

Numerical schemes for systems of conservation laws in 1D

In this part, we are interested in the numerical approximation by finite volume methods of the solution of the 1D problem:

$$\begin{cases} \partial_t u + \partial_x f(u) = 0, & x \in \mathbb{R}, t > 0, \\ u(x, 0) = u_0(x), \end{cases}$$

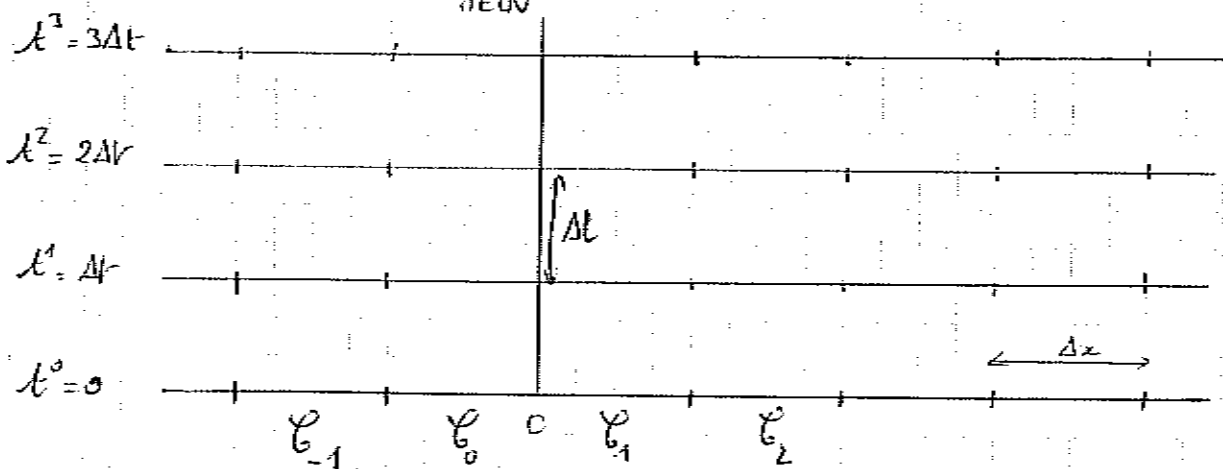
where $u = (u_1, \dots, u_p)^t \in \Omega \subset \mathbb{R}^p$ is the unknown and $f(u) \in \mathbb{R}^p$ is the flux function.

Let us first briefly recall the design principle of the finite volume approach.

We begin by introducing a time step Δt and a space step Δx that we assume to be constant for simplicity. We set $\lambda = \Delta t / \Delta x$ and define the mesh interfaces $x_{j+1/2} = j \Delta x$ for $j \in \mathbb{Z}$, and the intermediate times $t^n = n \Delta t$ for $n \in \mathbb{N}$. We thus have

$$\mathbb{R} = \bigcup_{j \in \mathbb{Z}} \mathcal{C}_j \quad \text{with} \quad \mathcal{C}_j = [x_{j-1/2}, x_{j+1/2}]$$

$$\text{and} \quad \mathbb{R}^+ = \bigcup_{n \in \mathbb{N}} [t^n, t^{n+1}]$$



At each time t^n , we look for an approximation u_j^n of the mean value $\frac{1}{\Delta x} \int_{\mathcal{C}_j} u(x, t^n) dx$ of the exact solution at time t^n on the cell \mathcal{C}_j .

We assume that the initial condition u_0 is known, and set for instance for $n=0$

$$u_j^0 = \frac{1}{\Delta x} \int_{\xi_j} u_0(x) dx, \quad \forall j \in \mathbb{Z}.$$

Being given the sequence $(u_j^n)_{j \in \mathbb{Z}}$ at time t^n , the main point is now to propose a definition of the sequence $(u_j^{n+1})_{j \in \mathbb{Z}}$ by a recurrence relation. To do so, we start from the governing equations to be solved, namely

$$\partial_t u + \partial_x f(u) = 0$$

and integrate them on the volume $\xi_j \times [t^n, t^{n+1}] = [x_{j-1/2}, x_{j+1/2}] \times [t^n, t^{n+1}]$

$$\int_{t^n}^{t^{n+1}} \int_{\xi_j} (\partial_t u + \partial_x f(u)) dx dt = 0,$$

which gives

$$\int_{\xi_j} u(x, t^{n+1}) dx - \int_{\xi_j} u(x, t^n) dx + \int_{t^n}^{t^{n+1}} f(u(x_{j+1/2}, t)) dt - \int_{t^n}^{t^{n+1}} f(u(x_{j-1/2}, t)) dt = 0.$$

Writing now

$$\int_{\xi_j} u(x, t^{n+1}) dx \approx \Delta x u_j^{n+1}, \quad j \in \mathbb{Z},$$

$$\int_{\xi_j} u(x, t^n) dx \approx \Delta x u_j^n, \quad j \in \mathbb{Z},$$

we then deduce the following form for a natural finite volume approximation of the solution of (3):

$$u_j^{n+1} = u_j^n - \frac{\Delta t}{\Delta x} \left(\frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f(u(x_{j+1/2}, t)) dt - \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f(u(x_{j-1/2}, t)) dt \right)$$

where $\frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f(u(x_{j+1/2}, t)) dt$ is the so-called numerical flux and represents an approximation of the exact flux that passes through the interface $x_{j+1/2}$ in the time interval $[t^n, t^{n+1}]$.

The choice of the numerical flux $f_{j+1/2}$ at each interface $x_{j+1/2}$ defines the numerical scheme under consideration. In the following, we propose to study several numerical schemes, focusing on the following form for the numerical flux:

$$\forall j: f_{j+1/2} := f(u_j^n, u_{j+1}^n) \quad (6)$$

The function $f(\cdot, \cdot)$ is called the numerical flux function and this leads to the following update formula:

$$u_j^{n+1} = u_j^n - \frac{\Delta t}{\Delta x} (f(u_j^n, u_{j+1}^n) - f(u_{j-1}^n, u_j^n)), \quad j \in \mathbb{Z}. \quad (7)$$

We observe that the definition of u_j^{n+1} only depends on the three values $u_{j-1}^n, u_j^n, u_{j+1}^n$. (7) is said to be a three-point explicit finite volume numerical scheme.

Definition

The scheme (7) is said to be consistent with (3) if $f(\cdot, \cdot)$ satisfies

$$f(u, u) = f(u) \quad \forall u \in \Omega.$$

Let us also briefly recall that we are basically interested in this course in entropy weak solutions of (1), that is weak solutions of (1) satisfying in addition the entropy condition

$$\partial_t S(u) + \partial_x G(u) \leq 0$$

in the distributional sense, where $(S(u), G(u))$ represents an entropy - entropy flux pair for (1) with $S: \Omega \rightarrow \mathbb{R}$ a convex function. We then introduce the following definition.

Definition

The scheme (7) is said to be consistent with (3)-(5) if (8) holds true and if there exists a numerical entropy flux function $G(\cdot, \cdot)$,

consistent with $G(\cdot)$ in the sense that

$$G(u, u) = G(u) \quad \forall u \in \Omega$$

and such that

$$(M) \quad S(u_j^{(n)}) - S(u_j^0) - \frac{\Delta t}{\Delta x} (G_j(u_j^0, u_{j+1}^0) - G_j(u_{j-1}^0, u_j^0)) \leq 0 \quad \forall j \in \mathbb{Z}$$

To conclude this introduction, we briefly recall the Lax-Wendroff theorem that asserts that when a consistent scheme of the form (7) converges in some sense to a function u , the limit is a weak solution of (3). In addition, if (7) is entropy consistent then the limit is an entropy weak solution of (3)-(9).

The Godunov method

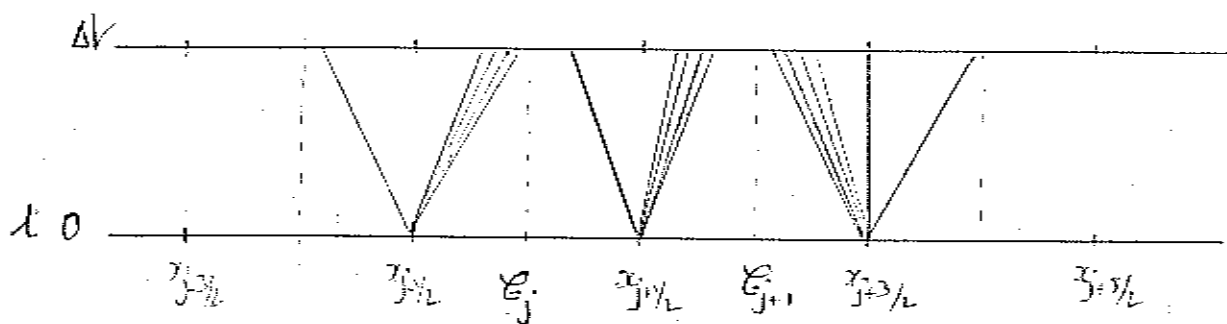
We begin with the celebrated Godunov's method. Let us define the piecewise constant approximate solution $x \rightarrow u_x(x, t^n)$ at time t^n given by $u_x(x, t^n) = u_j^0 \quad \forall x \in \mathcal{C}_j, j \in \mathbb{Z}, n \in \mathbb{N}$. The method is made of two steps.

Step 1: Evolution in time

In this first step, one solves the Cauchy problem (1)-(9) exactly with initial data $u(x) = u_x(x, t^n)$ and for times $t \in [0, \Delta t]$. Locally around each interface $x_{j+1/2}$, it is clear that this amounts to solve a Riemann problem. More precisely, one can even claim that under the so-called CFL condition involving the characteristic speeds $\lambda_i, i=1, \dots, p$ of (3)

$$\frac{\Delta t}{\Delta x} \max_u \{ |\lambda_i(u)|, i=1, \dots, p \} \leq 1/2, \quad (11)$$

for all the u under consideration, the solution of this Cauchy problem is known by glueing together the solutions of the Riemann problems set at each interface, see figure below.



More precisely, this solution is given by

$$\tilde{u}(x,t) = u_{\tilde{u}}\left(\frac{x-x_{j+1/2}}{t}; u_j^0, u_{j+1}^0\right) \text{ for all } (x,t) \in [x_j, x_{j+1}] \times [0, \Delta t],$$

where $x_j = \frac{1}{2}(x_{j-1/2} + x_{j+1/2})$ for all $j \in \mathbb{Z}$ and where $(x,t) \rightarrow u_{\tilde{u}}(x/t, u_L, u_R)$ denotes the self-similar solution of the Riemann problem (1)-(3)

$$\text{with initial data } u_{\tilde{u}}(x) = \begin{cases} u_L & \text{if } x < 0 \\ u_R & \text{if } x > 0. \end{cases}$$

Step 2: Projection

In this second step, we get back a piecewise constant approximate solution on each cell \mathcal{E}_j at time t^{n+1} by averaging the solution $u(x, \Delta t)$:

$$u_j^{n+1} = \frac{1}{\Delta x} \int_{\mathcal{E}_j} \tilde{u}(x, \Delta t) dx, \quad j \in \mathbb{Z}$$

Then, the process repeats.

In practice, this algorithm is considerably simplified by observing that the cell average (16) can be easily computed using the integral form of the governing equation (3), that is $\tau u + 2\tau |u| = 0$. Indeed, since (13) is assumed to be an exact solution, integrating over $\mathcal{E}_j \times [0, \Delta t]$

we know that

$$\int_{\mathcal{E}_j} \tilde{u}(x, \Delta t) dx - \int_{\mathcal{E}_j} \tilde{u}(x, 0) dx + \int_{t^n}^{t^{n+1}} \tau(\tilde{u}(x_{j+1/2}, t)) dt - \int_{t^n}^{t^{n+1}} \tau(\tilde{u}(x_{j-1/2}, t)) dt = 0$$

that is

$$u_j^{n+1} = u_j^n - \frac{\Delta t}{\Delta x} (\tau_{j+1/2} - \tau_{j-1/2}), \quad j \in \mathbb{Z}$$

where the numerical flux $f_{j+1/2}^{n+1} = f(u_j^n, u_{j+1}^n)$ is given by

$$\begin{aligned} f_{j+1/2}^{n+1} &= f(u_j^n, u_{j+1}^n) = \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f(\tilde{u}(x_{j+1/2}, t)) dt \\ &= \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f(u_{j+1/2}(0; u_j^n, u_{j+1}^n)) dt \end{aligned}$$

that is

$$f_{j+1/2}^{n+1} = f(u_j^n, u_{j+1}^n) = f(u_{j+1/2}(0; u_j^n, u_{j+1}^n)), \quad j \in \mathcal{I}$$

This shows that the Godunov method can be written in the form (7).

It is also clear that the method is consistent with (3) since

$$f(u, u) = f(u_{j+1/2}(0; u, u)) = f(u).$$

As far as the entropy consistency is concerned, since (3) is assumed to satisfy the entropy inequality $\partial_t S(u) + \partial_x G(u) \leq 0$, integrating over $\mathcal{I}_j \times [0, \Delta t]$ leads to

$$\int_{\mathcal{I}_j} S(\tilde{u}(x, \Delta t)) dx - \int_{\mathcal{I}_j} S(\tilde{u}(x, 0)) dx + \int_{t^n}^{t^{n+1}} G(\tilde{u}(x_{j+1/2}, t)) dt - \int_{t^n}^{t^{n+1}} G(\tilde{u}(x_{j-1/2}, t)) dt \leq 0$$

that is

$$\frac{1}{\Delta x} \int_{\mathcal{I}_j} S(\tilde{u}(x, \Delta t)) dx - S(u_j^n) + \frac{\Delta t}{\Delta x} (G(u_j^n, u_{j+1}^n) - G(u_{j-1}^n, u_j^n)) \leq 0 \quad (16)$$

$$\text{with } G_{j+1/2} = G(u_j^n, u_{j+1}^n) = G(u_{j+1/2}(0; u_j^n, u_{j+1}^n)), \quad j \in \mathcal{I} \quad (17)$$

which is clearly consistent with G since $G(u, u) = G(u_{j+1/2}(0; u, u)) = G(u)$.

But invoking the Jensen inequality and the convexity of $u \rightarrow S(u)$, we have

$$S\left(\frac{1}{\Delta x} \int_{\mathcal{I}_j} \tilde{u}(x, \Delta t) dx\right) \leq \frac{1}{\Delta x} \int_{\mathcal{I}_j} S(\tilde{u}(x, \Delta t)) dx$$

$$\text{that is } S(u_j^{n+1}) \leq \frac{1}{\Delta x} \int_{\mathcal{I}_j} S(\tilde{u}(x, \Delta t)) dx$$

Combining (16) and (18) proves that Godunov's method is consistent with (3)-(9).

The main drawback of Godunov's method is that it requires at each time iteration and at each interface $x_{j+1/2}$, the resolution of exact Riemann problems. Although these Riemann problems can be solved in theory, in practice it necessitates the resolution of nonlinear equations which can be expensive. We note also that most of the structure of the exact Riemann solutions is not used in Godunov's method since an averaging process is involved. This discussion suggests that it is probably not worthwhile and relevant calculating the Riemann solutions exactly, instead of that defining approximate Riemann solutions could be less expensive and give equally good numerical results, provided that these approximate solutions are consistent in some "averaged" sense.

Approximate Riemann solvers

Our objective is to propose an approximation of $(x, t) \rightarrow u_R(x/t, u_L, u_R)$. In this course, we will focus on simple approximate Riemann solvers of the following form

$$\tilde{u}_R(x/t, u_L, u_R) = \begin{cases} u_L & \text{if } x/t < \tilde{\lambda}_1 \\ u_k & \text{if } \tilde{\lambda}_{k-1} < x/t < \tilde{\lambda}_k, \quad k=2, \dots, l \\ u_R & \text{if } x/t > \tilde{\lambda}_l, \end{cases}$$

with u_k and $\tilde{\lambda}_k = \tilde{\lambda}_k(u_L, u_R)$, $k=1, \dots, l$, L-LP to be defined.

Note that this approximate Riemann solver is self-similar.

Following the general theory of Harten-Lax and van Leer, a suitable approximate Riemann solver is subject to the following consistency properties.

Definition

The simple approximate Riemann solver (19) is said to be consistent with the integral form of (3) if and only if for Δx and Δt such that

$$\max_{1 \leq k \leq l} |\tilde{\lambda}_k(u_L, u_R)| \frac{\Delta t}{\Delta x} \leq 1/2,$$

we have

$$\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \tilde{u}_k(x/\Delta t; u_L, u_R) dx = \frac{u_L + u_R}{2} - \frac{\Delta t}{\Delta x} (f(u_R) - f(u_L)),$$

(and if $\tilde{u}_k(x/\Delta t; u, u) = u \quad \forall (x, t)$)

Where does this equality come from and what does it mean? Actually, starting from $f(u) + f'(u) = 0$ and integrating over $[-\frac{\Delta x}{2}, \frac{\Delta x}{2}] \times [0, \Delta t]$ gives for the exact Riemann solution $u_k(\cdot, u_L, u_R)$

$$\begin{aligned} \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} u_k(x/\Delta t; u_L, u_R) dx &= \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} u_0(x) dx - \frac{\Delta t}{\Delta x} (f(u_R) - f(u_L)) \\ &= \frac{u_L + u_R}{2} - \frac{\Delta t}{\Delta x} (f(u_R) - f(u_L)) \end{aligned}$$

since u_0 is given by $u_0(x) = \begin{cases} u_L & \text{if } x < 0 \\ u_R & \text{if } x > 0 \end{cases}$.

Then, (21) says that the exact and approximate Riemann solutions may be different "in the details", but must coincide "in average" over a cell $[-\frac{\Delta x}{2}, \frac{\Delta x}{2}]$ with length Δx and centered on the position of the initial discontinuity.

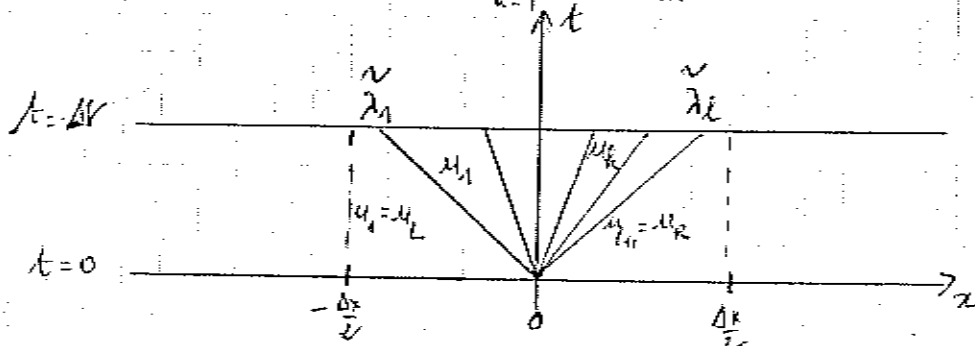
Remark

Since we clearly have (see figure below)

$$\begin{aligned} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \tilde{u}_k(x/\Delta t; u_L, u_R) dx &= \left(\frac{\Delta x}{2} + \tilde{\lambda}_1 \Delta t\right) u_L + (\tilde{\lambda}_2 - \tilde{\lambda}_1) \Delta t u_k + \dots + (\tilde{\lambda}_l - \tilde{\lambda}_{l-1}) \Delta t u_k \\ &\quad + \left(\frac{\Delta x}{2} - \tilde{\lambda}_l \Delta t\right) u_R \\ &= \frac{\Delta x}{2} (u_L + u_R) - \Delta t \sum_{k=1}^l \tilde{\lambda}_k (u_{k+1} - u_k) \end{aligned}$$

the consistency relation (21) also writes

$$f(u_R) - f(u_L) = \sum_{k=1}^l \tilde{\lambda}_k (u_{k+1} - u_k)$$



Similarly, we introduce the notion of consistency with the entropy inequality $\int S(u) + \int G(u) \leq 0$. integral form of the

Definition

The simple approximate Riemann solver (19) is said to be consistent with the integral form of the entropy inequality (3) if and only if for Δx and Δt satisfying (20) we have

$$\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} S(\tilde{u}_n(\frac{x}{\Delta}, u_L, u_R)) dx \leq \frac{S(u) + S(u_n)}{2} - \frac{\Delta t}{\Delta x} (G(u_n) - G(u_L)),$$

or equivalently

$$S(u_n) - G(u_L) \leq \frac{\ell}{2} \sum_{k=1}^n \tilde{\lambda}_k (S(u_{k+1}) - S(u_k)).$$

The associated Godunov-type method

The idea is to follow the same approach as in the original Godunov-method but replacing the exact Riemann solution $u_k(\cdot; u_L, u_R)$ with the approximate one $\tilde{u}_k(\cdot; u_L, u_R)$. Then we get:

Step 1: Approximate evolution in time

One solves approximately the Cauchy problem (1)-(3) with initial condition $u_0(x) = u_j(x, t^0)$ for times $t \in [0, \Delta t]$ with Δt such that

$$\frac{\Delta t}{\Delta x} \max_{k=1, \dots, \ell} |\lambda_k(u_j^0, u_{j+1}^0)| < \frac{1}{2} \quad (23)$$

and using the approximate Riemann solution, we get

$$\tilde{u}(x, t) = \tilde{u}_n\left(\frac{x - x_j^0 + t}{\Delta t}, u_j^0, u_{j+1}^0\right) \text{ for all } (x, t) \in [x_j^0, x_{j+1}^0] \times [0, \Delta t].$$

Step 2: Projection

This step is unchanged and we set

$$u_j^{n+1} = \frac{1}{\Delta x} \int_{I_j} \tilde{u}(x, \Delta t) dx, \quad j \in \mathbb{Z},$$

and the process repeats.

At this stage, it is natural to wonder if such a Godunov-type scheme admits an equivalent formulation of the form (7) with a given numerical flux $f_{j+1/2}$. Note that since the exact Riemann solution has been replaced with an approximate one, the calculations carried out in the context of the Godunov method and leading to (8) are not valid anymore here. However, we can write

$$\begin{aligned}
 u_j^{n+1} &= \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} \tilde{u}(x, \Delta t) dx = \\
 &= \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_j} \tilde{u}_n \left(\frac{x-x_{j-1/2}}{\Delta t}; u_{j-1}^n, u_j^n \right) dx + \frac{1}{\Delta x} \int_{x_j}^{x_{j+1/2}} \tilde{u}_n \left(\frac{x-x_{j+1/2}}{\Delta t}; u_j^n, u_{j+1}^n \right) dx \\
 &= \frac{1}{\Delta x} \int_0^{\frac{\Delta x}{2}} \tilde{u}_n \left(\frac{z}{\Delta t}; u_{j-1}^n, u_j^n \right) dz + \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^0 \tilde{u}_n \left(\frac{z}{\Delta t}; u_j^n, u_{j+1}^n \right) dz \\
 &= \frac{1}{2} \left(u^+(u_{j-1}^n, u_j^n) + u^-(u_j^n, u_{j+1}^n) \right)
 \end{aligned}$$

where we have set

$$\begin{cases}
 u^-(u_L, u_R) = \frac{2}{\Delta x} \int_{-\frac{\Delta x}{2}}^0 \tilde{u}_n \left(\frac{z}{\Delta t}; u_L, u_R \right) dz, \\
 u^+(u_L, u_R) = \frac{2}{\Delta x} \int_0^{\frac{\Delta x}{2}} \tilde{u}_n \left(\frac{z}{\Delta t}; u_L, u_R \right) dz.
 \end{cases}$$

Using the form (9) of the approximate Riemann solver, we can easily check that

$$\begin{aligned}
 u^-(u_L, u_R) &= u_L - \frac{2\Delta t}{\Delta x} \sum_{k=1}^l \tilde{\lambda}_k^- (u_{k+1} - u_L) & x^+ = \max(x, 0) \\
 u^+(u_L, u_R) &= u_R - \frac{2\Delta t}{\Delta x} \sum_{k=1}^l \tilde{\lambda}_k^+ (-u_{k+1} - u_R) & x^- = \min(x, 0)
 \end{aligned}$$

so that, with clear notations

$$u_j^{n+1} = u_j^n - \frac{\Delta t}{\Delta x} \left(\sum_{k=1}^l \tilde{\lambda}_{k, j+1/2}^- (u_{k, j+1}^{j+1/2} - u_k^{j+1/2}) + \sum_{k=1}^l \tilde{\lambda}_{k, j-1/2}^+ (u_{k, j-1}^{j-1/2} - u_k^{j-1/2}) \right)$$

Let us assume that the approximate Riemann solver under consideration is consistent with the integral form of (3).

Using the property $\begin{cases} x = x^+ + x^- \\ |x| = x^+ - x^- \end{cases}$, we can write

$$\begin{aligned} \sum_{k=1}^l \tilde{\lambda}_{j+\frac{1}{2}}^- (u_{k+1}^{j+\frac{1}{2}} - u_k^{j+\frac{1}{2}}) &= \frac{1}{2} \sum_{k=1}^l (\tilde{\lambda}_{k,j+\frac{1}{2}}^- - |\tilde{\lambda}_{k,j+\frac{1}{2}}^-|) (u_{k+1}^{j+\frac{1}{2}} - u_k^{j+\frac{1}{2}}) \\ &= \frac{1}{2} (f(u_{j+1}) - f(u_j)) - \frac{1}{2} \sum_{k=1}^l |\tilde{\lambda}_{k,j+\frac{1}{2}}^-| (u_{k+1}^{j+\frac{1}{2}} - u_k^{j+\frac{1}{2}}) \end{aligned}$$

Similarly we get

$$\sum_{k=1}^l \tilde{\lambda}_{j-\frac{1}{2}}^+ (u_{k+1}^{j-\frac{1}{2}} - u_k^{j-\frac{1}{2}}) = \frac{1}{2} (f(u_j) - f(u_{j+1})) + \frac{1}{2} \sum_{k=1}^l |\tilde{\lambda}_{k,j-\frac{1}{2}}^+| (u_{k+1}^{j-\frac{1}{2}} - u_k^{j-\frac{1}{2}})$$

so that we eventually have

$$u_j^{n+1} = u_j^n - \frac{\Delta t}{\Delta x} \left(\hat{f}_{j+\frac{1}{2}} - \hat{f}_{j-\frac{1}{2}} \right), \quad j \in \mathbb{Z}$$

with

$$\hat{f}_{j+\frac{1}{2}} = f(u_j^n, u_{j+1}^n) = \frac{1}{2} (f(u_j) + f(u_{j+1})) - \frac{1}{2} \sum_{k=1}^l |\tilde{\lambda}_{k,j+\frac{1}{2}}^-| (u_{k+1}^{j+\frac{1}{2}} - u_k^{j+\frac{1}{2}}), \quad j \in \mathbb{Z}.$$

The proposed Godunov-type method can then be written in the form (7) and is clearly consistent with (3) (keep in mind that we have assumed that the approximate Riemann solver is consistent in the integral sense to get this result).

Let us now focus on the validity of an entropy inequality similar to (11). For that, we assume that the approximate Riemann solver under consideration is consistent with the integral form of the entropy inequality (9), that is (22)' is valid:

$$G(u_R) - G(u_L) \leq \sum_{k=1}^l \tilde{\lambda}_k (S(u_{k+1}) - S(u_k)).$$

Let us first recall that

$$u_j^{n+1} = \frac{1}{2} (u^+(u_{j-1}^n, u_j^n) + u^-(u_j^n, u_{j+1}^n))$$

so that by a convexity argument we have

$$S(u_j^{n+1}) \leq \frac{1}{2} (S(u_{j-1}^n, u_j^n) + S(u_{j+1}^n, u_j^n)).$$

The next calculations now follow the same lines as before. Namely, we first note that using the form (9) of the approximate Riemann solver we can check that (using the Jensen inequality)

$$S(u_L, u_R) \leq \frac{2}{\Delta x} \int_{u_L}^{u_R} \tilde{u} \left(\frac{x}{\Delta t}, u_L, u_R \right) dx = S(u_L) - \frac{2\Delta t}{\Delta x} \sum_{k=1}^{\ell} \tilde{\lambda}_k^- (S(u_{k+1}) - S(u_k)),$$

$$S(u_L, u_R) \leq \frac{2}{\Delta x} \int_{u_L}^{u_R} \tilde{u} \left(\frac{x}{\Delta t}, u_L, u_R \right) dx = S(u_R) - \frac{2\Delta t}{\Delta x} \sum_{k=1}^{\ell} \tilde{\lambda}_k^+ (S(u_{k+1}) - S(u_k)),$$

so that with clear notations

$$S(u_j^{n+1}) \leq S(u_j^n) - \frac{\Delta t}{\Delta x} \left\{ \sum_{k=1}^{\ell} \tilde{\lambda}_{kj+\frac{1}{2}}^- (S(u_{k+1}^{j+\frac{1}{2}}) - S(u_k^{j+\frac{1}{2}})) + \sum_{k=1}^{\ell} \tilde{\lambda}_{kj-\frac{1}{2}}^+ (S(u_{k+1}^{j-\frac{1}{2}}) - S(u_k^{j-\frac{1}{2}})) \right\}$$

But by consistency with the entropy inequality we have

$$\sum_{k=1}^{\ell} \tilde{\lambda}_{kj+\frac{1}{2}}^- (S(u_{k+1}^{j+\frac{1}{2}}) - S(u_k^{j+\frac{1}{2}})) = \frac{1}{2} \sum_{k=1}^{\ell} (\tilde{\lambda}_{kj+\frac{1}{2}}^- - |\tilde{\lambda}_{kj+\frac{1}{2}}^-|) (S(u_{k+1}^{j+\frac{1}{2}}) - S(u_k^{j+\frac{1}{2}}))$$

$$\geq \frac{1}{2} (G(u_j) - G(u_{j-1})) - \frac{1}{2} \sum_{k=1}^{\ell} |\tilde{\lambda}_{kj+\frac{1}{2}}^-| (S(u_{k+1}^{j+\frac{1}{2}}) - S(u_k^{j+\frac{1}{2}}))$$

and

$$\sum_{k=1}^{\ell} \tilde{\lambda}_{kj-\frac{1}{2}}^+ (S(u_{k+1}^{j-\frac{1}{2}}) - S(u_k^{j-\frac{1}{2}})) \geq$$

$$\frac{1}{2} (G(u_j) - G(u_{j-1})) + \frac{1}{2} \sum_{k=1}^{\ell} |\tilde{\lambda}_{kj-\frac{1}{2}}^+| (S(u_{k+1}^{j-\frac{1}{2}}) - S(u_k^{j-\frac{1}{2}})).$$

This means that we eventually have

$$S(u_j^{n+1}) \leq S(u_j^n) + \frac{\Delta t}{\Delta x} (G(u_j, u_{j+1}^n) - G(u_{j-1}^n, u_j^n)), \quad j \in \mathbb{Z}$$

with

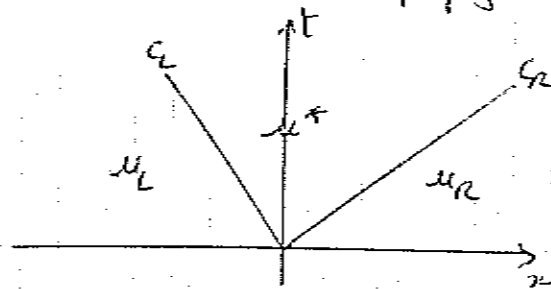
$$G_{j+\frac{1}{2}} = G(u_j^n, u_{j+1}^n) = \frac{1}{2} (G(u_j^n) + G(u_{j+1}^n)) - \frac{1}{2} \sum_{k=1}^{\ell} |\tilde{\lambda}_{kj+\frac{1}{2}}^-| (S(u_{k+1}^{j+\frac{1}{2}}) - S(u_k^{j+\frac{1}{2}})), \quad j \in \mathbb{Z}.$$

The proposed Godunov method is then consistent with (9) provided that the underlying approximate Riemann solver is consistent with the integral form of (9).

Examples of approximate Riemann solvers

The HLL approximate Riemann solver

One of the most simple and popular approximate Riemann solver consists in considering only two discontinuities propagating with speeds c_L and c_R :



so that only one intermediate state u^* has to be calculated. For that we use the consistency with the integral form of (3), namely (21) or equivalently (21)', that writes here

$$f(u_R) - f(u_L) = c_L (u^* - u_L) + c_R (u_R - u^*)$$

Then we have

$$u^* = \frac{c_R u_R - c_L u_L}{c_R - c_L} - \frac{1}{c_R - c_L} (f(u_R) - f(u_L))$$

Note that we have implicitly assumed $c_R - c_L > 0$.

We shall not give the details here, but it can be proved that this approximate Riemann solver is also consistent with the entropy inequality (9) provided that the so-called "subcharacteristic condition"

$$c_L \leq \lambda_k(u) \leq c_R, \quad k=1, \dots, p$$

is valid. It means that the eigenvalues of the system to be solved lie between the eigenvalues of the approximate Riemann solver. In other words, the information propagates faster in the approximate Riemann solver.

A local optimization of (20) is

$$c_L = \inf_{u=u_L, u_R} \inf_k \lambda_k(u) \quad \text{and} \quad c_R = \sup_{u=u_L, u_R} \sup_k \lambda_k(u)$$

In practice, the corresponding scheme is rather diffusive, but less than the next one.

The Lax-Friedrichs approximate Riemann solver

It consists in choosing $c_L = -c_R \equiv -c$, with

$$|\lambda_k(u)| \leq c, \quad k=1, \dots, p.$$

Again, stability properties can be proved in this context.

The Rusanov approximate Riemann solver

Here we simply take

$$c = \sup_{u=u_L, u_R} \sup_k |\lambda_k(u)|.$$

This works well in practice, except the excessive numerical diffusion on the waves associated with intermediate eigenvalues.

The Roe approximate Riemann solver

The idea of the Roe approximate solver is to exactly solve an linearized problem, associated with a so-called Roe linearization $A(u_L, u_R)$.

More precisely, the Roe approximate Riemann solver is defined by the exact Riemann solution of the linear system

$$\begin{cases} \partial_t u + A(u_L, u_R) \partial_x u = 0 \\ u(x, 0) = u_0(x) = \begin{cases} u_L & \text{if } x < 0, \\ u_R & \text{if } x > 0. \end{cases} \end{cases}$$

The matrix $A(u_L, u_R)$ is called a Roe-type linearization if the mapping $(u_L, u_R) \rightarrow A(u_L, u_R)$ from $\Omega \times \Omega$ into $\mathbb{R}^{p \times p}$ satisfies the following properties:

- (i) $f(u_R) - f(u_L) = A(u_L, u_R) (u_R - u_L)$ (consistency)
- (ii) $A(u_L, u_R)$ is \mathbb{R} -diagonalizable (hyperbolicity)
- (iii) $A(u, u) = \nabla_u f(u)$ (consistency)

In practice, (i) is the most difficult property to satisfy.

Let us now prove the existence of a Roe-type linearization when the system under consideration admits a strictly convex entropy.

Theorem (Harten-Lax, 1983)

Assume that (1) admits a strictly convex entropy S . Then there exists at least one Roe-type linearization.

Proof

The Roe-type linearization will be of the form

$$A(u_L, u_R) = R(u_L, u_R) P(u_L, u_R)$$

where R and P are symmetric matrices, P being in addition positive-definite. Then, $A = P^{-1/2} P^{1/2} R P^{1/2} P^{1/2}$ is similar to the matrix $P^{1/2} R P^{1/2}$ which is clearly symmetric, and therefore \mathbb{R} -diagonalizable.

Then, A is \mathbb{R} -diagonalizable as well, and property (ii) is immediately satisfied. Properties (i) and (iii) will follow from the definitions of R and P .

S being a strictly convex function, the change of variables $v(u) = \nabla S(u)$ is admissible and we can consider the function

$$g(v) = f(u(v))$$

Let us prove that $g'(v) = \nabla_u f(u(v)) \cdot u'(v)$ is symmetric.

Actually, since $u'(v) = v'(u)^{-1} = \nabla^2 S(u)^{-1}$, it amounts to prove that the matrix $\nabla_u f(u(v)) \cdot \nabla^2 S(u)^{-1} = \nabla^2 S(u)^{-1} (\nabla^2 S(u) \nabla_u f(u)) \nabla^2 S(u)^{-1}$ is symmetric, which clearly holds true since the similar matrix $\nabla^2 S(u) \nabla_u f(u)$ is symmetric by the strict convexity property of S .

For any given states u_L, u_R , let us set $v_L = \nabla S(u_L), v_R = \nabla S(u_R)$ and write

$$\begin{aligned} g(v_R) - g(v_L) &= \int_0^1 \frac{d}{d\theta} g(\theta v_R + (1-\theta)v_L) d\theta \\ &= \left\{ \int_0^1 g'(\theta v_R + (1-\theta)v_L) d\theta \right\} \cdot (v_R - v_L) \end{aligned}$$

We propose to set

$$R(u_L, u_R) = \int_{\sigma} g'(\theta u_R + (1-\theta)u_L) d\sigma$$

R is then symmetric, and satisfies

$$g(v_R) - g(v_L) = R(u_L, u_R) (v_R - v_L)$$

On the other hand

$$\begin{aligned} v_R - v_L &= \int_{\sigma} \frac{d}{d\theta} v(\theta u_R + (1-\theta)u_L) d\sigma \\ &= \left\{ \int_{\sigma} v'(\theta u_R + (1-\theta)u_L) d\sigma \right\} \cdot (u_R - u_L) \end{aligned}$$

We then naturally propose to set

$$P(u_L, u_R) = \int_{\sigma} v'(\theta u_R + (1-\theta)u_L) d\sigma$$

so that, since $v'(u) = \nabla S(u)$, P is clearly symmetric and positive definite, while we clearly have

$$\begin{aligned} A(u_L, u_R) (u_R - u_L) &= R(u_L, u_R) P(u_L, u_R) (u_R - u_L) \\ &= R(u_L, u_R) (v_R - v_L) \\ &= g(v_R) - g(v_L) \\ &= f(u_R) - f(u_L) \end{aligned}$$

Then property (i) is satisfied.

It remains to check (ii):

$$\begin{aligned} A(u, u) &= R(u, u) P(u, u) \\ &= g'(v) v'(u) \\ &= f'(u), \end{aligned}$$

which concludes the proof. \square

Remark:

It is easily checked that $P(u_L, u_R) = P(u_R, u_L)$ and $R(u_L, u_R) = R(u_R, u_L)$, so that we can reverse the roles of u_L and u_R ($A(u_L, u_R) = A(u_R, u_L)$) in the

Roe-type linearization proposed in the proof above.

Now that the Roe approximate Riemann solver can be defined, let us check that it gives rise to a simple approximate Riemann solver that is consistent with the integral form of (1).

First of all, we already saw that solving (33) gives rise to a simple approximate Riemann solver of the form (19), where $\tilde{\lambda}_k$, $k=1, \dots, \ell=p$, coincide with the eigenvalues of $A(u_L, u_R)$, and where the intermediate states u_k are given by

$$u_k = \sum_{l=1}^{k-1} (u_k, l_l) \pi_l + \sum_{l=k}^p (u_L, l_l) \pi_l$$

where π_l and l_l represent the right and left eigenvectors of the eigenvalues $\tilde{\lambda}_l$.

Then, remind that

$$\begin{aligned} u_{k+1} - u_k &= (u_k, l_k) \pi_k - (u_L, l_k) \pi_k \\ &= \{ (u_k, l_k) - (u_L, l_k) \} \pi_k \end{aligned}$$

so that

$$A(u_L, u_R) (u_{k+1} - u_k) = \tilde{\lambda}_k (u_{k+1} - u_k)$$

Then we successively obtain

$$\begin{aligned} f(u_{k+1}) - f(u_k) &= A(u_L, u_R) (u_{k+1} - u_k) && \text{(by (i))} \\ &= \sum_{l=1}^p A(u_L, u_R) (u_{k+1} - u_k) \\ &= \sum_{l=1}^p \tilde{\lambda}_l (u_{k+1} - u_k) \end{aligned}$$

which is nothing but the consistency property (21)' with the integral form of (1).

Then, the Roe approximate Riemann solver, with a Roe-type linearization $A(u_L, u_R)$ satisfying (i), (ii), (iii) above, gives rise to

an approximate Riemann solver which is consistent with the integral form of (1). The numerical flux function associated with the corresponding Godunov-type method is then, here again, given by (26), that is

$$f(u_L, u_R) = \frac{1}{2}(f(u_L) + f(u_R)) - \frac{1}{2} \sum_{k=1}^P |\tilde{\lambda}_k| (u_{k+1} - u_k).$$

A well-known drawback of this Roe's scheme is that it is not consistent with the entropy inequality (9), that is the consistency relation (22)', that is

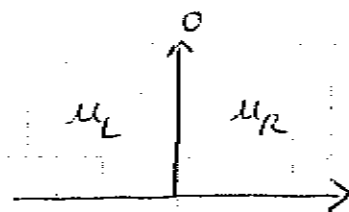
$$G(u_R) - G(u_L) \leq \sum_{k=1}^P \tilde{\lambda}_k (S(u_{k+1}) - S(u_k))$$

is not satisfied in general. To be convinced, it suffices to consider u_L and u_R such that $f(u_L) = f(u_R)$ and $G(u_L) \neq G(u_R)$.

Then, $A(u_L, u_R)(u_R - u_L) = 0$ which means that 0 is an eigenvalue of $A(u_L, u_R)$ and that $(u_R - u_L)$ is a corresponding eigenvector.

then

$$\begin{aligned} u_{k+1} - u_k &= \left\{ (u_R, l_k) - (u_L, l_k) \right\} \pi_k \\ &= \left\{ (u_R - u_L), l_k \right\} \pi_k \end{aligned}$$



equals 0 except across the discontinuity $\tilde{\lambda}_k = 0$, so that we immediately get in that case $\sum_{k=1}^P \tilde{\lambda}_k (S(u_{k+1}) - S(u_k)) = 0$.

Then, the consistency relation (22)' would write

$$G(u_R) - G(u_L) \leq 0.$$

This cannot be true since we can reverse the role of u_L and u_R in our argument just above.

This means that the Roe scheme may resolve nonphysical solutions.

Various entropy corrections have been proposed, for instance by Harten and Hyman in 1983, Roe in 1983, Roe and Pike in 1983, Huynh in 1995, ...

Another drawback of Roe's scheme is that it may generate in some particular test cases nonphysical intermediate states in the approximate Riemann solution, for instance with negative internal energies or densities in the frame of the Euler equations. A recent way to overcome this difficulty, as well as the one associated with the entropy inequality consistency, is given by the so-called Relaxation approach, that we now address. Note however that besides these drawbacks, the Roe scheme is again widely used in the industry.

The Relaxation approach