lax-Friedrichs method is a conservative method, which is consistent with the conservation law (4.11). Using Definition 4.2 and Taylor series expansions it follows that the local discretization error at the point  $(x, t)$  is given by

$$\begin{aligned} D_{\Delta t}(x, t) &= \frac{1}{2}\left(\Delta t\frac{\partial^2}{\partial t^2}u(x, t) - \frac{\Delta x^2}{\Delta t}\frac{\partial^2}{\partial x^2}u(x, t)\right) + \mathcal{O}(\Delta x^2) \\ &= \frac{\Delta x}{2\tau}(\sigma^2 - 1)\frac{\partial^2}{\partial x^2}u(x, t) + \mathcal{O}(\Delta x^2). \end{aligned} \quad (4.15)$$

**Example 4.7 (Basic upwind)** Upwind methods depend on the stream direction of the fluid. If in (4.11)  $a > 0$ , then the information is propagating in the positive  $x$ -direction. Thus the information in, say, point  $x_i$  has reached point  $x_{i-1}$  before. Therefore, in this case, it is meaningful to replace  $\partial u/\partial x$  by a backward difference.

Similarly  $\partial u / \partial x$  is replaced by a forward difference if  $a < 0$ . In both cases  $\partial u / \partial t$  is replaced by a forward difference. Thus the *basic upwind method* is given by the following difference schemes,

$$\begin{aligned} U_i^{n+1} &= U_i^n - \sigma(U_i^n - U_{i-1}^n) \quad \text{if } a > 0, \\ U_i^{n+1} &= U_i^n - \sigma(U_{i+1}^n - U_i^n) \quad \text{if } a < 0. \end{aligned} \quad (4.16)$$

If we define

$$\begin{aligned} a^+ &= \max(a, 0) \geq 0, \\ a^- &= \min(a, 0) \leq 0, \end{aligned}$$

then equation (4.16) can be rewritten as

$$U_i^{n+1} = U_i^n - \tau \{ a^+(U_i^n - U_{i-1}^n) + a^-(U_{i+1}^n - U_i^n) \}. \quad (4.17)$$

This scheme is also stable in the  $L_1$ -norm under the CFL condition (4.14) (cf. [9]). If the numerical flux-function  $F$  is defined by

$$F_{i+\frac{1}{2}} = a^+ U_i^n + a^- U_{i+1}^n,$$

then it follows immediately that the basic upwind method (4.17) is a conservative method, which is consistent with the conservation law (4.11).

The two second order methods presented below are based on the Taylor series expansion

$$u(x, t + \Delta t) = u(x, t) + \Delta t \frac{\partial}{\partial t} u(x, t) + \frac{1}{2} \Delta t^2 \frac{\partial^2}{\partial t^2} u(x, t) + \mathcal{O}(\Delta t^3). \quad (4.18)$$

Since it is assumed that  $u$  is three times continuously differentiable, it follows, from  $\partial u / \partial t = -a \partial u / \partial x$ , that

$$\frac{\partial^2}{\partial t^2} u(x, t) = -a \frac{\partial^2}{\partial t \partial x} u(x, t) = -a \frac{\partial^2}{\partial x \partial t} u(x, t) = a^2 \frac{\partial^2}{\partial x^2} u(x, t).$$

Using this equality, equation (4.18) becomes

$$u(x, t + \Delta t) = u(x, t) - \Delta t a \frac{\partial}{\partial x} u(x, t) + \frac{1}{2} \Delta t^2 a^2 \frac{\partial^2}{\partial x^2} u(x, t) + \mathcal{O}(\Delta t^3). \quad (4.19)$$

**Example 4.8 (Lax-Wendroff)** The *Lax-Wendroff method* results from retaining only the first three terms of (4.19) and using central difference approximations for the derivatives appearing there. Therefore, the corresponding finite difference scheme is

$$U_i^{n+1} = U_i^n - \frac{1}{2} \sigma (U_{i+1}^n - U_{i-1}^n) + \frac{1}{2} \sigma^2 (U_{i+1}^n - 2U_i^n + U_{i-1}^n). \quad (4.20)$$

The Lax-Wendroff scheme is stable under the CFL condition (4.14) (cf. [9]). If the numerical flux-function  $F$  is defined by

$$F_{i+\frac{1}{2}} = \frac{1}{2} a (U_{i+1}^n + U_i^n) - \frac{1}{2} \tau a^2 (U_{i+1}^n - U_i^n),$$

then it is obvious that the Lax-Wendroff method is a conservative method, which is consistent with the conservation law (4.11).



**Example 4.9 (Beam-Warming)** The *Beam-Warming method* is a one-sided version of the Lax-Wendroff method. It is also obtained from (4.19), but now using second order accurate one-sided approximations of the derivatives in (4.19). If  $a > 0$ , then the corresponding scheme is given by

$$U_i^{n+1} = U_i^n - \frac{1}{2}\sigma(3U_i^n - 4U_{i-1}^n + U_{i-2}^n) + \frac{1}{2}\sigma^2(U_i^n - 2U_{i-1}^n + U_{i-2}^n). \quad (4.21)$$

This scheme is stable under the condition (cf. [9])

$$0 \leq \sigma \leq 2.$$

With the numerical flux-function defined by

$$F_{i+\frac{1}{2}} = \frac{1}{2}a(3U_i^n - U_{i-1}^n) + \frac{1}{2}a\sigma(U_i^n - U_{i-1}^n),$$

the Beam-Warming method is a conservative method, which is consistent with the conservation law (4.11).

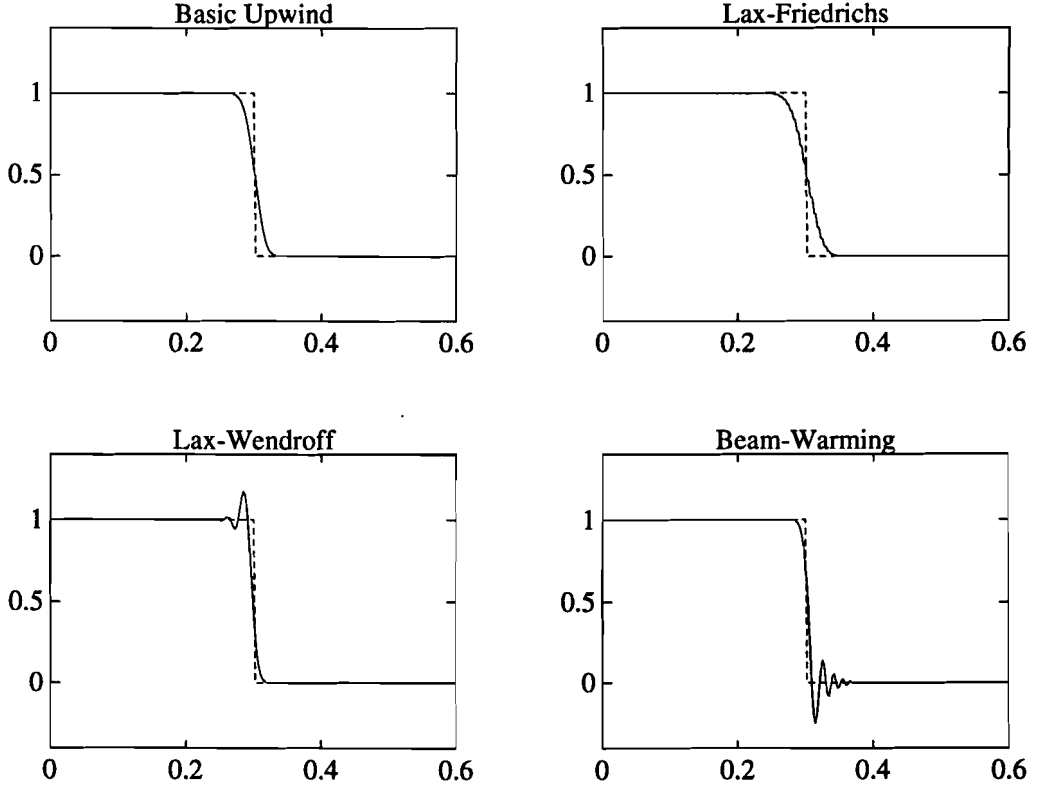


Figure 2: Numerical solution (solid line) and exact solution (dashed line) of (4.11) at  $t = 0.3$  with  $a = 1, \Delta t = 0.002, \sigma = 0.8$  and the initial condition  $u(x, 0) = 1$  if  $x < 0$  and  $u(x, 0) = 0$  if  $x > 0$ .

### 4.3 Modified equations

A useful technique for studying the behaviour of numerical solutions is to model the difference equation by a differential equation. Of course the difference equation was originally derived by approximating (2.3), but there are differential equations that are more accurately approximated than the original; these are the so-called *modified equations* (cf. [17]).

The modified equation is derived by a two-step procedure (cf. [31]). For the sake of simplicity, we describe this procedure for the scalar convection equation (4.11) only. In this analysis it is assumed that there exists a smooth function  $U = U(x, t)$ , which is an exact solution of a given finite difference scheme (4.7). The first step is to expand each term of the given finite difference scheme in a Taylor series expansion around  $U(x, t)$ . Substituting the Taylor series expansions in the scheme gives a partial differential equation which includes an infinite number of both space and time derivatives.

In the second step, all time derivatives appearing in the derived equation, are eliminated, with the exception of the  $\partial/\partial t$  term. The equation obtained after the second step, is called the modified equation. The two-step procedure is illustrated for the Lax-Wendroff scheme.

**Example 4.10** For the Lax-Wendroff scheme given in (4.20), the first step gives the differential equation (using  $\Delta t = \sigma \Delta x/a$ )

$$\begin{aligned} \frac{\partial}{\partial t}U(x, t) + a \frac{\partial}{\partial x}U(x, t) + \frac{\sigma}{2a} \Delta x \frac{\partial^2}{\partial t^2}U(x, t) - \frac{a\sigma}{2} \Delta x \frac{\partial^2}{\partial x^2}U(x, t) \\ + \frac{\sigma^2}{6a^2} \Delta x^2 \frac{\partial^3}{\partial t^3}U(x, t) + \frac{a}{6} \Delta x^2 \frac{\partial^3}{\partial x^3}U(x, t) + \frac{\sigma^3}{24a^3} \Delta x^3 \frac{\partial^4}{\partial t^4}U(x, t) \\ - \frac{a\sigma}{24} \Delta x^3 \frac{\partial^4}{\partial x^4}U(x, t) + \dots = 0. \end{aligned} \quad (4.22)$$

In the second step we have to eliminate all time derivatives appearing in (4.22). Suppose that we want to eliminate, for example, the term  $\partial^2 U/\partial t^2$  in this equation. Therefore, the operator  $-(\Delta t/2)\partial/\partial t$  is applied to (4.22), and the result is added to (4.22). The resulting new equation has a term  $-(\sigma \Delta x/2)\partial^2 U/\partial t \partial x$  which, in turn, can be eliminated by applying the operator  $(\sigma \Delta x/2)\partial/\partial x$  to equation (4.22) and adding the result to the new equation. Similarly, the other time derivatives appearing in (4.22) can be removed. Finally, the following equation is obtained

$$\frac{\partial}{\partial t}U(x, t) + a \frac{\partial}{\partial x}U(x, t) = \frac{a}{6}(\sigma^2 - 1) \Delta x^2 \frac{\partial^3}{\partial x^3}U(x, t) + \mathcal{O}(\Delta x^3). \quad (4.23)$$

Equation (4.23) is called the modified equation of the Lax-Wendroff method (cf. [17]).

In general, for a given finite difference scheme corresponding with (4.11), the procedure described above provides the modified equation

$$\frac{\partial}{\partial t}U(x, t) + a \frac{\partial}{\partial x}U(x, t) = \sum_{q=p+1}^{\infty} \mu(q) \Delta x^{q-1} \frac{\partial^q}{\partial x^q}U(x, t), \quad (4.24)$$

for a method of order  $p$ . The coefficients  $\mu(q)$  appearing in the sum denote the coefficients of the  $q$ th spatial derivatives. These derivatives do not occur in the original partial differential equation and constitute a form of discretization error introduced by the finite difference method.

**Example 4.11** The modified equation of the Lax-Friedrichs method (4.13) for the scalar convection equation (4.11) is given by

$$\frac{\partial}{\partial t}U(x,t) + a\frac{\partial}{\partial x}U(x,t) = \frac{1}{2\tau}(1-\sigma^2)\Delta x\frac{\partial^2}{\partial x^2}U(x,t) + \mathcal{O}(\Delta x^2).$$

The modified equation of the basic upwind method (4.17) for the scalar convection equation (4.11) with  $a > 0$  is given by (cf. [9])

$$\frac{\partial}{\partial t}U(x,t) + a\frac{\partial}{\partial x}U(x,t) = \frac{a}{2}(1-\sigma)\Delta x\frac{\partial^2}{\partial x^2}U(x,t) + \mathcal{O}(\Delta x^2).$$

For a first order scheme like Lax-Friedrichs or basic upwind, the modified equation is a *convection-diffusion equation* of the form

$$\frac{\partial}{\partial t}U(x,t) + a\frac{\partial}{\partial x}U(x,t) = \mu\Delta x\frac{\partial^2}{\partial x^2}U(x,t) + \mathcal{O}(\Delta x^2), \quad (4.25)$$

with a diffusion coefficient  $\mu\Delta x$ . The quantity  $\mu\Delta x\partial^2 U(x,t)/\partial x^2$  is called the *numerical diffusion* of the scheme. To study the behaviour of the numerical solution of these two methods, the solution  $U$  is developed in a Fourier series. Since linear schemes with constant coefficients are considered, it is sufficient to consider only a single Fourier mode of this series

$$U^k(x,t) = e^{ik(x-\tilde{a}t)}, \quad (4.26)$$

where  $k$  is called the *wave-number* and  $\tilde{a}$  is called the *numerical wave-speed*. Note that if  $k$  is large, then (4.26) is a highly oscillatory Fourier mode. These highly oscillatory Fourier modes appear around discontinuities. Furthermore, if  $U^k$  satisfies (4.11), then  $\tilde{a} = a$ . Substitution of the Fourier mode in the modified equation (4.25) gives

$$\tilde{a} = a - \mu\Delta x ik.$$

Note that  $\mu > 0$  is a necessary condition for stability. If  $\mu > 0$  (i.e. if  $|\sigma| < 1$ ), then it is seen that especially highly oscillatory Fourier modes at  $t = 0$  are damped as time evolves. Hence, it is expected that the solution of the modified equation is smeared out as time evolves (see Figure 2). This indicates why the Lax-Friedrichs and basic upwind method approximate discontinuities in solutions too smooth. In general, first order methods have the disadvantage to smear out the solution around discontinuities (cf. [9], [17]).

**Example 4.12** The modified equation of the Beam-Warming method (4.21) for the scalar convection equation (4.11) is given by

$$\frac{\partial}{\partial t}U(x,t) + a\frac{\partial}{\partial x}U(x,t) = \frac{a}{6}(\sigma-1)(\sigma-2)\Delta x^2\frac{\partial^3}{\partial x^3}U(x,t) + \mathcal{O}(\Delta x^3).$$

The modified equation for both the Lax-Wendroff method (see (4.23)) and also the Beam-Warming method is a *dispersive equation* of the form

$$\frac{\partial}{\partial t}U(x,t) + a\frac{\partial}{\partial x}U(x,t) = \mu\Delta x^2\frac{\partial^3}{\partial x^3}U(x,t) + \mathcal{O}(\Delta x^3). \quad (4.27)$$

The quantity  $\mu\Delta x^2\partial^3 U(x,t)/\partial x^3$  is called the *numerical dispersion* of the scheme. Using the same arguments as in the case of the convection-diffusion equation, again a

single Fourier mode (4.26) is considered. If this mode is substituted in the modified equation (4.27) it is seen that

$$\tilde{a} = a + \mu \Delta x^2 k^2.$$

Suppose that  $a > 0$ . If  $\mu < 0$  (i.e. if  $|\sigma| < 1$  for the Lax-Wendroff method and if  $1 < \sigma < 2$  for the Beam-Warming method), then  $\tilde{a} < a$ . Hence, the (highly oscillatory) Fourier modes propagate with a numerical velocity less than the exact velocity  $a$ . Thus, oscillations occur behind the discontinuity (see Figure 2). If  $\mu > 0$  (i.e. if  $0 < \sigma < 1$  for the Beam-Warming method), then  $\tilde{a} > a$ , and the (highly oscillatory) Fourier modes travel too fast. Thus, the oscillations are ahead of the discontinuity (see Figure 2). In general second order methods produce oscillations around discontinuities.

#### 4.4 Numerical entropy stability

Suppose that a conservative numerical method, which is consistent with a conservation law, converges to some function  $\mathbf{u}$ . In this subsection we are looking at some conditions which guarantee that the limit function  $\mathbf{u}$  is an entropy stable weak solution of the conservation law (see Definition 2.5).

Therefore, a numerical variant of Definition 2.5 is given (cf. [29], [30]).

**Definition 4.13** *A conservative numerical method is called an entropy stable method if, for all convex entropy functions  $\eta$  and corresponding entropy fluxes  $\psi$ , the inequality*

$$\eta(\mathbf{U}_i^{n+1}) \leq \eta(\mathbf{U}_i^n) - \tau(\Psi_{i+\frac{1}{2}} - \Psi_{i-\frac{1}{2}}) \quad (4.28)$$

*is satisfied. Here  $\Psi$  is a function of the values of  $\mathbf{U}$  at  $2k$  points, i.e.*

$$\Psi_{i+\frac{1}{2}} = \Psi(\mathbf{U}_{i+k}^n, \dots, \mathbf{U}_{i-k+1}^n). \quad (4.29)$$

$\Psi$  is called the numerical entropy flux. It is assumed that the numerical entropy flux is consistent with the entropy flux, i.e.  $\Psi$  is Lipschitz continuous and

$$\Psi(\mathbf{u}, \dots, \mathbf{u}) = \psi(\mathbf{u}). \quad (4.30)$$

The following theorem shows the importance of the concepts that are introduced in this section. It is a simple extension of the well known theorem of Lax and Wendroff (cf. [17]) and is proved in [8].

**Theorem 4.14** *Suppose that the conservative method with difference scheme (4.4) is consistent with the conservation law (2.3). Furthermore, assume that the method is entropy stable. Let  $\mathbf{U}_i^n$  be a solution of (4.4) with given initial values  $\mathbf{U}_i^0 = \bar{\mathbf{u}}_i^0$ , as defined in (4.1). Define the piecewise constant function  $\mathbf{U}_{\Delta t}$  as in (4.2). Suppose that for some sequence  $\Delta t \rightarrow 0$  the limit*

$$\lim_{\Delta t \rightarrow 0} \mathbf{U}_{\Delta t}(x, t) = \mathbf{u}(x, t)$$

*exists in the sense of bounded,  $L_1^{\text{loc}}$  convergence (i.e.  $\mathbf{U}_{\Delta t}(\cdot, t)$  converges towards  $\mathbf{u}$  in  $L_1^{\text{loc}}$  as  $\Delta t \rightarrow 0$ , for any fixed  $t \geq 0$ , and  $\mathbf{U}_{\Delta t}(\cdot, t)$  is bounded in  $L_\infty$ ). Then the limit  $\mathbf{u}$  is an entropy stable solution of (2.3).*

If a certain amount of numerical diffusion is added to the numerical method, then entropy stable methods are obtained, at the cost of smearing out the physical discontinuities (cf. [19], [29]).

In the scalar case, there exists an easier requirement for a numerical method to converge to the entropy stable solution (cf. [20]).

**Definition 4.15** *Consider a conservative,  $(2k + 1)$ -point finite difference method with 2 time levels, which is consistent with the scalar conservation law (2.3). If the corresponding numerical flux  $F_{i+1/2}$  of the method satisfies*

$$\operatorname{sgn}(U_{i+1} - U_i)(F_{i+1/2} - f(u)) \leq 0, \quad (4.31)$$

*for all  $u$  between  $U_i$  and  $U_{i+1}$ , then the method is called an E-method.*

The following theorem is proved by Osher in [20] and clarifies why E-methods are useful.

**Theorem 4.16** *Suppose that the conservative difference scheme (4.4) is consistent with the conservation law (2.3). If the method is an E-method, then the method is convergent in the sense of bounded,  $L_1^{\text{loc}}$  convergence and its limit is the unique entropy stable solution of the scalar conservation law (2.3).*

E-methods have the following important disadvantage (cf. [20]).

**Theorem 4.17** *An E-method is consistent of at most order one.*

## 5 Godunov-type methods

### 5.1 Introduction

An important class of numerical methods for hyperbolic conservation laws are the *Godunov-type methods*. In Godunov-type methods, the numerical solution is considered piecewise constant in each mesh cell  $[x_{i-1/2}, x_{i+1/2})$  at a certain time level  $t_n = n\Delta t$ . The evolution of the solution to the next time level  $t_{n+1}$  results from the wave interactions originating at the boundaries between adjacent cells. The cell interface at  $x_{i+1/2}$  separates two constant states  $U_i$  at the left and  $U_{i+1}$  at the right side, thus the resulting local interaction can be resolved exactly, since the initial conditions at the time  $t_n$  correspond to the Riemann problem (3.1)-(3.2). As was shown in Subsection 3.2, this problem has an exact solution consisting of constant states separated by shocks, contact discontinuities or simple waves (see Theorem 3.6). The new piecewise constant approximation at time  $t_{n+1}$  is then obtained by averaging over each cell, the exact solution of the Riemann problem.

However, the computational costs to obtain this exact solution are high in general (cf. [26]). Therefore, *approximate Riemann solutions* are considered in order to reduce the computational work. The *approximate Riemann solvers* to be described in this section have been developed by Osher (cf. [22]) and Roe (cf. [24]).

Since we shall apply the theory of Section 3, it is assumed that each eigenvector is either linearly degenerate or genuinely nonlinear and that all eigenvalues are distinct. Only conservative methods are considered. Such methods are completely determined by their numerical flux-function (see (4.4)). Therefore, we restrict ourselves to the computation of the numerical flux-function for all the considered methods.

### 5.2 The basic Godunov method

Three steps are involved in the *basic Godunov method* in order to calculate the numerical solution at time level  $t_{n+1}$  from the known numerical solution at time level  $t_n$  (cf. [10]).

In the first step the numerical solution  $U_i^n$  is used to define a piecewise constant function  $\tilde{U}_i^n$  by

$$\tilde{U}_i^n(x, t_n) = U_i^n, \quad x \in [x_{i-1/2}, x_{i+1/2}). \quad (5.1)$$

At time  $t_n$  this function is equal to the piecewise constant function  $U_{\Delta t}$  that has already been introduced (see (4.2)). Unlike  $U_{\Delta t}$ , the function  $\tilde{U}_i^n$  will not be constant for  $t_n \leq t < t_{n+1}$ . Because of the piecewise constant approximation of  $u$ , the Godunov method is first order accurate in space.

In the second step the solution of the local Riemann problem at the cell interfaces is computed. This Riemann problem is given by (see (3.1)-(3.2))

$$\frac{\partial}{\partial t} \tilde{U}_i^n(x, t) + \frac{\partial}{\partial x} f(\tilde{U}_i^n(x, t)) = 0 \quad (5.2)$$

with

$$\tilde{U}_i^n(x, t_n) = \begin{cases} U_i^n, & x < x_{i+1/2}, \\ U_{i+1}^n, & x > x_{i+1/2}. \end{cases} \quad (5.3)$$

Let the solution (see Subsection 3.2) be denoted by

$$\tilde{U}_i^n(x, t) = \tilde{U}^{(R)}((x - x_{i+1/2})/(t - t_n); U_i^n, U_{i+1}^n), \quad (5.4)$$

for all  $t > t_n$ . For simplicity it is assumed that adjacent Riemann problems do not interfere. If the inequality

$$\Delta t |\lambda|_{\max} < \frac{1}{2} \Delta x \quad (5.5)$$

holds, where  $|\lambda|_{\max} = \max(|\lambda_1|, |\lambda_2|, \dots, |\lambda_m|)$ , then this assumption is fulfilled (cf. [10]).

Finally, in the third step the approximate solution  $U_i^{n+1}$  at time level  $t_{n+1}$  is defined by averaging the exact solution  $\tilde{U}_i^n$  at time  $t_{n+1}$ , thus

$$U_i^{n+1} = \frac{1}{\Delta x} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \tilde{U}_i^n(x, t_{n+1}) dx. \quad (5.6)$$

Note that in this latter equation two different Riemann problems are involved. Hence, this equation can be rewritten as

$$U_i^{n+1} = \frac{1}{\Delta x} \int_0^{\frac{1}{2}\Delta x} \tilde{U}^{(R)}\left(\frac{y}{\Delta t}; U_{i-1}^n, U_i^n\right) dy + \frac{1}{\Delta x} \int_{-\frac{1}{2}\Delta x}^0 \tilde{U}^{(R)}\left(\frac{y}{\Delta t}; U_i^n, U_{i+1}^n\right) dy$$

with respectively  $y = x - x_{i-1/2}$  in the first integral, and  $y = x - x_{i+1/2}$  in the second integral. The values computed in (5.6) are then used to define a new piecewise constant function  $\tilde{U}_i^{n+1}$  (see (5.1)) and the procedure is repeated.

The numerical flux  $F^{(G)}$  of the Godunov scheme can be computed from an integral form of the conservation law (5.2). Since  $\tilde{U}_i^n$  is the exact solution of (5.2), it is easy to see that

$$\begin{aligned} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \tilde{U}_i^n(x, t_{n+1}) dx &= \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \tilde{U}_i^n(x, t_n) dx + \int_{t_n}^{t_{n+1}} f(\tilde{U}_i^n(x_{i-\frac{1}{2}}, t)) dt \\ &\quad - \int_{t_n}^{t_{n+1}} f(\tilde{U}_i^n(x_{i+\frac{1}{2}}, t)) dt, \end{aligned}$$

by integrating (5.2) with respect to space and time. Dividing by  $\Delta x$ , using (5.1) and (5.6) this equation reduces to

$$U_i^{n+1} = U_i^n - \tau \left\{ \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} f(\tilde{U}_i^n(x_{i+\frac{1}{2}}, t)) dt - \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} f(\tilde{U}_i^n(x_{i-\frac{1}{2}}, t)) dt \right\}.$$

Hence, the Godunov method can be written in conservation form (4.4) with the numerical flux

$$F_{i+\frac{1}{2}}^{(G)} = F_{i+\frac{1}{2}}^{(G)}(U_i^n, U_{i+1}^n) = \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} f(\tilde{U}_i^n(x_{i+\frac{1}{2}}, t)) dt.$$

Using (5.4), it is easy to see that the integrand in the above equation is independent of  $t$ . Therefore, the numerical flux can be rewritten as

$$F_{i+\frac{1}{2}}^{(G)} = f(\tilde{U}^{(R)}(0; U_i^n, U_{i+1}^n)). \quad (5.7)$$

If every Riemann solution  $\tilde{U}^{(R)}$  is entropy stable, then it can be shown (cf. [8], [17]) that, under certain assumptions, the Godunov method is an entropy stable method. From (5.7) it directly follows that the Godunov method is consistent. Thus, all hypotheses of Theorem 4.14 are satisfied and therefore, if the method is convergent in the sense of bounded,  $L_1^{loc}$  convergence, then the limit is an entropy stable solution of (2.3).

**Example 5.1** In this example we consider the scalar Burgers' equation

$$\frac{\partial}{\partial t}u(x, t) + \frac{\partial}{\partial x}\left(\frac{1}{2}(u(x, t))^2\right) = 0. \quad (5.8)$$

Using Subsection 3.2, it is easy to see that the corresponding Riemann problem has the following solution. If  $u_L < u_R$ , then the solution is a simple wave, which is given by

$$u(x, t) = u^{(R)}\left(\frac{x}{t}\right) = \begin{cases} u_L & \text{if } x/t < u_L, \\ x/t & \text{if } u_L < x/t < u_R, \\ u_R & \text{if } u_R < x/t. \end{cases} \quad (5.9)$$

If  $u_L > u_R$ , then the solution is a shock wave propagating at speed  $\bar{s} = \frac{1}{2}(u_L + u_R)$  (see (2.14)). This solution is given by

$$u(x, t) = u^{(R)}\left(\frac{x}{t}\right) = \begin{cases} u_L & \text{if } x/t < \bar{s}, \\ u_R & \text{if } x/t > \bar{s}. \end{cases} \quad (5.10)$$

Note that the scalar flux-function in (5.8) is  $f(u) = \frac{1}{2}u^2$ . Using this and (5.7) it is not difficult to see that Godunov's numerical flux is given by (cf. [10], [16])

$$F_{i+\frac{1}{2}}^{(G)} = \begin{cases} \frac{1}{2}(U_{i+1}^n)^2 & \text{if } U_i^n < 0, U_{i+1}^n < 0, \\ \frac{1}{2}(U_i^n)^2 & \text{if } U_i^n > 0, U_{i+1}^n > 0. \end{cases} \quad (5.11)$$

If  $U_i^n$  and  $U_{i+1}^n$  have opposite signs, then the numerical flux is given by

$$F_{i+\frac{1}{2}}^{(G)} = \begin{cases} 0 & \text{if } U_i^n < 0 < U_{i+1}^n, \\ \frac{1}{2}(U_i^n)^2 & \text{if } U_i^n > 0 > U_{i+1}^n \text{ and } \bar{s}_{i+\frac{1}{2}}^n > 0, \\ \frac{1}{2}(U_{i+1}^n)^2 & \text{if } U_i^n > 0 > U_{i+1}^n \text{ and } \bar{s}_{i+\frac{1}{2}}^n < 0, \end{cases} \quad (5.12)$$

where  $\bar{s}_{i+\frac{1}{2}}^n$  is the propagation speed of the shock, defined by  $\bar{s}_{i+\frac{1}{2}}^n = \frac{1}{2}(U_i^n + U_{i+1}^n)$ .

### 5.3 Osher's method

The idea of *Osher's approximate Riemann solver* is to split the numerical flux (5.7) in a forward and a backward flux (cf. [10]). Osher's method is first introduced for a nonlinear scalar conservation law. Subsequently, the method will be extended for systems of nonlinear hyperbolic conservation laws. Let the scalar conservation law

$$\frac{\partial}{\partial t}u(x, t) + \frac{\partial}{\partial x}f(u(x, t)) = 0$$

be given. Let  $a(u)$  be defined by  $a(u) = f'(u)$  (see (2.5)) and, furthermore, define the functions  $a^+$  and  $a^-$  by

$$\begin{aligned} a^+(\xi) &= \max(a(\xi), 0) \geq 0, \\ a^-(\xi) &= \min(a(\xi), 0) \leq 0. \end{aligned}$$

Note that  $a = a^+ + a^-$ . Next, the forward flux-function  $f^+$  and the backward flux-function  $f^-$  are given by, respectively,

$$\frac{d}{du}f^+(u) = a^+(u), \quad \frac{d}{du}f^-(u) = a^-(u). \quad (5.13)$$



It is assumed that the equality

$$f = f^+ + f^-$$

holds. The exact Riemann flux  $f(\tilde{U}^{(R)}(0; U_i^n, U_{i+1}^n))$  (see (5.7)) is now approximated by (cf. [10], [26])

$$\begin{aligned} F_{i+\frac{1}{2}}^{(O)}(U_i^n, U_{i+1}^n) &= f^+(U_i^n) + f^-(U_{i+1}^n) \\ &= f(U_i^n) - f^-(U_i^n) + f^-(U_{i+1}^n) \\ &= f(U_i^n) + \int_{U_i^n}^{U_{i+1}^n} a^-(u) du. \end{aligned} \quad (5.14)$$

To generalize this definition for systems the following concepts are useful. Define

$$\begin{aligned} \lambda_k^+(\mathbf{u}) &= \max(\lambda_k(\mathbf{u}), 0) \geq 0, \\ \lambda_k^-(\mathbf{u}) &= \min(\lambda_k(\mathbf{u}), 0) \leq 0, \end{aligned} \quad (5.15)$$

for all  $k$ , with  $1 \leq k \leq m$ . As in Definition 2.1, the diagonal matrices  $\Lambda^+(\mathbf{u})$ ,  $\Lambda^-(\mathbf{u})$  and  $|\Lambda|(\mathbf{u})$  are defined by

$$\begin{aligned} \Lambda^+(\mathbf{u}) &= \text{diag}(\lambda_1^+(\mathbf{u}), \lambda_2^+(\mathbf{u}), \dots, \lambda_m^+(\mathbf{u})), \\ \Lambda^-(\mathbf{u}) &= \text{diag}(\lambda_1^-(\mathbf{u}), \lambda_2^-(\mathbf{u}), \dots, \lambda_m^-(\mathbf{u})), \\ |\Lambda|(\mathbf{u}) &= \text{diag}(|\lambda_1(\mathbf{u})|, |\lambda_2(\mathbf{u})|, \dots, |\lambda_m(\mathbf{u})|). \end{aligned} \quad (5.16)$$

The matrices  $A^+$ ,  $A^-$  and  $|A|$ , related to the Jacobian matrix  $A$ , are logically defined by (see (2.6))

$$\begin{aligned} A^+(\mathbf{u}) &= R(\mathbf{u})\Lambda^+(\mathbf{u})R^{-1}(\mathbf{u}), \\ A^-(\mathbf{u}) &= R(\mathbf{u})\Lambda^-(\mathbf{u})R^{-1}(\mathbf{u}), \\ |A|(\mathbf{u}) &= R(\mathbf{u})|\Lambda|(\mathbf{u})R^{-1}(\mathbf{u}). \end{aligned} \quad (5.17)$$

Now we are able to generalize (5.13)-(5.14) for systems of equations. Suppose that there exist continuously differentiable functions  $\mathbf{f}^+$  and  $\mathbf{f}^-$  such that

$$\frac{\partial}{\partial \mathbf{u}} \mathbf{f}^+(\mathbf{u}) = A^+(\mathbf{u}), \quad \frac{\partial}{\partial \mathbf{u}} \mathbf{f}^-(\mathbf{u}) = A^-(\mathbf{u}), \quad (5.18)$$

and the equality

$$\mathbf{f} = \mathbf{f}^+ + \mathbf{f}^- \quad (5.19)$$

holds. Then the exact Riemann flux  $\mathbf{f}(\tilde{U}^{(R)}(0; U_i^n, U_{i+1}^n))$  is approximated by

$$\begin{aligned} \mathbf{F}_{i+\frac{1}{2}}^{(O)}(U_i^n, U_{i+1}^n) &= \mathbf{f}^+(U_i^n) + \mathbf{f}^-(U_{i+1}^n) \\ &= \mathbf{f}(U_i^n) - \mathbf{f}^-(U_i^n) + \mathbf{f}^-(U_{i+1}^n) \\ &= \mathbf{f}(U_i^n) + \int_{U_i^n}^{U_{i+1}^n} A^-(\mathbf{u}) d\mathbf{u}. \end{aligned} \quad (5.20)$$

Equivalent representations of  $\mathbf{f}_{i+1/2}^{(O)}$  in terms of  $A^+$  or  $|A|$  are easy to find. Unfortunately, in general no functions  $\mathbf{f}^+$  and  $\mathbf{f}^-$  exist such that (5.18) and (5.19) hold (cf. [26]). This is equivalent with the observation that the integrals

$$\int_{U_i^n}^{U_{i+1}^n} A^-(\mathbf{u}) d\mathbf{u} \quad \text{and} \quad \int_{U_i^n}^{U_{i+1}^n} A^+(\mathbf{u}) d\mathbf{u}$$

depend on their integration path. To remain consistent with the scalar case, this fact is simply ignored. Therefore, Osher's numerical flux is given by (5.20), where the integration path is chosen in such a way that the evolution of the integral in (5.20) is easy.

To explain how the integration path is chosen, Osher's flux is rewritten in the more general form (cf. [26])

$$\mathbf{F}^{(O)}(\mathbf{u}_L, \mathbf{u}_R) = \mathbf{f}(\mathbf{u}_L) + \int_{\mathbf{u}_L}^{\mathbf{u}_R} A^-(\mathbf{u}) d\mathbf{u}, \quad (5.21)$$

where  $\mathbf{u}_L \in \mathbb{R}^m$  and  $\mathbf{u}_R \in \mathbb{R}^m$  are some given constant states. Suppose that these states can be interconnected by an integration path  $\Gamma_k$ , which is tangential to the eigenvector  $\mathbf{r}^{(k)}$ , i.e. (see (3.5))

$$\begin{cases} \frac{d}{d\xi} \mathbf{u}(\xi) = \mathbf{r}^{(k)}(\mathbf{u}(\xi)), \\ \mathbf{u}(0) = \mathbf{u}_L, \mathbf{u}(\xi_R) = \mathbf{u}_R. \end{cases} \quad (5.22)$$

Thus, the constant states are connected by a simple wave, contact discontinuity or compression wave. This is an important property of Osher's method. Using (5.17) and (5.22) it is seen that

$$\begin{aligned} \int_{\mathbf{u}_L}^{\mathbf{u}_R} A^-(\mathbf{u}) d\mathbf{u} &= \int_0^{\xi_R} A^-(\mathbf{u}(\xi)) \frac{d}{d\xi} \mathbf{u}(\xi) d\xi \\ &= \int_0^{\xi_R} A^-(\mathbf{u}(\xi)) \mathbf{r}^{(k)}(\mathbf{u}(\xi)) d\xi \\ &= \int_0^{\xi_R} \lambda_k^-(\mathbf{u}(\xi)) \mathbf{r}^{(k)}(\mathbf{u}(\xi)) d\xi. \end{aligned}$$

Note that it is assumed that  $\mathbf{r}^{(k)}(\mathbf{u})$  is either genuinely nonlinear or linearly degenerate. If the eigenvector  $\mathbf{r}^{(k)}(\mathbf{u})$  is genuinely nonlinear, then  $\lambda_k$  is monotone along  $\Gamma_k$  (see (3.6)), and if  $\mathbf{r}^{(k)}(\mathbf{u})$  is linearly degenerate, then  $\lambda_k$  is constant along  $\Gamma_k$  (see (3.10)). Thus we can distinguish between two possibilities: either  $\lambda_k$  does not change its sign along the integration path  $\Gamma_k$  or  $\lambda_k$  changes its sign only once along the integration path  $\Gamma_k$ .

Assume first that  $\lambda_k$  does not change its sign along  $\Gamma_k$ . If  $\lambda_k(\mathbf{u}(\xi)) \geq 0$  for all  $\xi \in [0, \xi_R]$ , then  $\lambda_k^- = 0$  and

$$\int_{\mathbf{u}_L}^{\mathbf{u}_R} A^-(\mathbf{u}) d\mathbf{u} = 0.$$

If  $\lambda_k(\mathbf{u}(\xi)) < 0$  for all  $\xi \in [0, \xi_R]$ , then  $\lambda_k^- = \lambda_k$  and

$$\begin{aligned} \int_{\mathbf{u}_L}^{\mathbf{u}_R} A^-(\mathbf{u}) d\mathbf{u} &= \int_0^{\xi_R} A(\mathbf{u}(\xi)) \mathbf{r}^{(k)}(\mathbf{u}(\xi)) d\xi \\ &= \int_0^{\xi_R} \frac{\partial}{\partial \mathbf{u}} \mathbf{f}(\mathbf{u}(\xi)) \frac{d}{d\xi} \mathbf{u}(\xi) d\xi \\ &= \int_{\mathbf{u}_L}^{\mathbf{u}_R} \frac{\partial}{\partial \mathbf{u}} \mathbf{f}(\mathbf{u}) d\mathbf{u} = \mathbf{f}(\mathbf{u}_R) - \mathbf{f}(\mathbf{u}_L). \end{aligned}$$

Next assume that  $\lambda_k$  changes its sign once along  $\Gamma_k$ . Suppose that  $\lambda_k(\mathbf{u}(\xi_S)) = 0$  with  $0 < \xi_S < \xi_R$  and define  $\mathbf{u}_S = \mathbf{u}(\xi_S)$ . The point  $\mathbf{u}_S$  is called a *sonic point*. If

$\lambda_k(\mathbf{u}(\xi)) > 0$  for all  $\xi \in [0, \xi_S)$  and  $\lambda_k(\mathbf{u}(\xi)) < 0$  for all  $\xi \in (\xi_S, \xi_R]$ , then  $\lambda_k^-(\mathbf{u}(\xi)) = 0$  for all  $\xi \in [0, \xi_S)$  and

$$\int_{\mathbf{u}_L}^{\mathbf{u}_R} A^-(\mathbf{u}) d\mathbf{u} = \int_{\xi_S}^{\xi_R} \lambda_k(\mathbf{u}(\xi)) \mathbf{r}^{(k)}(\mathbf{u}(\xi)) d\xi = \mathbf{f}(\mathbf{u}_R) - \mathbf{f}(\mathbf{u}_S).$$

If  $\lambda_k(\mathbf{u}(\xi)) < 0$  for all  $\xi \in [0, \xi_S)$  and  $\lambda_k(\mathbf{u}(\xi)) > 0$  for all  $\xi \in (\xi_S, \xi_R]$ , then  $\lambda_k^-(\mathbf{u}(\xi)) = 0$  for all  $\xi \in (\xi_S, \xi_R]$  and

$$\int_{\mathbf{u}_L}^{\mathbf{u}_R} A^-(\mathbf{u}) d\mathbf{u} = \int_0^{\xi_S} \lambda_k(\mathbf{u}(\xi)) \mathbf{r}^{(k)}(\mathbf{u}(\xi)) d\xi = \mathbf{f}(\mathbf{u}_S) - \mathbf{f}(\mathbf{u}_L).$$

Thus, if the states  $\mathbf{u}_L$  and  $\mathbf{u}_R$  can be interconnected by an integration path as defined in (5.22), Osher's numerical flux (5.21) becomes (cf. [22], [26])

$$\mathbf{F}^{(O)}(\mathbf{u}_L, \mathbf{u}_R) = \begin{cases} \mathbf{f}(\mathbf{u}_L) & \text{if } \lambda_k \geq 0 \text{ along } \Gamma_k, \\ \mathbf{f}(\mathbf{u}_R) & \text{if } \lambda_k \leq 0 \text{ along } \Gamma_k, \\ \mathbf{f}(\mathbf{u}_R) - \mathbf{f}(\mathbf{u}_S) + \mathbf{f}(\mathbf{u}_L) & \text{if } \lambda_k(\mathbf{u}_L) > 0, \lambda_k(\mathbf{u}_R) < 0, \lambda_k(\mathbf{u}_S) = 0, \\ \mathbf{f}(\mathbf{u}_S) & \text{if } \lambda_k(\mathbf{u}_L) < 0, \lambda_k(\mathbf{u}_R) > 0, \lambda_k(\mathbf{u}_S) = 0. \end{cases} \quad (5.23)$$

Now consider Osher's flux as given by (5.20). A general pair  $(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  can be connected by a continuous integration path  $\Gamma$  which is subdivided into  $m$  subcurves  $\Gamma_k$ , i.e.

$$\Gamma = \bigcup_{k=1}^m \Gamma_k,$$

where each subcurve  $\Gamma_k$  is tangential to the right eigenvector  $\mathbf{r}^{(k)}$  (see (5.22)). The subcurve  $\Gamma_1$  starts in  $\mathbf{U}_i^n$  and the subcurve  $\Gamma_m$  ends in  $\mathbf{U}_{i+1}^n$ . Define the  $m-1$  points of intersection  $\mathbf{U}_{i+\frac{k}{m}}^n$ ,  $k = 1, \dots, m-1$  by

$$\mathbf{U}_{i+\frac{k}{m}}^n = \Gamma_k \cap \Gamma_{k+1}.$$

These intersection points can easily be found by using Riemann invariants (see Definition 3.8). In Subsection 3.3 it is shown that along the subcurve  $\Gamma_k$  the  $m-1$   $k$ -Riemann invariants  $w_k^1, w_k^2, \dots, w_k^{m-1}$  remain constant. Therefore, the following equalities must hold

$$\begin{aligned} w_k^1(\mathbf{U}_{i+\frac{k-1}{m}}^n) &= w_k^1(\mathbf{U}_{i+\frac{k}{m}}^n), \\ w_k^2(\mathbf{U}_{i+\frac{k-1}{m}}^n) &= w_k^2(\mathbf{U}_{i+\frac{k}{m}}^n), \\ &\vdots \\ w_k^{m-1}(\mathbf{U}_{i+\frac{k-1}{m}}^n) &= w_k^{m-1}(\mathbf{U}_{i+\frac{k}{m}}^n), \end{aligned} \quad (5.24)$$

for  $k = 1, \dots, m$ . In this way a nonsingular system of  $m(m-1)$  equations is obtained for the  $m(m-1)$  unknowns  $\mathbf{U}_{i+k/m}^n$ ,  $k = 1, \dots, m-1$ . If along a particular subcurve,  $\Gamma_k$ , say, a sonic point  $\mathbf{U}_S$  occurs, then this sonic point can be computed by adding to

the system (5.24) the  $m$  equations

$$\begin{aligned} w_{k*}^1(\mathbf{U}_{i+\frac{k*}{m}}^n) &= w_{k*}^1(\mathbf{U}_S^n), \\ w_{k*}^2(\mathbf{U}_{i+\frac{k*}{m}}^n) &= w_{k*}^2(\mathbf{U}_S^n), \\ &\vdots \\ w_{k*}^{m-1}(\mathbf{U}_{i+\frac{k*}{m}}^n) &= w_{k*}^{m-1}(\mathbf{U}_S^n), \\ \lambda_{k*}(\mathbf{U}_S^n) &= 0. \end{aligned}$$

A nonsingular system of  $m^2$  equations for  $m^2$  unknowns is therefore obtained.

Osher's flux is now defined by

$$\mathbf{F}_{i+\frac{1}{2}}^{(O)}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = \mathbf{f}(\mathbf{U}_i^n) + \sum_{k=1}^m \int_{\Gamma_k} A^-(\mathbf{u}) d\mathbf{u}. \quad (5.25)$$

Here each integral along a subcurve  $\Gamma_k$  is evaluated in the manner described by (5.21)-(5.23).

In [22] it is shown that under fairly general hypotheses Osher's method is entropy stable.

**Example 5.2** In this example again the scalar Burgers' equation (5.8) is considered, with a solution given by (5.9) or (5.10). For  $f(u) = \frac{1}{2}u^2$  equation (5.20) becomes

$$F_{i+\frac{1}{2}}^{(O)}(U_i^n, U_{i+1}^n) = \frac{1}{2}(U_i^n)^2 + \int_{U_i^n}^{U_{i+1}^n} u^- du.$$

Using this equality and (5.23), it is not difficult to see that Osher's numerical flux for the scalar Burgers' equation is given by (cf. [10], [16])

$$F_{i+\frac{1}{2}}^{(O)} = \begin{cases} \frac{1}{2}(U_{i+1}^n)^2 & \text{if } U_i^n < 0, U_{i+1}^n < 0, \\ \frac{1}{2}(U_i^n)^2 & \text{if } U_i^n > 0, U_{i+1}^n > 0. \end{cases}$$

When  $U_i^n$  and  $U_{i+1}^n$  have opposite signs, the flux is given by

$$F_{i+\frac{1}{2}}^{(O)} = \begin{cases} 0 & \text{if } U_i^n < 0 < U_{i+1}^n, \\ \frac{1}{2}(U_i^n)^2 + \frac{1}{2}(U_{i+1}^n)^2 & \text{if } U_i^n > 0 > U_{i+1}^n. \end{cases}$$

Compared to the Godunov scheme (see (5.11), (5.12)), this scheme only differs by the representation of transonic shocks (i.e.  $\lambda_k$  changes its sign across the shock). As shown by van Leer in [16], the Osher scheme replaces the shock in the exact Riemann solution by a multivalued compression wave.

**Example 5.3** In this example Osher's method is applied to the shocktube problem for the Euler equations. Hence,  $m = 3$  and a sonic point possibly occurs along  $\Gamma_1$  or  $\Gamma_3$ . For every pair  $(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  the equations derived in Example 3.9 are used to compute the two intersection points  $\mathbf{U}_{i+1/3}^n$  and  $\mathbf{U}_{i+2/3}^n$ . In, for example, the special case that  $(u_i^n < c_i^n)$ ,  $(-c_{i+2/3}^n < u_{i+1/3}^n < 0)$  and  $(-u_{i+1}^n < c_{i+1}^n)$  hold, Osher's flux is given by

$$\begin{aligned} \mathbf{F}_{i+\frac{1}{2}}^{(O)} &= \mathbf{f}(\mathbf{U}_i^n) + \int_{U_i^n}^{U_{i+1/3}^n} A^-(\mathbf{u}) d\mathbf{u} + \int_{U_{i+1/3}^n}^{U_{i+2/3}^n} A^-(\mathbf{u}) d\mathbf{u} + \int_{U_{i+2/3}^n}^{U_{i+1}^n} A^-(\mathbf{u}) d\mathbf{u} \\ &= \mathbf{f}(\mathbf{U}_i^n) + \mathbf{f}(\mathbf{U}_{i+1/3}^n) - \mathbf{f}(\mathbf{U}_i^n) + \mathbf{f}(\mathbf{U}_{i+2/3}^n) - \mathbf{f}(\mathbf{U}_{i+1/3}^n) + 0 \\ &= \mathbf{f}(\mathbf{U}_{i+2/3}^n). \end{aligned}$$

Similar calculations give the numerical flux-function in all other cases (cf. [26]). The numerical results in Figure 3 illustrate clearly that Osher's method is a first order method, since the discontinuities are smeared out.

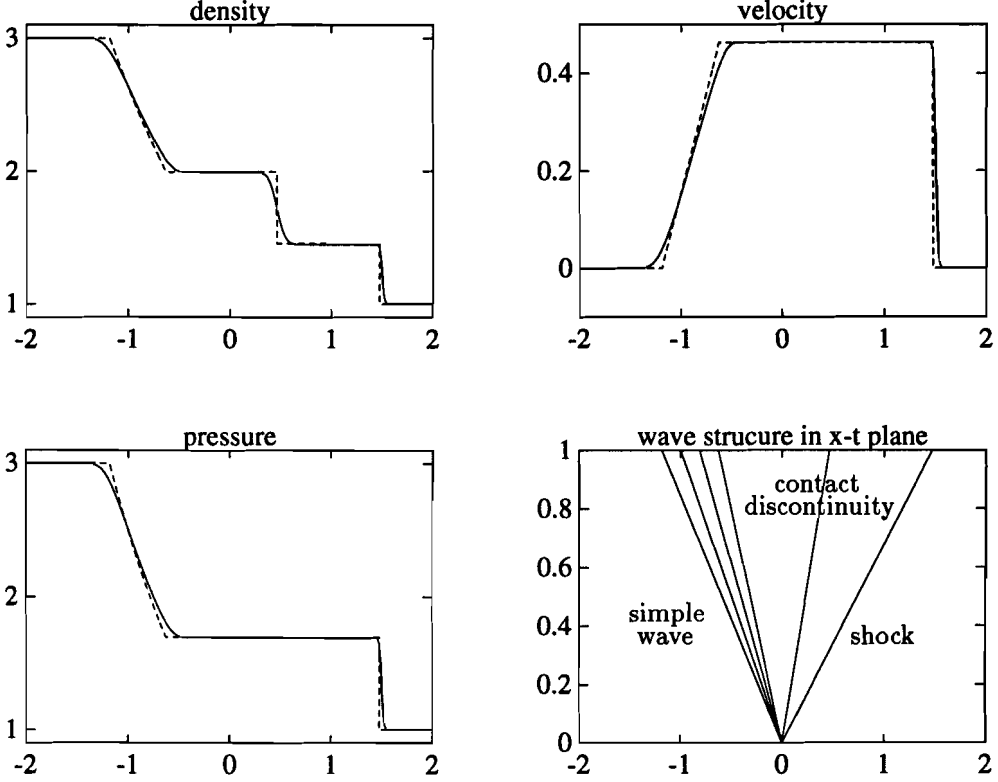


Figure 3: Numerical solution (solid line) computed with Osher's method (with  $\tau = 0.1$ ) and exact solution (dashed line) at time  $t = 1$  of a shock tube problem for the one-dimensional Euler equations (2.7) with initial conditions  $p(x, 0) = 3$ ,  $\rho(x, 0) = 3$ ,  $u(x, 0) = 0$  if  $x < 0$  and  $p(x, 0) = 1$ ,  $\rho(x, 0) = 1$ ,  $u(x, 0) = 0$  if  $x > 0$ .

#### 5.4 Roe's method

Another approach to decrease the computational cost of the basic Godunov method is to solve an approximate Riemann problem at the cell interfaces instead of (5.2)-(5.3). Therefore, consider the following Riemann problem

$$\frac{\partial}{\partial t} \hat{U}_i^n(x, t) + \frac{\partial}{\partial x} \hat{f}(\hat{U}_i^n(x, t)) = 0 \quad (5.26)$$

with

$$\hat{U}_i^n(x, t_n) = \begin{cases} U_i^n, & x < x_{i+\frac{1}{2}}, \\ U_{i+1}^n, & x > x_{i+\frac{1}{2}}. \end{cases} \quad (5.27)$$

Here  $\hat{f}$  is an approximation of  $f$ . Let the solution (see Subsection 3.2) be denoted by

$$\hat{U}_i^n(x, t) = \hat{U}^{(R)}((x - x_{i+\frac{1}{2}})/(t - t_n); U_i^n, U_{i+1}^n), \quad (5.28)$$

for all  $t > t_n$ . The approximate solution  $\mathbf{U}_i^{n+1}$  at time level  $t_{n+1}$  is defined by averaging the exact solution  $\hat{\mathbf{U}}_i^n$  at time  $t_{n+1}$ , thus

$$\mathbf{U}_i^{n+1} = \frac{1}{\Delta x} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \hat{\mathbf{U}}_i^n(x, t_{n+1}) dx. \quad (5.29)$$

The method is conservative if the solution  $\hat{\mathbf{U}}$  of the approximate Riemann problem has the following property (cf. [8], [17])

$$\begin{aligned} \int_{x_i}^{x_{i+\frac{1}{2}}} \hat{\mathbf{U}}_i^n(x, t_{n+1}) dx &= \int_{x_i}^{x_{i+\frac{1}{2}}} \hat{\mathbf{U}}_i^n(x, t_n) dx + \int_{t_n}^{t_{n+1}} \mathbf{f}(\hat{\mathbf{U}}_i^n(x_i, t)) dt \\ &\quad - \int_{t_n}^{t_{n+1}} \mathbf{f}(\hat{\mathbf{U}}_i^n(x_{i+\frac{1}{2}}, t)) dt. \end{aligned}$$

Using (5.1), (5.27) and the assumption that adjacent Riemann problems do not interfere (see (5.5)), this equation can be rewritten as

$$\hat{\mathbf{F}}_{i+\frac{1}{2}}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = \mathbf{f}(\mathbf{U}_i^n) - \frac{1}{\Delta t} \int_{x_i}^{x_{i+\frac{1}{2}}} \hat{\mathbf{U}}^{(R)}\left(\frac{x - x_{i+\frac{1}{2}}}{\Delta t}; \mathbf{U}_i^n, \mathbf{U}_{i+1}^n\right) dy + \frac{\Delta x}{2\Delta t} \mathbf{U}_i^n, \quad (5.30)$$

with the numerical flux-function  $\hat{\mathbf{F}}$  given by

$$\hat{\mathbf{F}}_{i+\frac{1}{2}} = \hat{\mathbf{F}}_{i+\frac{1}{2}}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \mathbf{f}(\hat{\mathbf{U}}_i^n(x_{i+\frac{1}{2}}, t)) dt.$$

A popular approximate Riemann problem is due to Roe (cf. [24]). The idea is to determine  $\hat{\mathbf{U}}_i^n$  by solving a constant coefficient linear system of conservation laws. Therefore, let  $\hat{\mathbf{f}}$  be given by

$$\hat{\mathbf{f}}(\hat{\mathbf{U}}_i^n(x, t)) = \hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) \hat{\mathbf{U}}_i^n(x, t),$$

where  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  is a constant  $m \times m$ -matrix. Thus, the system (5.26) can be rewritten as

$$\frac{\partial}{\partial t} \hat{\mathbf{U}}_i^n(x, t) + \hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) \frac{\partial}{\partial x} (\hat{\mathbf{U}}_i^n(x, t)) = 0. \quad (5.31)$$

Roe requires that the matrix  $\hat{A}$  has the following properties:

- (i) if  $\mathbf{U}_i^n, \mathbf{U}_{i+1}^n \rightarrow \bar{\mathbf{u}}$ , then  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) \rightarrow A(\bar{\mathbf{u}})$  with  $\bar{\mathbf{u}}$  some point between  $\mathbf{U}_i^n$  and  $\mathbf{U}_{i+1}^n$ ;
- (ii)  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)(\mathbf{U}_{i+1}^n - \mathbf{U}_i^n) = \mathbf{f}(\mathbf{U}_{i+1}^n) - \mathbf{f}(\mathbf{U}_i^n)$ ;
- (iii)  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  is diagonalizable with real eigenvalues.

Condition (i) is a necessary condition in order to recover smoothly the linearized algorithm from the non-linear version. Condition (ii) has two effects. Firstly it is sufficient to guarantee that the scheme is conservative, and secondly, in the special case that  $\mathbf{U}_i^n$  and  $\mathbf{U}_{i+1}^n$  are connected by a single shock wave or contact discontinuity, the approximate Riemann solution agrees with the exact Riemann solution (cf. [17], [24]). Finally, condition (iii) is clearly required in order that the problem is hyperbolic and solvable.

Since, instead of the original Riemann problem (5.2)-(5.3), Roe considers a linear Riemann problem, the approximate Riemann solver recognizes only discontinuities (cf. [10]).

It is very easy to construct  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  such that condition (i) is satisfied (cf. [24]). Condition (iii) can be easily checked a posteriori. The difficulty arises entirely from

condition (ii). In [8] it is shown that for a general system with an entropy function, a complicated averaging of the Jacobian matrix can be used for  $\hat{A}$ . This shows that a matrix  $\hat{A}$  exists that satisfies the conditions (i)-(iii), but, unfortunately, it appears that the computed matrix is too complicated to use in practice. Fortunately, for special systems of equations it is possible to derive suitable matrices that are very efficient to use relative to the exact Riemann solution. For example in [24] a suitable matrix  $\hat{A}$  is derived for the Euler equations and in [17] a matrix is derived for the isothermal Euler equations.

In the following it is assumed that there exists a matrix  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$ , such that the conditions (i)-(iii) are satisfied.

Condition (iii) implies that there exists a real diagonal matrix  $\hat{\Lambda}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  and a non-singular real matrix  $\hat{R}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  such that

$$\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) \hat{R}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = \hat{R}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) \hat{\Lambda}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n).$$

Here  $\hat{\Lambda}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = \text{diag}(\hat{\lambda}_1(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n), \hat{\lambda}_2(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n), \dots, \hat{\lambda}_m(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n))$  is the matrix of the eigenvalues of  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$ , where the eigenvalues are labeled in increasing order. The matrix  $\hat{R}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = (\hat{\mathbf{r}}^{(1)}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n), \hat{\mathbf{r}}^{(2)}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n), \dots, \hat{\mathbf{r}}^{(m)}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n))$  is the matrix of the corresponding right eigenvectors of  $\hat{A}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$ . For shortness of notation,  $\hat{\lambda}_k(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  and  $\hat{\mathbf{r}}^{(k)}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$  will simply be denoted by  $\hat{\lambda}_k$  and  $\hat{\mathbf{r}}^{(k)}$ . For all  $k$  with  $1 \leq k \leq m$ ,  $\hat{\lambda}_k^+$  and  $\hat{\lambda}_k^-$  are defined analogously to (5.15). Further, the diagonal matrices  $\hat{\Lambda}^+$ ,  $\hat{\Lambda}^-$  and  $|\hat{\Lambda}|$  are defined as in (5.16). The matrices  $\hat{A}^+$ ,  $\hat{A}^-$  and  $|\hat{A}|$  are defined in the same way as  $A^+$ ,  $A^-$  and  $|A|$  in (5.17).

For the linear Riemann problem (5.31) the solution is given in Example 3.7. Since all the eigenvectors are linearly independent the initial states  $\mathbf{U}_i^n$  and  $\mathbf{U}_{i+1}^n$  of (5.31) can be decomposed as

$$\mathbf{U}_{i+1}^n - \mathbf{U}_i^n = \sum_{k=1}^m \hat{\alpha}_k \hat{\mathbf{r}}^{(k)}, \quad (5.32)$$

where  $\hat{\alpha}_k \in \mathbb{R}$  for all  $k$  with  $1 \leq k \leq m$ . The solution of (5.31) is then given by (analogous to (3.15))

$$\hat{\mathbf{U}}^{(R)}(x, t) = \mathbf{U}_i^n + \sum_{k=1}^m \hat{\alpha}_k H(x - x_{i+\frac{1}{2}} - \hat{\lambda}_k(t - t_n)) \hat{\mathbf{r}}^{(k)}, \quad (5.33)$$

for all  $t > t_n$ . After substituting this solution in (5.30), Roe's numerical flux is derived

$$\mathbf{F}_{i+\frac{1}{2}}^{(R)} = \mathbf{F}_{i+\frac{1}{2}}^{(R)}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n) = \mathbf{f}(\mathbf{U}_i^n) + \sum_{k=1}^m \hat{\lambda}_k^- \hat{\alpha}_k \hat{\mathbf{r}}^{(k)}. \quad (5.34)$$

It is known that Roe's method can include a physically inadmissible expansion shock. This is a direct consequence of the admission of an expansion shock in the underlying approximate Riemann solution. Roe has proposed a modification of the numerical flux function for a transonic expansion that excludes expansion shocks (cf. [10], [16], [17]).

**Example 5.4** In this example again Burgers' scalar equation (5.8) is considered, with a solution given by (5.9) or (5.10). Let  $\bar{s}_{i+1/2}^n = \frac{1}{2}(U_i^n + U_{i+1}^n)$  (see Example 5.1). Note that  $A(u) = u$  and let  $\hat{A}(U_i^n, U_{i+1}^n)$  be defined by

$$\hat{A}(U_i^n, U_{i+1}^n) = \frac{1}{2}(U_i^n + U_{i+1}^n) = \bar{s}_{i+\frac{1}{2}}^n.$$

It is easy to see that the conditions (i)-(iii) are satisfied. For  $f(u) = \frac{1}{2}u^2$  equation (5.34) becomes

$$F_{i+\frac{1}{2}}^{(R)}(U_i^n, U_{i+1}^n) = \frac{1}{2}(U_i^n)^2 + (\bar{s}_{i+\frac{1}{2}}^n)^-(U_{i+1}^n - U_i^n).$$

Using, this it is not difficult to see that Roe's numerical flux for the scalar Burgers' equation is given by (cf. [10], [16])

$$F_{i+\frac{1}{2}}^{(R)} = \begin{cases} \frac{1}{2}(U_{i+1}^n)^2 & \text{if } \bar{s}_{i+\frac{1}{2}}^n < 0, \\ \frac{1}{2}(U_i^n)^2 & \text{if } \bar{s}_{i+\frac{1}{2}}^n > 0. \end{cases}$$

The numerical flux-function deviates from the Godunov flux-function (see (5.11),(5.12)) only in the case of a transonic expansion wave. In [16] it is shown that in this case Roe's method replaces the transonic expansion wave by a so-called *expansion shock*.

**Example 5.5** In the last example of this section Roe's method is applied to the shock-tube problem for the Euler equations. The matrix  $\hat{A}(U_i^n, U_{i+1}^n)$  is derived in [24], and is for every pair  $(U_i^n, U_{i+1}^n)$  given by

$$\hat{A}(U_i^n, U_{i+1}^n) = \begin{pmatrix} 0 & 1 & 0 \\ \frac{\gamma-3}{2}\hat{u}^2 & (3-\gamma)\hat{u} & \gamma-1 \\ \hat{u}(\frac{\gamma-1}{2}\hat{u}^2 - \hat{H}) & \hat{H} - (\gamma-1)\hat{u}^2 & \gamma\hat{u} \end{pmatrix},$$

where the quantities  $\hat{u}$  and  $\hat{H}$  are defined as

$$\hat{u} = \frac{(u\sqrt{\rho})_{i+1}^n + (u\sqrt{\rho})_i^n}{\sqrt{\rho_{i+1}^n} + \sqrt{\rho_i^n}}, \quad \hat{H} = \frac{(H\sqrt{\rho})_{i+1}^n + (H\sqrt{\rho})_i^n}{\sqrt{\rho_{i+1}^n} + \sqrt{\rho_i^n}}. \quad (5.35)$$

In order to derive the eigenvectors and the eigenvalues of the matrix  $\hat{A}(U_i^n, U_{i+1}^n)$  the following quantities are useful. Define

$$\hat{\rho} = \sqrt{\rho_{i+1}^n \rho_i^n}, \quad \hat{c}^2 = (\gamma-1)(\hat{H} - \frac{\hat{u}^2}{2}). \quad (5.36)$$

Now the computation of the eigenvalues and the eigenvectors is not difficult. They are given by

$$\hat{\lambda}_1(U_i^n, U_{i+1}^n) = \hat{u} - \hat{c}, \quad \hat{\lambda}_2(U_i^n, U_{i+1}^n) = \hat{u}, \quad \hat{\lambda}_3(U_i^n, U_{i+1}^n) = \hat{u} + \hat{c}, \quad (5.37)$$

and

$$\begin{aligned} \hat{\mathbf{r}}^{(1)}(U_i^n, U_{i+1}^n) &= -\frac{\hat{\rho}}{2\hat{c}}(1, \hat{u} - \hat{c}, \hat{H} - \hat{u}\hat{c})^T, \\ \hat{\mathbf{r}}^{(2)}(U_i^n, U_{i+1}^n) &= (1, \hat{u}, \frac{\hat{u}^2}{2})^T, \\ \hat{\mathbf{r}}^{(3)}(U_i^n, U_{i+1}^n) &= \frac{\hat{\rho}}{2\hat{c}}(1, \hat{u} + \hat{c}, \hat{H} + \hat{u}\hat{c})^T, \end{aligned} \quad (5.38)$$

Hence, for every pair  $(U_i^n, U_{i+1}^n)$  Roe's numerical flux  $F_{i+1/2}$  at the cell interface  $x_{i+1/2}$  is derived by a three-step procedure. The first step is the computation of the quantities defined in (5.35) and (5.36). In the second step the eigenvalues (5.37) and eigenvectors (5.38) are computed. Finally, in the third step, (5.32) is used to compute  $\hat{\alpha}_1$ ,  $\hat{\alpha}_2$  and  $\hat{\alpha}_3$ .



The computation of Roe's numerical flux (5.34) is now straightforward. The numerical results (see Figure 4) illustrates clearly that Roe's method is also a first order method, since again the discontinuities are smoothed.

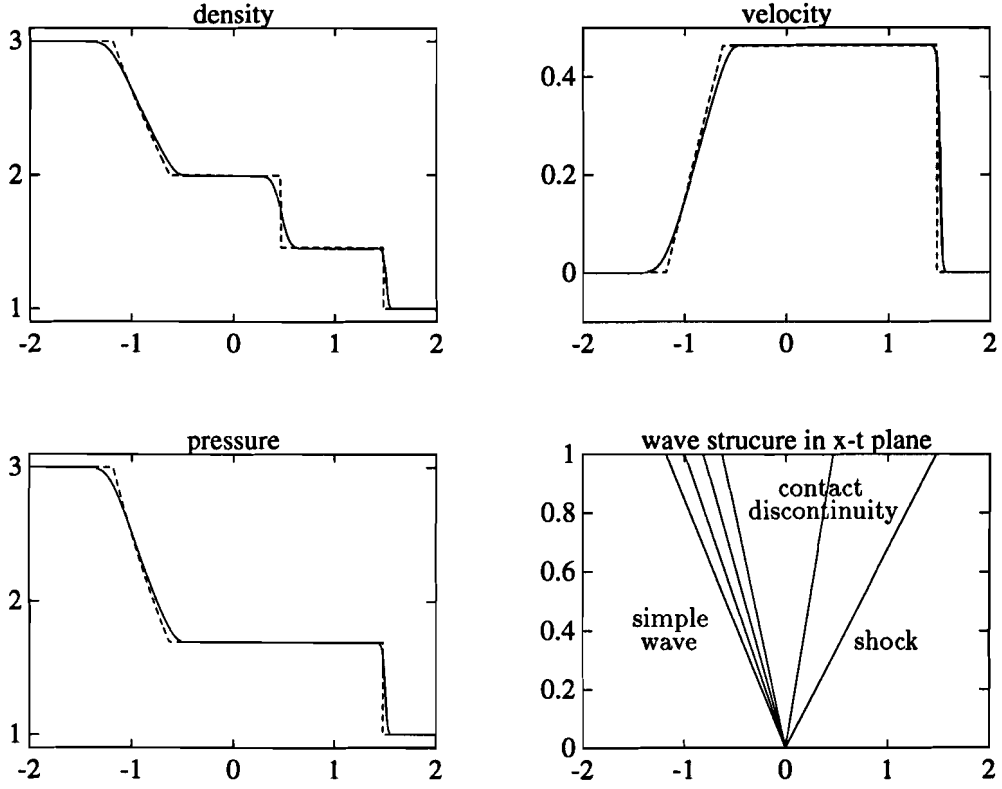


Figure 4: Numerical solution (solid line) computed with Roe's method (with  $\tau = 0.1$ ) and exact solution (dashed line) at time  $t = 1$  of a shock tube problem for the one-dimensional Euler equations (2.7) with initial conditions  $p(x,0) = 3$ ,  $\rho(x,0) = 3$ ,  $u(x,0) = 0$  if  $x < 0$  and  $p(x,0) = 1$ ,  $\rho(x,0) = 1$ ,  $u(x,0) = 0$  if  $x > 0$ .

## 6 High Resolution Methods

### 6.1 Some convergence results

In this section *high resolution methods* are introduced. High resolution methods are numerical methods which are second order accurate in regions where the solution is smooth, and give good results (no oscillations) around shocks.

In the preceding sections we have not investigated whether a numerical method converges, only that if a sequence of approximations converge, then the limit is a weak solution. In this subsection theorems will be presented for the scalar case, which, under certain assumptions, guarantee convergence of a method.

Consider the non-linear scalar conservation law

$$\frac{\partial}{\partial t}u(x, t) + \frac{\partial}{\partial x}f(u(x, t)) = 0. \quad (6.1)$$

Let  $a(u)$  be defined by  $a(u) = f'(u)$ . To calculate solutions of (6.1) numerically, we consider only conservative,  $(2k + 1)$ -point finite difference methods with 2 time levels, which are consistent with the conservation law (6.1). Let the function  $U_{\Delta t}$  be defined by (4.2), then the numerical scheme can be written as in (4.7). From now on such a method is simply denoted by

$$U_{\Delta t}(\cdot, t + \Delta t) = \mathcal{H}_{\Delta t}U_{\Delta t}(\cdot, t). \quad (6.2)$$

Let  $T > 0$  be a given constant. First some new concepts are introduced. For a given function  $u = u(x, t)$  the *total variation* over  $[0, T]$  is defined by

$$\begin{aligned} \text{TV}_T(u) = & \limsup_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \int_0^T \int_{-\infty}^{+\infty} |u(x + \varepsilon, t) - u(x, t)| dx dt \\ & + \limsup_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \int_0^T \int_{-\infty}^{+\infty} |u(x, t + \varepsilon) - u(x, t)| dx dt. \end{aligned} \quad (6.3)$$

The total variation over  $[0, T]$  of the function  $U_{\Delta t}$  is derived after substituting this function in (6.3) with  $T = N\Delta t$  for some integer  $N$ , which gives

$$\text{TV}_T(U_{\Delta t}) = \sum_{n=0}^N \sum_{i=-\infty}^{+\infty} \left\{ \Delta t |U_{i+1}^n - U_i^n| + \Delta x |U_i^{n+1} - U_i^n| \right\}.$$

Analogous to (6.3), the one-dimensional total variation at time  $t$  is defined by

$$\text{TV}(u(\cdot, t)) = \limsup_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \int_{-\infty}^{+\infty} |u(x + \varepsilon, t) - u(x, t)| dx. \quad (6.4)$$

The total variation of the function  $U_{\Delta t}$  at time  $t_n$  is defined by substituting this function in (6.4), which gives

$$\text{TV}(U_{\Delta t}(\cdot, t_n)) = \sum_{i=-\infty}^{+\infty} |U_{i+1}^n - U_i^n|.$$

Next, a convergence theorem is given, which is proved in [6].

**Theorem 6.1** *Let  $T > 0$  be a given constant and suppose that  $U_{\Delta t}$  is generated by the numerical method (6.2). Suppose that the method is entropy stable. If for each initial data  $u_0 = u(x, 0)$  there exist some  $k_0, R > 0$  such that*

$$\text{TV}(U_{\Delta t}(\cdot, t_n)) \leq R, \text{ for all } n, \Delta t \text{ with } \Delta t < k_0, n\Delta t \leq T, \quad (6.5)$$

*then the method is convergent (for  $\Delta t \rightarrow 0$ ) in the sense of bounded,  $L_1^{\text{loc}}$  convergence and its limit is the unique entropy stable solution of (6.1).*

In the remainder of this section we assume that (6.1) has initial data  $u(x, 0) = u_0(x)$ , such that the total variation of  $u_0$  is finite. An easy way to ensure that condition (6.5) is fulfilled, is to require that the total variation is nonincreasing as time evolves, so that the total variation of  $U_{\Delta t}$  at any time  $t > 0$  is bounded by the total variation of the initial data. This requirement gives rise to the following definition (cf. [6], [17]).

**Definition 6.2** *The numerical method (6.2) is called total variation diminishing (abbreviated TVD) if*

$$\text{TV}(U_{\Delta t}(\cdot, t_{n+1})) \leq \text{TV}(U_{\Delta t}(\cdot, t_n)),$$

*for all grid functions  $U_{\Delta t}(\cdot, t_n)$ .*

Thus, if a TVD method is used, then the following inequalities hold

$$\text{TV}(U_{\Delta t}(\cdot, t_n)) \leq \text{TV}(U_{\Delta t}(\cdot, 0)) \leq \text{TV}(u_0),$$

for all  $n \geq 1$ . Since the initial function  $u_0$  is assumed to have a finite total variation, (6.5) holds. Therefore, a TVD method is convergent. Another argument to consider TVD methods is that the exact solution to the scalar conservation law (6.1) has also this TVD property (cf. [6], [17]). Any weak solution of (6.1) satisfies

$$\text{TV}(u(\cdot, t_2)) \leq \text{TV}(u(\cdot, t_1)), \text{ for all } t_2 \geq t_1.$$

In [5] some examples of TVD methods are given.

It has been shown earlier that one difficulty associated with numerical approximations of discontinuous solutions is that oscillations may appear near a discontinuity. It can be proved that the exact solution does not have these oscillations. More precisely, if  $u$  is a weak solution of the scalar conservation law (6.1) with initial data  $u_0$  with finite total variation, then  $u$  has the following *monotonicity preserving* properties as a function of  $t$  (cf. [5]):

- (i) no new extrema in  $x$  are created;
- (ii) the value of a local minimum is nondecreasing and the value of a local maximum is nonincreasing.

Since the exact solution has this property, it seems natural to require that the numerical solution has this same property (cf. [17]).

**Definition 6.3** *The numerical method (6.2) is called monotonicity preserving if the following statement hold. If  $u_0$  is monotone (either nonincreasing or nondecreasing), then  $U_{\Delta t}(\cdot, t)$  is also monotone for all  $t > 0$ .*

In [26] it is shown that a linear  $(2k + 1)$ -point finite difference scheme

$$U_i^{n+1} = \sum_{j=-k}^k \alpha_j U_{i+j}^n$$

is monotonicity preserving if and only if  $\alpha_j \geq 0$  for all  $j$  with  $-k \leq j \leq k$ .

Another useful property of an entropy stable solution of (6.1) is given by the following theorem (cf. [12]).

**Theorem 6.4** *Suppose that  $u$  and  $v$  are two entropy stable solutions of (6.1). If  $u(\cdot, 0) - v(\cdot, 0) \in L_1$ , then*

$$\|u(\cdot, t_2) - v(\cdot, t_2)\|_1 \leq \|u(\cdot, t_1) - v(\cdot, t_1)\|_1, \quad (6.6)$$

for all  $t_1, t_2$ , with  $t_2 \geq t_1 \geq 0$ . Here  $\|\cdot\|_1$  denotes the  $L_1$ -norm in the space variable.

The property (6.6) is called  $L_1$ -contraction. In analogy to this, an  $L_1$ -contracting numerical method is defined as follows (cf. [17]).

**Definition 6.5** *The numerical method (6.2) is called  $L_1$ -contracting if, for any two  $U_{\Delta t}(\cdot, t_n)$  and  $V_{\Delta t}(\cdot, t_n)$  satisfying (6.2), for which  $U_{\Delta t}(\cdot, t_n) - V_{\Delta t}(\cdot, t_n)$  has compact support, the following inequality holds:*

$$\|U_{\Delta t}(\cdot, t_{n+1}) - V_{\Delta t}(\cdot, t_{n+1})\|_1 \leq \|U_{\Delta t}(\cdot, t_n) - V_{\Delta t}(\cdot, t_n)\|_1.$$

The last property of the entropy stable weak solution of (6.1) that is used is the following (cf. [14]).

**Theorem 6.6** *If  $u$  and  $v$  are two entropy stable solutions of (6.1) with initial data that satisfy  $v_0(x) \geq u_0(x)$  for all  $x$ , then the solutions  $u$  and  $v$  satisfy  $v(x, t) \geq u(x, t)$  for all  $x$  and  $t > 0$ .*

A numerical method which has the same property is called a *monotone method* and is defined as follows (cf. [17]).

**Definition 6.7** *The numerical method (6.2) is called monotone if the following statement holds*

$$V_{\Delta t}(x, t_n) \geq U_{\Delta t}(x, t_n) \text{ for all } x \Rightarrow V_{\Delta t}(x, t_{n+1}) \geq U_{\Delta t}(x, t_{n+1}) \text{ for all } x.$$

To prove the monotonicity of a method, it is sufficient to check whether the difference operator  $\mathcal{H}_{\Delta t}$  is a nondecreasing function of each argument (cf. [2]). Examples of monotone methods are the basic Godunov method, basic upwind method or the Lax-Friedrichs method. In [20] it is shown that every E-method (see Definition 4.15) is monotone.

The relations between all the concepts which are introduced in this subsection are given by the following theorem (cf. [17]).

**Theorem 6.8** *If the numerical method (6.2) is monotone, then it is  $L_1$ -contracting. A numerical method (6.2) which is  $L_1$  contracting, is always TVD, and furthermore, a numerical method (6.2) which is TVD, is always monotonicity preserving.*

These relations can be summarized as

$$\text{monotone} \Rightarrow L_1\text{-contracting} \Rightarrow \text{TVD} \Rightarrow \text{monotonicity preserving}.$$

An easy application of Theorem 6.1 and Theorem 6.8 shows that a monotone method converges. Monotone numerical methods have the satisfying property that we do not have to worry about entropy stability, since a monotone method contains enough numerical diffusion to converge always to the entropy stable solution. The following theorem shows this property (cf. [7]).

**Theorem 6.9** *If the numerical method (6.2) is monotone, then the method is convergent in the sense of bounded,  $L_1^{loc}$  convergence and its limit is the unique entropy stable solution of (6.1).*

Although the monotonicity requirement is easy to check and monotone methods always converge to the entropy stable solution, the class of monotone methods is seriously restricted as the following theorem shows (cf. [7]).

**Theorem 6.10** *A monotone numerical method is consistent of at most order one.*

A monotone method is not accurate enough in regions where the solution is smooth. Therefore, TVD methods are used more frequently. To derive a higher order TVD method is not trivial. In [30] it is shown that any 3-point TVD method is at most first order accurate. This shows that methods with more than 3 points are required to achieve second-order accuracy. Also in [30] a 5-point TVD method is derived, which is entropy stable and second order accurate in the regions where the solution is smooth. In the next subsection we will describe *flux limiter methods*. These methods are second order accurate in regions where the solution  $u$  is smooth.

## 6.2 Flux limiter methods

In the flux limiter approach, we choose a high order flux (e.g. the Lax-Wendroff flux) that works well in regions where the solution is smooth, and a low order flux (e.g. the flux from some monotone method) that behaves well near discontinuities. The main idea is the hybridization of these two flux-functions into a single flux in such a way that this single flux reduces to the high order flux in smooth regions and to the low order flux near discontinuities. This idea is elaborated in this subsection.

Let a conservative 3-point method be given, which is consistent with the conservation law (6.1). The corresponding finite difference scheme is given by (see (4.4))

$$U_i^{n+1} = U_i^n - \tau(F_{i+\frac{1}{2}}^{(E)} - F_{i-\frac{1}{2}}^{(E)}), \quad (6.7)$$

where  $\tau = \Delta t / \Delta x$  and  $F_{i+1/2}^{(E)} = F(U_i^n, U_{i+1}^n)$  denotes the numerical flux of some arbitrary E-method (see Definition 4.15). For shortness of notation we define  $\delta_+ y_i = \delta y_{i+1/2} = \delta_- y_{i+1} = y_{i+1} - y_i$ . Furthermore, the following flux differences are defined

$$\begin{aligned} (\delta f_{i+\frac{1}{2}}^n)^+ &= f(U_{i+1}^n) - F_{i+\frac{1}{2}}^{(E)}, \\ (\delta f_{i+\frac{1}{2}}^n)^- &= F_{i+\frac{1}{2}}^{(E)} - f(U_i^n). \end{aligned} \quad (6.8)$$

Note that  $(\delta f_{i+1/2}^n)^+ + (\delta f_{i+1/2}^n)^- = \delta f_{i+1/2}^n$ . These flux differences in turn are used to define the local CFL numbers,

$$\sigma_{i+\frac{1}{2}}^+ = \frac{\tau(\delta f_{i+\frac{1}{2}}^n)^+}{\delta U_{i+\frac{1}{2}}^n}, \quad \sigma_{i+\frac{1}{2}}^- = \frac{\tau(\delta f_{i+\frac{1}{2}}^n)^-}{\delta U_{i+\frac{1}{2}}^n}. \quad (6.9)$$

It is not difficult to see that the inequality (4.31), which defines an E-method, implies (cf. [5], [28])

$$\sigma_{i+\frac{1}{2}}^- \leq 0 \leq \sigma_{i+\frac{1}{2}}^+. \quad (6.10)$$

Let the following 3-point finite difference method be given to approximate (6.1) numerically,

$$U_i^{n+1} = U_i^n - (D_{i+\frac{1}{2}} \delta U_{i+\frac{1}{2}}^n - C_{i-\frac{1}{2}} \delta U_{i-\frac{1}{2}}^n), \quad (6.11)$$

where  $C_{i-1/2}$  and  $D_{i+1/2}$  are data dependent coefficients, i.e.  $C_{i-1/2} = C(U_i^n, U_{i-1}^n)$  and  $D_{i+1/2} = D(U_{i+1}^n, U_i^n)$ . In [5] the following theorem is proved, which gives sufficient conditions for the above method to be TVD.

**Theorem 6.11** *If the coefficients in (6.11) satisfy*

$$C_{i+\frac{1}{2}} \leq 0, \quad D_{i+\frac{1}{2}} \leq 0, \quad -(C_{i+\frac{1}{2}} + D_{i+\frac{1}{2}}) \leq 1,$$

*then the numerical method (6.11) is a TVD method.*

From (6.8) it is seen that

$$F_{i+\frac{1}{2}}^{(E)} - F_{i-\frac{1}{2}}^{(E)} = (\delta f_{i+\frac{1}{2}}^n)^- + (\delta f_{i-\frac{1}{2}}^n)^+ = \frac{1}{\tau}(\sigma_{i+\frac{1}{2}}^- \delta U_{i+\frac{1}{2}}^n + \sigma_{i-\frac{1}{2}}^+ \delta U_{i-\frac{1}{2}}^n),$$

and therefore, one possibility of writing a general scheme (6.7) in the form (6.11) is

$$U_i^{n+1} = U_i^n - (\sigma_{i+\frac{1}{2}}^- \delta U_{i+\frac{1}{2}}^n + \sigma_{i-\frac{1}{2}}^+ \delta U_{i-\frac{1}{2}}^n),$$

i.e. taking  $C_{i+1/2} = -\sigma_{i+1/2}^+$  and  $D_{i+1/2} = \sigma_{i+1/2}^-$ . Using (6.10) and Theorem 6.11 it is obvious that (6.7) is a TVD method, if it is an E-method and the CFL-like condition

$$\sigma_{i+\frac{1}{2}}^+ - \sigma_{i+\frac{1}{2}}^- \leq 1 \quad (6.12)$$

is fulfilled. Let  $F_{i+1/2}^{(LW)}$  and  $F_{i+1/2}^{(BU)}$  denote the numerical flux corresponding to, respectively, the Lax-Wendroff method and the basic upwind method, both applied to (6.1). It can be shown that the Lax-Wendroff flux can be rewritten as

$$F_{i+\frac{1}{2}}^{(LW)} = F_{i+\frac{1}{2}}^{(BU)} + \frac{1}{2}(1 - \sigma_{i+\frac{1}{2}}^+)(\delta f_{i+\frac{1}{2}}^n)^+ - \frac{1}{2}(1 + \sigma_{i+\frac{1}{2}}^-)(\delta f_{i+\frac{1}{2}}^n)^-,$$

(cf. [10], where the basic upwind flux is replaced by a general upwind flux). Hence, The Lax-Wendroff flux-function is composed by a first order basic upwind flux plus an additional flux, which is given by

$$\frac{1}{2}(1 - \sigma_{i+\frac{1}{2}}^+)(\delta f_{i+\frac{1}{2}}^n)^+ - \frac{1}{2}(1 + \sigma_{i+\frac{1}{2}}^-)(\delta f_{i+\frac{1}{2}}^n)^-. \quad (6.13)$$

This extra flux is often called an *antidiffusive flux*.

Since it is well known that the Lax-Wendroff scheme is not TVD, we try to remedy this by adding only a limited amount of the antidiffusive flux (6.13) to the first order scheme, i.e.

$$F_{i+\frac{1}{2}} = F_{i+\frac{1}{2}}^{(BU)} + \varphi(\theta_i^+) \frac{1}{2} (1 - \sigma_{i+\frac{1}{2}}^+) (\delta f_{i+\frac{1}{2}}^n)^+ - \varphi(\theta_{i+1}^-) \frac{1}{2} (1 + \sigma_{i+\frac{1}{2}}^-) (\delta f_{i+\frac{1}{2}}^n)^-, \quad (6.14)$$

where the function  $\varphi$  is called a *limiter*. To detect where the amount of the antidiffusive flux is large, the limiters are considered as functions of the following ratios (cf. [10], [28])

$$\theta_i^+ = \frac{(1 - \sigma_{i-\frac{1}{2}}^+) (\delta f_{i-\frac{1}{2}}^n)^+}{(1 - \sigma_{i+\frac{1}{2}}^+) (\delta f_{i+\frac{1}{2}}^n)^+}, \quad \theta_i^- = \frac{(1 + \sigma_{i+\frac{1}{2}}^-) (\delta f_{i+\frac{1}{2}}^n)^-}{(1 + \sigma_{i-\frac{1}{2}}^-) (\delta f_{i-\frac{1}{2}}^n)^-}. \quad (6.15)$$

The limiter  $\varphi$  is taken to be nonnegative, so that the sign of the antidiffusive flux (6.13) is maintained, i.e.

$$\varphi(\theta) \geq 0 \text{ for all } \theta. \quad (6.16)$$

Next (6.14) is generalized as follows

$$F_{i+\frac{1}{2}}^{(FL)} = F_{i+\frac{1}{2}}^{(E)} + \varphi(\theta_i^+) \frac{1}{2} (1 - \sigma_{i+\frac{1}{2}}^+) (\delta f_{i+\frac{1}{2}}^n)^+ - \varphi(\theta_{i+1}^-) \frac{1}{2} (1 + \sigma_{i+\frac{1}{2}}^-) (\delta f_{i+\frac{1}{2}}^n)^-. \quad (6.17)$$

This is a generalization, since the basic upwind method is a particular example of an E-method. The numerical method defined by the flux (6.17) is called a *flux limiter method*. An easy calculation, using Taylor series expansions shows that this flux defines a numerical method, which is second order consistent in space if  $\varphi = 1$ .

If we want to apply Theorem 6.11, then the numerical method given by the flux (6.17), has to be rewritten in the same form as (6.11). One possibility is to take  $C_{i+1/2}$  and  $D_{i+1/2}$  as

$$\begin{aligned} C_{i+\frac{1}{2}} &= -\sigma_{i+\frac{1}{2}}^+ \left\{ 1 + \frac{1}{2} (1 - \sigma_{i+\frac{1}{2}}^+) (\varphi(\theta_{i+1}^+)/\theta_{i+1}^+ - \varphi(\theta_i^+)) \right\} \\ D_{i+\frac{1}{2}} &= \sigma_{i+\frac{1}{2}}^- \left\{ 1 + \frac{1}{2} (1 + \sigma_{i+\frac{1}{2}}^-) (\varphi(\theta_i^-)/\theta_i^- - \varphi(\theta_{i+1}^-)) \right\}. \end{aligned}$$

Suppose that there exists a constant  $\Phi$ , with  $0 < \Phi \leq 2$  such that

$$\left| \frac{\varphi(\theta_i^\pm)}{\theta_i^\pm} - \varphi(\theta_{i-1}^\pm) \right| \leq \Phi. \quad (6.18)$$

If the following CFL-like condition is satisfied (see (6.12))

$$\sigma_{i+\frac{1}{2}}^+ - \sigma_{i+\frac{1}{2}}^- \leq \frac{2}{2 + \Phi}, \quad (6.19)$$

then all assumptions of Theorem 6.11 are fulfilled and the method is TVD. If in addition to (6.16) it is also required that

$$\varphi(\theta) = 0 \text{ for all } \theta \leq 0,$$

then the bound (6.18) reduces to

$$0 \leq \frac{\phi(\theta)}{\theta} \leq \Phi, \quad 0 \leq \phi(\theta) \leq \Phi \text{ for all } \theta. \quad (6.20)$$

The last condition on the limiter is given by the following theorem (cf. [21]).

**Theorem 6.12** *The flux limiter method with flux (6.17) is consistent with the conservation law (6.1) provided  $\varphi$  is a bounded function. It is a second order TVD method (on smooth solutions with  $\partial u/\partial x$  bounded away from zero) provided  $\varphi$  satisfies (6.20),  $\varphi(1) = 1$  and  $\varphi$  is Lipschitz continuous at  $\theta = 1$ .*

In [28] it appears that the best choice for  $\varphi$  is a convex combination of 1 and  $\theta$ , i.e.

$$\varphi(\theta) = 1 + \zeta(\theta)(\theta - 1), \quad (6.21)$$

with  $0 \leq \zeta(\theta) \leq 1$  for all  $\theta$ . Other choices apparently give too much compression, i.e. smooth initial data such as a sine wave tends to turn into a square wave as time evolves (cf. [28]). Note that with this choice of  $\varphi$  the condition  $\varphi(1) = 1$  is automatically satisfied.

**Example 6.13** Roe chooses  $\varphi(\theta)$  as large as possible such that all conditions of Theorem 6.12 are fulfilled. This limiter is called the *superbee limiter* and is given by (cf. [17]).

$$\varphi(\theta) = \max(0, \min(1, 2\theta), \min(\theta, 2)).$$

A smoother limiter function is introduced by van Leer (cf. [15]) and is given by

$$\varphi(\theta) = \frac{|\theta| + \theta}{1 + |\theta|}.$$

In Figure 5, some numerical results of van Leer's limiter and Roe's superbee limiter are given. In this results the underlying E-method is simply the basic upwind method.

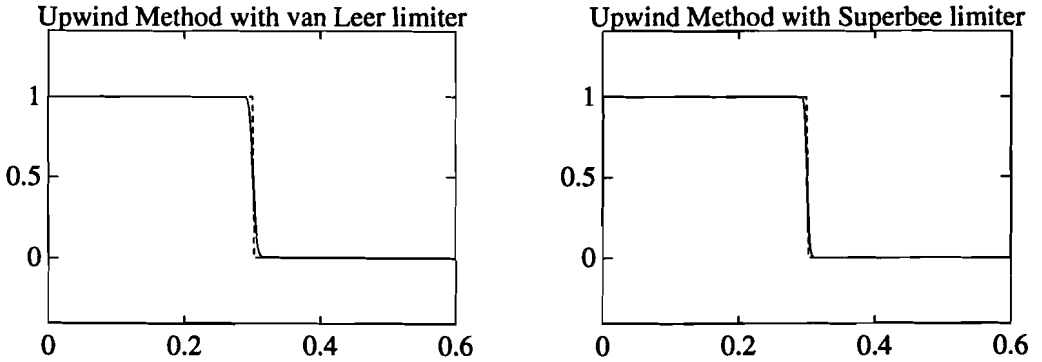


Figure 5: Numerical solution (solid line) and exact solution (dashed line) of (6.1) at  $t = 0.3$  with  $f(u) = u$ ,  $\Delta t = 0.002$ ,  $\sigma = 0.8$  and the initial condition  $u(x, 0) = 1$  if  $x < 0$  and  $u(x, 0) = 0$  if  $x > 0$ .

In [28] some other examples of limiters and the corresponding numerical results are given.

Two questions remain open. Is the flux limiter method an entropy stable method, and can flux limiter methods be extended to a system of nonlinear conservation laws. We refer to [21], where a particular flux limiter method is described for systems. Further, in the scalar case it is proved, that this method converges to the unique entropy



stable solution. Another way to generalize the scalar flux limiter method to systems of equations is to linearize the nonlinear system (cf. [17]). The generalization is then obtained by diagonalizing the system and applying the flux limiter method to each of the resulting scalar equations.

**Acknowledgement.** We wish to thank prof. R.M.M. Mattheij for critically reading the manuscript.

## References

- [1] Ames, W.F., Numerical Methods for Partial Differential Equations, Academic Press, New York (1977).
- [2] Crandall, M.G. and A. Majda, *Monotone difference approximations for scalar conservation laws*, Math. Comp. **34** (1980), pp. 1-21.
- [3] Engquist, B. and S. Osher, *One-sided difference approximations for nonlinear conservation laws*, Math. Comp. **36** (1981), pp. 321-351.
- [4] Garabedian, P.R., Partial Differential Equations, J. Wiley, New York (1964).
- [5] Harten, A., *High resolution schemes for hyperbolic conservation laws*, J. Comput. Phys. **49** (1983), pp. 357-393.
- [6] Harten, A., *On a class of high resolution total-variation-stable finite-difference schemes*, SIAM J. Numer. Anal. **21** (1984), pp. 1-23.
- [7] Harten, A., J.M. Hyman and P.D. Lax, *On finite-difference approximations and entropy conditions for shocks*, Comm. Pure Appl. Math. **29** (1976), pp. 297-322.
- [8] Harten, A., P.D. Lax and B. van Leer, *On upstream differencing and Godunov-type schemes for hyperbolic conservation laws*, SIAM Review **25** (1983), pp. 35-61.
- [9] Hirsch, C., Numerical Computation of Internal and External Flows, volume 1: Fundamentals of Numerical Discretization, J. Wiley, Chichester (1988).
- [10] Hirsch, C., Numerical Computation of Internal and External Flows, volume 2: Computational Methods for Inviscid and Viscous Flows, J. Wiley, Chichester (1990).
- [11] Jeffrey, A. and T. Taniuti, Non-linear Wave Propagation, Academic Press, New York (1964).
- [12] Kreiss, H.O., *Difference approximations for initial-boundary value problems for hyperbolic differential equations*, in Numerical Solutions of Nonlinear Partial Differential Equations, D. Greenspan (ed.), J. Wiley, New York, 1966, pp. 140-166.
- [13] Krushkov, S.N., *First order quasi-linear equations in several independent variables*, Math. USSR Sb. **10** (1970), pp. 217-243.
- [14] Lax, P.D., Hyperbolic Systems of Conservation Laws and the Mathematical Theory of Shock Waves, SIAM Regional Conference Series in Applied Mathematics, volume 11, SIAM, Philadelphia (1973).
- [15] Leer, B. van, *Towards the ultimate conservative difference scheme II. Monotonicity and conservation combined in a second order scheme*, J. Comput. Phys. **14** (1974), pp. 361-370.
- [16] Leer, B. van, *On the relation between the upwind-differencing schemes of Godunov, Engquist-Osher and Roe*, SIAM J. Sci. Stat. Comput. **5** (1984), pp. 1-20.

- [17] LeVeque, R.J., *Numerical Methods for Conservation Laws*, Lectures in Mathematics, Birkhäuser Verlag, Basel (1990).
- [18] Liepmann, H.W. and A. Roshko, *Elements of Gas dynamics*, J. Wiley, London (1957).
- [19] Majda, A. and S. Osher, *Numerical viscosity and the entropy condition*, Comm. Pure Appl. Math. **32** (1979), pp. 797-838.
- [20] Osher, S., *Riemann solvers, the entropy condition, and difference approximations*, SIAM J. Numer. Anal. **21** (1984), pp. 217-235.
- [21] Osher, S. and S. Chakravarthy, *High resolution schemes and the entropy condition*, SIAM J. Numer. Anal. **21** (1984), pp. 955-984.
- [22] Osher, S. and F. Solomon, *Upwind difference schemes for hyperbolic systems of conservation laws*, Math. Comp. **38** (1982), pp. 339-374.
- [23] Osher, S. and E. Tadmor, *On the convergence of difference approximations to scalar conservation laws*, Math. Comp. **50** (1988), pp. 19-51.
- [24] Roe, P.L., *Approximate Riemann solvers, parameter vectors and difference schemes*, J. Comput. Phys. **43** (1981), pp. 357-372.
- [25] Smoller, J., *Shock Waves and Reaction-Diffusion Equations*, volume 258 of Grundlehren der mathematischen Wissenschaften, Springer-Verlag, New York (1983).
- [26] Spekreijse, S.P., *Multigrid Solution of the Steady Euler Equations*, volume 46 of CWI tracts, CWI, Amsterdam (1988).
- [27] Steger, J.L. and R.F. Warming, *Flux vector splitting of the inviscid gas dynamic equations with applications to finite-difference methods*, J. Comput. Phys. **40** (1981), pp. 263-293.
- [28] Sweby, P.K., *High resolution schemes using flux limiters for hyperbolic conservation laws*, SIAM J. Numer. Anal. **21** (1984), pp. 995-1011.
- [29] Tadmor, E., *The numerical viscosity of entropy stable schemes for systems of conservation laws 1*, Math. Comp. **49** (1987), pp. 91-103.
- [30] Vila, J.P., *High-order schemes and entropy condition for nonlinear hyperbolic systems of conservation laws*, Math. Comp. **50** (1988), pp. 53-73.
- [31] Warming, R.F. and B.J. Hyett, *The modified equation approach to the stability and accuracy analysis of finite difference methods*, J. Comput. Phys. **14** (1974), pp. 159-179.
- [32] Whitham, G.B., *Linear and Nonlinear Waves*, J. Wiley, New York (1974).