

# SPATIAL DISCRETIZATION

Shengtai Li and Hui Li

Los Alamos National Laboratory

## 1 OVERVIEW

## 2 FINITE DIFFERENCE METHOD

- Finite Difference for Linear Advection Equation

## 3 CONSERVATION LAWS

- Modern Finite Difference Methods for Conservation Laws

## 4 FINITE VOLUME METHOD

- System of PDEs
- Multi-dimensional Scheme

## 5 STAGGERED GRID AND CONSTRAINT TRANSPORT

## 6 SERIES EXPANSION METHOD

- Spectral Method
- Finite Element Method

# GENERAL CONCEPT TO NUMERICAL PDEs

Numerical discretization of a time-dependent partial differential equation (PDE) includes the following approaches:

- Fully-discretized in both time and space
- Method of lines (MOL) approach: discretize in space first, transform PDEs to ordinary differential equations (ODEs); and then using the time integration method of ODEs to advance in time. For MOL approach, we
  - do not worry about the time discretization and stability issue
  - can reuse available ODE solvers and softwares to solve a PDE problem.

Several spatial discretization methods:

- finite difference method (FDM): using point-value solution
- finite volume method (FVM): using cell-average value solution
- finite element method (FEM)
- spectral and pseudo-spectral method.

# CONVENTION FOR NOTATION FOR THIS LECTURE

- Time level: we use superscript  $n$  -  $u^n = u$  at time level  $n$ .
  - $\Delta t$  = time step =  $t^{n+1} - t^n$ . We often use  $k = \Delta t$  as our constant time-step.
  - $t = n\Delta t = nk$  where  $n = 0, 1, 2, \dots, N$ .  $T = N\Delta t$  = final time
  - $t^n$  = present;  $t^{n-1}$  = past;  $t^{n+1}$  = future
- Spatial Location: we use subscript  $i, j$  for  $x, y$ 
  - $\Delta x$  = spatial length scale.  $h = \Delta x$  denotes constant spatial length scale.
  - $x_i = x_L + i\Delta x$ ,  $i = 0, 1, 2, \dots, N$ , where  $x_L$  is the left boundary.  $x_R = x_N$  is the right boundary.
  - $\text{cell}_i = [x_{i-1}, x_i]$ . For FVM,  $\text{cell}_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$
  - a stencil means a range of consecutive points:  $[x_{i-s}, x_{i+r}]$  is a  $r + s + 1$ -point stencil.

# PRACTICAL EVALUATION CRITERIA OF NUMERICAL SCHEMES FOR PDEs

- **Solution accuracy:** quantify the numerical error, which is difference between the numerical solutions and the exact solutions. It depends on the initial data and boundary condition treatment as well as the following two criteria:
  - order of accuracy of the numerical schemes
  - preservation of the physical properties of the continuum PDEs: conservation laws (mass, momentum, and energy), density and pressure positivity, divergence-free condition of the magnetic fields, time-reversibility, etc.
- **Efficiency:** computer CPU time used to calculate the numerical solutions.
- **Algorithm robustness**

# INTRODUCTION TO FINITE DIFFERENCE SCHEMES

The basic idea of finite difference schemes is to replace derivatives by finite differences.

Methods for obtaining finite-difference expressions

- **Taylor series expansion:** the most common, but purely mathematical.
- **Polynomial fitting or interpolation:** the most general ways. Taylor series is a subset of this method.
- **Control volume approach:** also called finite volume method (FVM). We solve the equations in integral form rather than differential form. Popular in engineering where complex geometries and coordinate transformations are involved. For Cartesian grids, simplest FVM = FDM.

We will look at the first two approaches. The FVM will be covered in more details later.

# GENERAL METHOD FOR DERIVING FINITE DIFFERENCE SCHEMES

For a 3-point stencil  $[x_{i-1}, x_{i+1}]$ , we can write a generic expression as

$$\frac{\partial u}{\partial x}|_{x_i} = au_{i-1} + bu_i + cu_{i+1} + O(h^m)$$

where  $a$ ,  $b$ , and  $c$  are unknowns to be determined, and  $m$  is the order of approximation. (*General rule: if stencil spans  $n$  points, we can derive a  $n - 1$  order finite difference scheme*).

Using the Taylor series for  $u_{i\pm 1}$ , we can write

$$au_{i-1} + bu_i + cu_{i+1} = (a+b+c)u_i + (-a+c)hu_x + (a+c)\frac{h^2}{2}u_{xx} + (-a+c)\frac{h^3}{6}u_{xxx} + \dots$$

To have a second-order scheme, we set

$$a + b + c = 0, \quad -a + c = 1, \quad a + c = 0$$

from which we obtain  $b = 0$ ,  $c = -a = 1/(2h)$ , Therefore

$$\frac{\partial u}{\partial x}|_{x_i} = \frac{u_{i+1} - u_{i-1}}{2h} + O(h^2) \text{ is a second order scheme.}$$

# POLYNOMIAL FITTING

We assume that the solution of the PDE can be approximated by a polynomial, and the values at the mesh points are exact. Over a three-point stencil  $[x_{i-1}, x_{i+1}]$ , we assume the local fitting function is  $\bar{u}(x) = ax^2 + bx + c$ . Applying the polynomials to the three points gives

$$u_{i-1} = ax_{i-1}^2 + bx_{i-1} + c, \quad u_i = ax_i^2 + bx_i + c, \quad , u_{i+1} = ax_{i+1}^2 + bx_{i+1} + c$$

Solve for  $a, b, c$ , we obtain the *Lagrange Interpolation Polynomial*,

$$u(x) = u_{i-1} \frac{(x - x_i)(x - x_{i+1})}{(x_{i-1} - x_i)(x_{i-1} - x_{i+1})} + u_i \frac{(x - x_{i-1})(x - x_{i+1})}{(x_i - x_{i-1})(x_i - x_{i+1})} + u_{i+1} \frac{(x - x_{i-1})(x - x_i)}{(x_{i+1} - x_{i-1})(x_{i+1} - x_i)}$$

Note that grid spacing does not have to be uniform.  
Differentiating  $\bar{u}(x)$  with respect to (w.r.t)  $x$ , we obtain

$$u_x = (u_{i+1} - u_{i-1})/(2h), \quad u_{xx} = (u_{i+1} - 2u_i + u_{i-1})/h^2$$

where we assume a uniform grid is used.



# DIFFERENCE FORMULA FOR ONE-DIMENSION UNIFORM GRID

Using either of the above approaches, we can obtain the finite-difference formula for several derivatives in both one-dimensional and two-dimensional problem

The first-derivatives:

$$u'(x_0) = \frac{u_1 - u_0}{h} + O(h),$$

$$u'(x_0) = \frac{u_1 - u_{-1}}{2h} + O(h^2),$$

$$u'(x_0) = \frac{-u_2 + 4u_1 - 3u_0}{2h} + O(h^2)$$

$$u'(x_0) = \frac{-u_2 + 8u_1 - 8u_{-1} + u_{-2}}{12h} + O(h^4)$$

- Second derivative:

$$u''(x_0) = \frac{u_2 - 2u_1 + u_0}{h^2} + O(h),$$

$$u''(x_0) = \frac{u_1 - 2u_0 + u_{-1}}{h^2} + O(h^2),$$

$$u''(x_0) = \frac{-u_3 + 4u_2 - 5u_1 + 2u_0}{h^2} + O(h^2)$$

$$u''(x_0) = \frac{-u_2 + 16u_1 - 30u_0 + 16u_{-1} - u_{-2}}{12h^2} + O(h^4)$$

- Third derivatives:

$$u'''(x_0) = \frac{u_3 - 3u_2 + 3u_1 - u_0}{h^3} + O(h),$$

$$u'''(x_0) = \frac{u_2 - 2u_1 + 2u_{-1} - u_{-2}}{2h^3} + O(h^2),$$

- Fourth derivatives:

$$u^{(4)}(x_0) = \frac{u_4 - 4u_3 + 6u_2 - 4u_1 + u_0}{h^4} + O(h),$$

$$u^{(4)}(x_0) = \frac{u_2 - 4u_1 + 6u_0 - 4u_{-1} + u_{-2}}{h^4} + O(h^2),$$

# DIFFERENCE FORMULA FOR TWO-DIMENSION UNIFORM GRID

- Mixed derivatives

$$u_{xy}(x_0, y_0) = \frac{1}{4h^2}(u_{1,1} - u_{1,-1} - u_{-1,1} + u_{-1,-1}) + O(h^2)$$

- Laplacian

$$\nabla^2 u(x_0, y_0) = \frac{1}{h^2}(u_{1,0} + u_{0,1} + u_{-1,0} + u_{0,-1} - 4u_{0,0}) + O(h^2)$$

$$\begin{aligned} \nabla^2 u(x_0, y_0) = & \frac{1}{12h^2}(-60u_{0,0} + 16(u_{1,0} + u_{0,1} + u_{-1,0} + u_{0,-1}) \\ & -(u_{2,0} + u_{0,2} + u_{-2,0} + u_{0,-2})) + O(h^4) \end{aligned}$$

# LINEAR ADVECTION EQUATION

We now consider the finite-difference scheme for a simplest case:  
linear advection equation

$$u_t + au_x = 0$$

$$u(x, 0) = u_0(x)$$

where  $a$  is constant. We also assume  $a > 0$ . Its solution is simply a translation of the  $u_0$ ,

$$u(x, t) = u_0(x - at)$$

# BASIC FINITE DIFFERENCE SCHEMES FOR LINEAR ADVECTION EQUATION

Combine the spatial and temporal discretization together, we list some finite-difference scheme below

- 
- Upwind (first-order)  $\frac{u_i^{n+1} - u_i^n}{k} + a \frac{u_i^n - u_{i-1}^n}{h} = 0$
- Lax-Friedrichs scheme,  $\frac{u_i^{n+1} - \frac{1}{2}(u_{i+1}^n + u_{i-1}^n)}{k} + a \frac{u_{i+1}^n - u_{i-1}^n}{2h} = 0$
- Lax-Wendroff scheme
$$u_i^{n+1} = u_i^n - \frac{ak}{2h}(u_{i+1}^n - u_{i-1}^n) + \frac{a^2 k^2}{2h^2}(u_{i+1}^n - 2u_i^n + u_{i-1}^n)$$
- Beam-Warming ( $2^{\text{nd}}$  order upwind):
$$u_i^{n+1} = u_i^n - \frac{ak}{2h}(3u_i^n - 4u_{i-1}^n + u_{i-2}^n) + \frac{a^2 k^2}{2h^2}(u_{i+1}^n - 2u_i^n + u_{i-1}^n)$$
- Leapfrog in time central space,  $\frac{u_i^{n+1} - u_i^{n-1}}{2k} + a \frac{u_{i+1}^n - u_{i-1}^n}{2h} = 0$

An important trick, called “Lax-Wendroff” approach, to construct high order scheme is to replace the high-order time derivatives with high-order spatial derivatives with the help of PDE.

In general, an explicit finite difference scheme can be expressed as

$$u_i^{n+1} = F(u_{i-l}^n, u_{i-l+1}^n, \dots, u_{i+r}^n) = \sum_{k=-l}^r a_k u_{i+k}^n$$

### Homework:

- 1 write a computing program using the above listed schemes to the linear advection equation. Use periodic boundary conditions and Gaussian initial conditions.

# QUANTITATIVE PROPERTIES OF FINITE DIFFERENCE SCHEMES

- **Consistency:** An finite-difference discretization of a PDE is consistent if the finite-difference equations converge to the PDE, i.e., the truncation error vanishes as grid spacing and time step tend to zero.
- **Stability:** the errors from any source will not grow unbounded with time
- **Convergence:** the solution of the finite-difference equations converge to the true solution of the PDE as grid spacing and time step tend to zero.



# STABILITY: COURANT-FRIDRICHS-LEVY (CFL) CONDITION

For a finite-difference scheme,  $u_i^{n+1} = F(u_{i-l}^n, u_{i-l+1}^n, \dots, u_{i+r}^n)$ , the numerical domain dependence of  $(x_i, t_n)$  is  $[x_{i-l}, x_{i+r}]$ . To ensure the finite-difference scheme stable, a natural condition is

physical domain dependence  $\in$  numerical domain dependence

This gives a Courant-Fridrichs-Levy (CFL) condition on the ratio of  $h$  and  $k$ . For the linear advection equation with  $a > 0$ , the CFL condition is

$$0 \leq \frac{ak}{lh} \leq 1$$

# CONSISTENCY AND TRUNCATION ERROR

Let us express the difference scheme as  $u^{n+1} = Fu^n$ . The truncation error is defined as  $e^n = \frac{u^{n+1} - Fu^n}{k}$ . A finite-difference scheme is called *consistent* if  $e^n \rightarrow 0$  as  $k \rightarrow 0$  and  $h \rightarrow 0$ . For finite difference scheme  $u_i^{n+1} = \sum_{m=-l}^r a_m u_{i+m}^n$ , the necessary condition for consistency is

$$\sum_{m=-l}^r a_m = 1.$$

because the constant is a solution to the PDE.

If  $e^n = O(k^p, h^q)$ , then the scheme is called of order  $(p, q)$ . It is easy to check the upwind method has accuracy of order (1,1), and the Lax-Wendroff method has accuracy of order (2,2).

**Exercise:** Find the truncation error of the schemes listed above.

# LAX-RICHTMYER EQUIVALENCE THEOREM

## THEOREM

**The Lax-Richtmyer equivalence theorem.** *A consistent finite difference scheme for a partial differential equation for which the initial value problem is well-posed is convergent if and only if it is stable*

This implies

consistency + Stability  $\Rightarrow$  Convergence

consistency + Convergence  $\Rightarrow$  Stability

- **Consistency** implies that the solution of the PDE is an approximate solution of the finite difference scheme.
- **Convergence** means that a solution of the finite difference scheme approximates a solution of the PDE.

# EXAMPLES FOR CONVERGENCE AND CONSISTENCY

**The Lax-Friedrichs Scheme** For the Lax-Friedrichs scheme,

$$\phi_i^{n+1} = \frac{1}{2}(\phi_{i+1}^n + \phi_{i-1}^n) - \frac{ak}{2h}(\phi_{i+1}^n - \phi_{i-1}^n)$$

The CFL condition is  $\frac{ak}{h} \leq 1$ . From the Taylor series (derivatives are evaluated at  $(x_i, t^n)$ ) we have

$$\frac{1}{2}(\phi_{i+1}^n + \phi_{i-1}^n) = \phi_i^n + \frac{h^2}{2}\phi_{xx} + O(h^4), \quad \frac{\phi_{i+1}^n - \phi_{i-1}^n}{2h} = \phi_x + \frac{h^2}{6}\phi_{xxx} + O(h^4).$$

Substituting these expressions in the scheme, we obtain

$$e^n = \frac{1}{2}k\phi_{tt} - \frac{1}{2}k^{-1}h^2\phi_{xx} + \frac{1}{6}ah^2\phi_{xxx} + O(k^4 + k^{-1}h^4 + k^2)$$

So  $e^n \rightarrow 0$  as  $h, k \rightarrow 0$ , i.e., it is consistent, as long as  $k^{-1}h^2 \rightarrow 0$ . Note that reducing time step ( $k \rightarrow 0$ ) with a fixed spatial length  $h$  leads to inconsistency for this scheme.

# A CONSISTENT SCHEME MAY NOT CONVERGE

Consider the forward time forward space scheme

$\frac{u_i^{n+1} - u_i^n}{k} + a \frac{u_{i+1}^n - u_i^n}{h} = 0$ . It is easy to verify it is NOT stable using the domain dependence. This scheme can be rewritten as

$$u_i^{n+1} = u_i^n - \frac{ak}{h}(u_{i+1}^n - u_i^n) = (1 + \lambda)u_i^n - \lambda u_{i+1}^n \quad (1)$$

where  $\lambda = ak/h$ . It is easy to show it is consistent (**Exercise**). Assume  $a = 1$ . The solution of the PDE is a shift of  $u_0$  to the right by  $t$ . If we take the initial condition as

$$u_0(x) = \begin{cases} 1 & \text{if } -1 \leq x \leq 0 \\ 0 & \text{elsewhere} \end{cases} \Rightarrow u_i^0 = \begin{cases} 1 & \text{if } -1 \leq ih \leq 0 \\ 0 & \text{elsewhere} \end{cases}$$

As equation (1) shows, the solution  $(t_n, x_i)$  depends only on  $x_m$  for  $m \geq i$  at previous times. Therefore,  $u_i^n = 0, \forall i > 0, n \geq 0$ , which mean  $u_i^n$  cannot converge to  $u(t, x)$ , since for positive  $t$  and  $x$ , the function  $u(t, x)$  is not identically zero, yet  $u_i^n$  is zero.

# STABILITY ANALYSIS: VON NEUMANN METHOD

## THEOREM

A finite-difference scheme with constant coefficients

$$u_i^{n+1} = \sum_{m=-l}^r a_m u_{i+m}^n \text{ is stable if and only if } \hat{G}(\xi) := \sum_{m=-l}^r a_m e^{-ik\xi}$$

satisfies  $\max_{-\pi \leq \xi \leq \pi} |\hat{G}(\xi)| \leq 1$ .

As a simple example, we show that forward time central space scheme,  $\frac{u_i^{n+1} - u_i^n}{k} + a \frac{u_{i+1}^n - u_{i-1}^n}{2h} = 0$ , is unstable. The corresponding  $\hat{G}(\xi) = 1 + i\lambda \sin \xi$ , which cannot be bounded by 1 in magnitude.

**Exercise:** Compute the  $\hat{G}(\xi)$  for the schemes: Lax-Friedrichs, Lax-Wendroff, Leap-Frog, Beam-Warming.

# METHOD OF LINES AND SPATIAL DISCRETIZATION

For wave equation  $u_t + au_x = 0$ , the central-differencing in space results in

$$\frac{du_i}{dt} + a \frac{u_{i+1} - u_{i-1}}{2h} = 0$$

For heat equation  $u_t = cu_{xx}$ , the central-differencing in space results in

$$\frac{du_i}{dt} = c \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2}$$

# DISSIPATION, DISPERSION, AND THE MODIFIED EQUATION

The truncation error in the finite-difference formula can be described as dissipation (amplitude) error or dispersion (phase-speed) error.

- One-sided, first-order

$$\frac{u_i - u_{i-1}}{h} = \frac{\partial u}{\partial x} - \frac{h}{2} \frac{\partial^2 u}{\partial x^2} + \frac{h^2}{6} \frac{\partial^3 u}{\partial x^3} + O(h^3)$$

- centered, second-order

$$\frac{u_{i+1} - u_{i-1}}{h} = \frac{\partial u}{\partial x} + \frac{h^2}{6} \frac{\partial^3 u}{\partial x^3} + O(h^4)$$

- centered, fourth-order

$$\frac{-u_2 + 8u_1 - 8u_{-1} + u_{-2}}{12h} = \frac{\partial u}{\partial x} - \frac{h^4}{30} \frac{\partial^5 u}{\partial x^5} + O(h^6)$$



If the spatial accuracy is  $O(h^m)$ , the MOL approach will approximate *modified equation*

$$u_t + au_x = bh^m \frac{\partial^{m+1} u}{\partial x^{m+1}} + ch^{m+1} \frac{\partial^{m+2} u}{\partial x^{m+2}} \quad (2)$$

to  $O(h^{m+2})$ , where  $b$  and  $c$  are rational numbers determined by the particular finite-difference scheme.

- The even derivative on the right-side of (2) produces amplitude error, or *numerical dissipation*.
- The odd-order derivative on the right side of (2) produces a wave-number-dependent phase error known as *numerical dispersion*.

# COMPACT DIFFERENCING

- High order scheme depends on a larger stencil. The sixth-order scheme

$$u_x = \frac{3}{2}\delta_{2h}u - \frac{3}{5}\delta_{4h}u + \frac{1}{10}\delta_{6h}u + O(h^6) \quad (3)$$

where  $\delta_{mh}u = (u_{i+m/2} - u_{i-m/2})/(mh)$ , provides only marginal improvement over the fourth-order scheme.

- Rewrite the second-order central difference into

$$\delta_{2h}u = \left(1 + \frac{h^2}{6}\delta_h^2\right) u_x + O(h^4),$$

where  $\delta_h^2 = \delta_h(\delta_h u) = \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2}$ . Expanding the finite-difference operator yields the following  $O(h^4)$  accurate scheme

$$\frac{u_{i+1} - u_{i-1}}{2h} = \frac{1}{6} ((u_x)_{i+1} + 4(u_x)_i + (u_x)_{i-1}) \quad (4)$$

which is calculated on three-point stencil. At intermediate numerical resolution, the compact fourth-order scheme (4) is typically more accurate than the sixth-order scheme (3).

# CONSERVATION LAWS IN IDEAL MHD EQUATIONS

Ideal MHD equations:

$$\rho_t + \nabla \cdot (\rho \mathbf{u}) = 0, \text{ (Mass conservation)}$$

$$(\rho \mathbf{u})_t + \nabla \cdot [\rho \mathbf{u} \mathbf{u}^T + (p) \mathbf{I} - \mathbf{B} \mathbf{B}^T] = 0, \text{ (Momentum conservation)}$$

$$e_t + \nabla \cdot [(e + p) \mathbf{u} - \mathbf{B}(\mathbf{u} \cdot \mathbf{B})] = 0, \text{ (Energy conservation)}$$

$$\mathbf{B}_t - \nabla \cdot (\mathbf{u} \mathbf{B}^T - \mathbf{B} \mathbf{u}^T) = 0, \text{ (Magnetic Flux Conservation)}$$

We will first describe the numerical schemes for scalar conservation law:

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

# WEAK SOLUTIONS

If a function contains a discontinuity, it cannot be the solution to a PDE in the conventional sense, because derivatives are not defined at the discontinuity. Take scalar conservation law  $u_t + f(u)_x = 0$  as an example. Use the integral form (*weak form*), and apply some trial function  $\phi$ ,

$$\begin{aligned} & \int_0^\infty \int_{-\infty}^\infty [u_t + f(u)_x] \phi(x, t) dx dt = 0 \\ \Rightarrow & \int_0^\infty \int_{-\infty}^\infty [u \phi_t + f(u) \phi_x] dx dt = - \int_0^\infty u_0(x) \phi_0(x) dx. \quad (5) \end{aligned}$$

where  $\phi$  is any  $C_0^1$  function, which is continuously differential and have a compact support, meaning it is zero outside of some bounded region. Note that (5) is valid even if  $u$  is discontinuous.

## DEFINITION

The function  $u(x, t)$  is a *weak solution* of the conservation law with given initial data  $u_0(x)$  if (5) holds for all  $C_0^1$  functions  $\phi$ .

# NUMERICAL ISSUES OF FINITE-DIFFERENCE SCHEMES FOR CONSERVATION LAWS

- Spurious oscillation appear around discontinuities in every high order schemes (dispersion error).
- The convergent solution may not be a weak solution. For example, the shock position is totally wrong for some non-conservative scheme.
- The convergent weak solution may not be an entropy solution (non-physical). It requires entropy-fix.

# CONSERVATIVE SCHEMES

## DEFINITION

A finite-difference approximation to the scalar conservation law is in conservative form if it can be written as  $\frac{u_i^{n+1} - u_i^n}{k} + \frac{\hat{f}_{i+\frac{1}{2}} - \hat{f}_{i-\frac{1}{2}}}{h} = 0$  where  $\hat{f}_{i+\frac{1}{2}}$  is numerical flux approximation of the form  $\hat{f}_{i+\frac{1}{2}} = \hat{f}(u_{i-r}, \dots, u_{i+s})$ ,  $r, s$  are integers, and satisfy the consistent condition  $\hat{f}(u_0, u_0, \dots, u_0) = f(u_0)$ .

**Exercise:** verify that Lax-Friedrichs and Lax Wendroff schemes are conservative schemes and derive the form of  $\hat{f}_{i+\frac{1}{2}}$ .

## THEOREM

*The conservative scheme will converge to a weak solution if it converges (Wendroff and Lax 1960)*

The nonconservative scheme may not converge to a weak solution if the solution is discontinuous.

# EXAMPLES

- Conservative scheme: two-step Lax-Wendroff

$$\hat{f}_{i+\frac{1}{2}}^{n+\frac{1}{2}} = f(u_{i+\frac{1}{2}}^{n+\frac{1}{2}}), \quad u_{i+\frac{1}{2}}^{n+\frac{1}{2}} = \frac{1}{2} \left( u_i^n + u_{i+1}^n + \frac{k}{h} (f(u_i^n) - f(u_{i+1}^n)) \right)$$

- Nonconservative scheme for  $f(u) = \frac{1}{2}u^2$ :

$$\left( \frac{1}{2}u^2 \right)_x = uu_x := u_i \frac{u_{i+1} - u_{i-1}}{2h}$$

# OVERVIEW OF THE MODERN SCHEMES FOR CONSERVATION LAWS

Roughly speaking, modern schemes can be classified into two categories:

- **flux-splitting methods**, The basic idea is to add a switch such that the scheme becomes first order near discontinuity and remains high order in the smooth region.
  - artificial viscosity methods, (ZEUS code)
  - flux-correction transport (FCT)
  - total variation diminishing (TVD) or bounded (TVB)
  - central scheme
  - relaxation schemes
- **high-order Godunov methods**: MUSCL, piecewise parabolic method (PPM), essentially non-oscillatory (ENO) schemes, and wave-propagation method of Leveque.

For the first category, we will describe the TVD and FCT schemes. For the second, we will describe it in the section of finite-volume method.



# FLUX-CORRECTED TRANSPORT (FCT)

Consider the conservative scheme  $\frac{u_i^{n+1} - u_i^n}{k} + \frac{f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}}}{h} = 0$ . The fluxes  $f_{i+\frac{1}{2}}$  is computed as follows

- 1 Compute low-order fluxes  $f_{i+\frac{1}{2}}^l$  using a low-order upwind scheme
- 2 Compute high-order fluxes  $f_{i+\frac{1}{2}}^h$  using a high-order scheme
- 3 Compute the antidiffusive fluxes  $A_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}^h - f_{i+\frac{1}{2}}^l$
- 4 Predict the solution at  $(n+1)k$  (also known as “transported and diffused” solution) using the low-order fluxes

$$u_i^{td} = u_i^n - \frac{k}{h}(f_{i+\frac{1}{2}}^l - f_{i-\frac{1}{2}}^l)$$

- 5 Correct the  $A_{i+\frac{1}{2}}$  so that the final “antidiffusion” step does not generate new maxima or minima:  $A_{i+\frac{1}{2}}^c = C_{i+\frac{1}{2}} A_{i+\frac{1}{2}}$ ,  $0 \leq C_{i+\frac{1}{2}} \leq 1$ .
- 6 Perform the “antidiffusion” step

$$u_i^{n+1} = u_i^{td} - \frac{k}{h}(A_{i+\frac{1}{2}}^c - A_{i-\frac{1}{2}}^c)$$

# FLUX-CORRECTED TRANSPORT: FLUX CORRECTION

- Original proposal of Boris and Book (1973):

$$\frac{k}{h} A_{i+\frac{1}{2}}^c = \max \left( 0, \min \left[ A_{i+\frac{1}{2}} \frac{k}{h}, u_{i+2}^{td} - u_{i+1}^{td}, u_i^{td} - u_{i-1}^{td} \right] \right)$$

- Zalesak (1979) corrector: has several improvement over the original FCT correction. See Zalesak, JCP(1979), for more information.

# TOTAL VARIATION DIMINISHING (TVD)

Consider the linear advection equation where  $f(u) = au$  with  $a > 0$ . The numerical flux is

$$f_{i+\frac{1}{2}} = au_i + C_{i+\frac{1}{2}} \left( \frac{1}{2}a(1 - ak/h)(u_{i+1} - u_i) \right),$$

where  $C_{i+\frac{1}{2}} := C(\theta_{i+\frac{1}{2}})$  is a flux-limiter that depends on  $\theta_{i+\frac{1}{2}} := \frac{u_i - u_{i-1}}{u_{i+1} - u_i}$ .

It is proved that the scheme is TVD if  $0 \leq \frac{C(\theta)}{\theta} \leq 2$  and  $0 \leq C(\theta) \leq 2$ .

- $C(\theta) = 1$  gives the Lax-Wendroff scheme.  $C(\theta) = \theta$  gives Beam-Warming scheme.
- van Leer limiter:  $C(\theta) = \frac{\theta + |\theta|}{1 + |\theta|}$
- Monotone central limiter:  
 $C(\theta) = \max[0, \min(2\theta, (1 + \theta)/2, 2)]$ .
- superbee limiter:  $C(\theta) = \max(0, \min(1, 2\theta), \min(\theta, 2))$

# FINITE VOLUME SCHEMES

In a *finite volume method* the average values of a function over local mesh cells are taken as the unknowns. Discrete approximations of the divergence, gradient, and curl operators are defined using general forms of Stokes' Theorem. Take the 1D conservation law  $u_t + f(u)_x = 0$  as an example. Integrating over the interval  $[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$  yields

$$\frac{d\bar{u}(x_i, t)}{dt} = -\frac{1}{\Delta x_i} \left( f(u(x_{i+\frac{1}{2}}, t)) - f(u(x_{i-\frac{1}{2}}, t)) \right)$$

where  $\bar{u}(x_i, t) = \frac{1}{\Delta x_i} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} u(\xi, t) d\xi$  is the cell average. The finite volume discretization is to approximate flux  $f(u(x_{i+\frac{1}{2}}, t))$  using a numerical flux  $\hat{f}_{i+\frac{1}{2}} = h(u_{i+\frac{1}{2}}^-, u_{i+\frac{1}{2}}^+)$  where  $u_{i+\frac{1}{2}}^\pm$  are the approximate values of  $u$  at interface from both side. The two-argument function  $h$  is a monotone flux that satisfies consistent condition  $h(u, u) = f(u)$ .

# PROCEDURE TO CONSTRUCT FLUX FOR FINITE VOLUME SCHEME

- Reconstruction: from the cell-centered average value  $\bar{u}_i$ , reconstruct an interpolate formula and then evaluate the point value  $u_{i+\frac{1}{2}}^\pm$  at the cell-interface.
- Flux evaluation:
  - Godunov flux: 
$$h(a, b) = \begin{cases} \min_{a \leq u \leq b} f(u), & \text{if } a \leq b \\ \max_{b \leq u \leq a} f(u), & \text{if } a > b \end{cases}$$
  - Lax-Friedrichs flux:  $h(a, b) = \frac{1}{2} (f(a) + f(b) - \alpha(b - a))$ , where  $\alpha \geq \max_u |f'(u)|$  is a constant
  - Solve the Riemann problem to evaluate the flux.

# RECONSTRUCTION

Given cell averages  $\bar{u}_i$ , find a polynomial  $p_i(x)$ , of degree at most  $q - 1$ , for each cell  $I_i = [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$ , such that it is a  $q$ -th order accurate approximation to  $u(x)$  inside  $I_i$ .

- Piecewise constant:  $q = 1$ , and  $p_i(x) = \bar{u}_i$
- Piecewise linear:  $q = 2$ , and  $p_i(x) = \bar{u}_i + \sigma_i(x - x_i)$ , where  $\sigma$  is a slope of solution  $u$  over  $I_i$ .
- Piecewise parabolic:  $q = 3$  (see Woodward and Colella JCP(1984))
- Essentially Non-Oscillatory scheme (ENO) (see Shu and Osher JCP(1989))
- Weighted ENO scheme (WENO) (see Jiang and Shu JCP (1997))

# PIECEWISE LIMITED LINEAR RECONSTRUCTION

To make the scheme *total variation diminishing* (TVD), the slope in the piecewise linear reconstruction must be limited.

- Generalized Minmod limiter:

$$\Delta x \sigma_i = \text{minmod} \left( \theta(u_{i+1} - u_i), \frac{1}{2}(u_{i+1} - u_{i-1}), \theta(u_i - u_{i-1}) \right)$$

where  $\text{minmod}(a, b, c) = \text{minmod}(\text{minmod}(a, b), c)$  and

$$\text{minmod}(x, y) = \begin{cases} 0, & \text{if } xy \leq 0 \\ \text{sgn}(x) \min(|x|, |y|), & \text{Otherwise} \end{cases}$$

where  $\theta \in [1, 2]$ . The larger the  $\theta$ , the less dissipative the scheme.

- minmod limiter :  $\theta = 1$ ,
- monotone central limiter :  $\theta = 2$

# HIGH ORDER FINITE-VOLUME RECONSTRUCTION

**Problem:** Given  $k = r + s + 1$  cell average,  $\bar{u}_{i-r}, \dots, \bar{u}_{i+s}$ , find a polynomial  $p_i(x)$ , of degree at most  $k - 1$ , such that it is a  $k$ -th order accurate approximation to the function  $u(x)$ , and

$$\frac{1}{\Delta x_j} \int_{x_{j-1/2}}^{x_{j+1/2}} p(\xi) d\xi = \bar{u}_j, \quad j = i - r, \dots, i + s.$$

**Algorithm:** Consider primitive function  $U(x) = \int_{-\infty}^x u(\xi) d\xi$ . We have

$$U(x_{j+1/2}) = \sum_{j=-\infty}^i \int_{x_{j-1/2}}^{x_{j+1/2}} u(\xi) d\xi = \sum_{j=-\infty}^i \bar{u}_j \Delta x_j$$

where  $-\infty$  can be replaced by any fixed number. We can construct a unique polynomial  $P(x)$  of degree  $k$ , which interpolate  $U(x_{j+1/2})$  at  $k + 1$  points:  $x_{i-r-1/2}, \dots, x_{i+s+1/2}$ . Then we define

$$p(x) := P'(x)$$

Knowing  $p(x)$ , we can obtain the interface values as

$$u_{i+1/2}^- = p_i(x_{i+1/2}), \quad u_{i-1/2}^+ = p_i(x_{i-1/2}), \quad i = 1, \dots, N.$$

**Exercise:** verify  $\frac{1}{\Delta x_j} \int_{x_{j-1/2}}^{x_{j+1/2}} p(\xi) d\xi = \bar{u}_j$ .



# HIGH ORDER CONSERVATIVE FINITE-DIFFERENCE RECONSTRUCTION

**Problem:** Given the point values of a function  $f(x)$ :  
 $f_i = f(x_i), i = 1, 2, \dots, N$ , find a numerical flux function  
 $\hat{f}_{i+\frac{1}{2}} = \hat{f}(f_{i-r}, \dots, f_{i+s}), i = 0, 1, \dots, N$ , such that the flux difference  
approximates the derivative  $f_x(x)$  to  $k$ -th order accuracy:

$$\frac{\hat{f}_{i+\frac{1}{2}} - \hat{f}_{i-\frac{1}{2}}}{\Delta x_i} = f_x(x_i) + O(\Delta x^k), \quad i = 0, 1, \dots, N.$$

**Algorithm:** We assume the grid is uniform with  $h = \Delta x$  and try to find a function  $F(x)$  such that  $f(x_i)$  is cell-average of  $F(x)$ . Then clearly

$$f(x) = \frac{1}{h} \int_{x-h/2}^{x+h/2} F(\xi) d\xi \Rightarrow f'(x) = \frac{1}{h} (F(x+h/2) - F(x-h/2))$$

Note that  $f(x_i)$  is the cell-average of the unknown function  $F(x)$ , we can use the procedure of finite-volume reconstruction to find  $F(x)$ .

# ENO RECONSTRUCTION IN ONE DIMENSION

The Newton form of the  $k$ -th degree interpolation polynomial  $P(x)$ , which interpolates  $U(x)$  at  $k + 1$  points can be expressed using the divided differences by

$$P(x) = \sum_{j=0}^k U[x_{i-r-\frac{1}{2}}, \dots, x_{i+s+\frac{1}{2}}] \prod_{m=0}^{j-1} (x - x_{i-r+m-\frac{1}{2}})$$

**Basic Idea of ENO:** find the “smoothest” stencil among all the possible stencils with  $k + 1$  consecutive points. **Algorithm:** Start with two point stencil and polynomial

$P^1(x) = U[x_{i-1/2}] + U[x_{i-1/2}, x_{i+1/2}](x - x_{i-1/2})$ , add one point to the stencil in each step.

- If  $|U[x_{i-1/2}, x_{i+1/2}, x_{i+3/2}]| < |U[x_{i-3/2}, x_{i-1/2}, x_{i+1/2}]|$ , choose the 3-point stencil as  $S_3(i) = \{x_{i-1/2}, x_{i+1/2}, x_{i+3/2}\}$ ; otherwise, take  $S_3(i) = \{x_{i-3/2}, x_{i-1/2}, x_{i+1/2}\}$
- The procedure is continued, according to the smaller of the absolute values of the two relevant divided differences, until the desired number of points in the stencil is reached.

# WENO (WEIGHTED ENO) RECONSTRUCTION

**Basic Idea of WENO:** instead of using only one of the candidate stencils to form the reconstruction, one uses a convex combination of all of them.

Suppose the  $k$  candidate stencils

$S_r(i) = \{x_{i-r}, \dots, x_{i-r+k-1}\}$ ,  $r = 0, k-1$ , produce  $k$  different reconstructions to the value  $u_{i+\frac{1}{2}}$ , we denote

$$u_{i+\frac{1}{2}}^{(r)} = \sum_{j=0}^{k-1} c_{rj} \bar{u}_{i-r+j}, \quad r = 0, \dots, k-1$$

WENO reconstruction would take a convex combination of all  $u_{i+1/2}^{(r)}$  as a new approximation to the  $u(x_{i+1/2})$ .

$$u_{i+1/2} = \sum_{r=0}^{k-1} w_r u_{i+1/2}^{(r)},$$

where the weight  $w_r$  satisfies  $w_r > 0$ ,  $\sum_{r=0}^{k-1} w_r = 1$ .

# SYSTEM OF PDEs

- The Jacobian  $\partial_u f(u)$  for system of conservation laws has  $m$  real eigenvalues. For example, the one dimensional ideal MHD system has **seven** eigenvalues and a complete set of eigenvectors.
- Exact Riemann solver is costly. Approximate Riemann solver is usually adopted.
- Reconstruction procedure:
  - Component-wise: apply the scalar algorithm to each component of the solution.
  - Characteristic-wise: project the solution to the characteristic space, do the reconstruction using the scalar algorithm, and project back to the solution space.

Characteristic-wise is more expensive but more accurate and oscillation-free.

# EXAMPLE: EIGENVALUES FOR IDEAL MHD

## SEVEN WAVES

$$\lambda_{1,7} = v_x \pm c_f, \quad \lambda_{2,6} = v_x \pm c_a, \quad \lambda_{3,5} = v_x \pm c_s, \quad \lambda_4 = v_x,$$

where  $c_a = \sqrt{\frac{B_x^2}{\rho}}$  is the speed of Alfvén waves, and

$$c_{f,s} = \left[ \frac{1}{2} \left( a^2 \pm \frac{B^2}{\rho} + \sqrt{\left( a^2 + \frac{B^2}{\rho} \right)^2 - 4a^2 c_a^2} \right) \right]^{\frac{1}{2}},$$

are the speeds of fast and slow magneto-sonic waves, and  $a = \sqrt{\frac{\gamma p}{\rho}}$  is the speed of the acoustic wave.

## EIGHT WAVES

$$\lambda_{1,8} = v_x \pm c_f, \quad \lambda_{2,7} = v_x \pm c_a, \quad \lambda_{3,6} = v_x \pm c_s, \quad \lambda_{4,5} = v_x,$$

# EXAMPLE: MUSCL-HANCOCK SCHEME FOR IDEAL MHD

- Reconstruct the primitive variables using piece-wise linear limited reconstruction: Minmod, van Leer, monotone-central limiters.

$$u_i(x) = u_i + (x - x_i) \overline{\Delta u_i} / \Delta x,$$

- Predict the time-center values using the Hancock predictor

$$u_i^{n+\frac{1}{2}} = u_i^n + \frac{1}{2} \frac{\Delta t}{\Delta x} \left( F(u_i^n + \frac{1}{2} \overline{\Delta_x u_i^n}) - F(u_i^n - \frac{1}{2} \overline{\Delta_x u_i^n}) \right)$$

and

$$u_{i+\frac{1}{2}}^L = u_i^{n+\frac{1}{2}} + \frac{1}{2} \overline{\Delta_x u_i^n}, \quad u_{i+\frac{1}{2}}^R = u_{i+1}^{n+\frac{1}{2}} - \frac{1}{2} \overline{\Delta_x u_{i+1}^n}.$$

- Evaluate the flux via Riemann solver: Roe's

$$\mathbf{F}_{i+\frac{1}{2}} = \frac{1}{2}(\mathbf{F}_L + \mathbf{F}_R) - \frac{1}{2} \sum_i \alpha_i |\lambda_i| K_i$$

where  $K_i$  are right eigenvectors and  $\alpha_i$  are signal strength.

# MULTI-DIMENSION SCHEMES

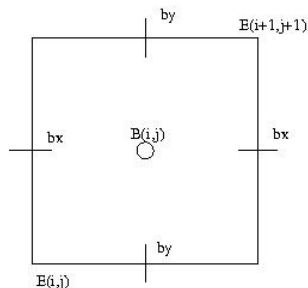
- Dimensional Split
  - Godunov split:  $u^{n+1} = L_{dt}^x L_{dt}^y(u^n)$
  - Strang split  $u^{n+2} = L_{dt}^x L_{dt}^y L_{dt}^y L_{dt}^x(u^n)$
- unsplit scheme:
  - Method of lines (MOL) approach: PDE→ODE, and then Runge-Kutta time integration.
  - Lax-Wendroff Approach: replace the time-derivatives with spatial derivatives, then using the spatial discretization.

The Lax-Wendroff approach has smaller stencil than the MOL approach.

The dimensional split achieves only up to second order schemes. But the CFL limit is more relax for the split scheme than for the unsplit scheme.

# STAGGERED GRID AND CONSTRAINT TRANSPORT

Define the magnetic fields at the face-center, and other physical variables at the cell-center. Define the electro-motive-force (EMF)  $E = v \times B$  at the edge-center. The Yee's staggered grid method preserves the divergence-free condition.



$$b_{x,i+1/2,j}^{n+1} = b_{x,i+1/2,j}^n - \frac{\Delta t}{\Delta y} (E_{i,j+1} - E_{i,j})$$

$$b_{y,i,j+1/2}^{n+1} = b_{y,i,j+1/2}^n + \frac{\Delta t}{\Delta x} (E_{i+1,j} - E_{i,j})$$

$$\nabla \cdot b^n = 0 \Rightarrow \nabla \cdot b^{n+1} = 0,$$

where  $(\nabla \cdot b)_{i,j} :=$

$$(b_{x,i+1/2,j} - b_{x,i-1/2,j})/\Delta x + \\ (b_{y,i,j+1/2} - b_{y,i,j-1/2})/\Delta y$$



# SERIES EXPANSION METHOD

For equation  $L(u) = f(x)$ ,  $a \leq x \leq b$ , where  $L$  is an operator involving partial derivatives of  $u$ . Consider a set of linearly independent basis-function  $\phi_j(x)$ , so that  $u(x) \approx \sum_{j=1}^N U_j \phi_j(x)$ . The residual is

$$R_N \approx L\left(\sum_{j=1}^N U_j \phi_j(x)\right) - f(x).$$

- Galerkin method: The residual is required to be orthogonal to each basis function, which yields

$$\int_a^b R_N \phi_i dx = 0, \Rightarrow \int_a^b \phi_i L\left(\sum_{j=1}^N U_j \phi_j(x)\right) - \int_a^b \phi_i f(x) dx = 0, \quad i = 1, \dots, N$$

which is a set of algebraic equations of  $U_j$ . If  $L$  includes time derivatives, the system becomes ODE's of  $U_j$ .

- Collocation method:  $R_N(x)$  is zero at a set of discrete grid points, i.e.,  $R_N(x_i) = 0$ ,  $i = 1, \dots, N$

# CHOICES OF BASIS FUNCTION

- Spectral methods use orthogonal global series as the basis function:
  - Fourier series:  $e^{-ikx}$
  - Bessel, Chebyshev, Legendre series
- Finite Element Method use local basis function: has local compact support (nonzero region is bounded)
  - piecewise linear: hat-function
  - piecewise quadratic or cubic basis function
  - piecewise higher order Gauss-Legendre polynomials

# SPECTRAL METHOD BY FFT

The fast Fourier Transform (FFT) is widely-used by the spectral method. Consider a periodic function  $u(x)$  on  $[0, 2\pi]$ ,

$$u(x) \approx \sum_{k=1}^N \hat{u}_k e^{ikx}, \text{ where } \hat{u}_k = \frac{1}{N} \sum_{j=1}^N u(x_j) e^{-ikx_j}, \quad x_j = 2j\pi/N, j = 1, \dots, N$$

From the property of the Fourier transform, we can transform the PDEs into an system of ODEs. For a convection-reaction-diffusion equation  $u_t + a(x)u_x + bu = cu_{xx}$ , the corresponding ODEs are

$$\dot{\hat{u}}_k - ik \sum_{k=p+q} \hat{a}_p \hat{u}_q + b\hat{u}_k = -k^2 \hat{u}_k, \quad k = 1, \dots, N$$

where can be integrated by Runge-Kutta method.

The trick part is to evaluate the convolution  $\sum_{k=p+q} \hat{a}_p \hat{u}_q$  efficiently.

# CONVOLUTION EVALUATION IN SPECTRAL METHOD

Let us consider one-dimensional convolution sum

$\hat{w}_k = \sum_{k=p+q} \hat{u}_p \hat{v}_q$ . It involves  $O(N^2)$  operations. If we evaluate  $w(x) = u(x)v(x)$  first in physical space and then apply FFT to  $w(x)$  we obtain

$$\tilde{w}_k = \frac{1}{N} \sum_{x_j} w(x_j) e^{-ikx_j}$$

where requires only  $O(N \log(N))$  operations by FFT.

De-aliasing: Note that

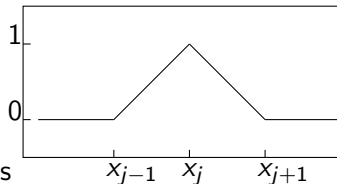
$$\tilde{w}_k = \frac{1}{N} \left[ \sum_{k=p+q} \hat{u}_p \hat{v}_q + \sum_{k+N=p+q} \hat{u}_p \hat{v}_q + \sum_{k-N=p+q} \hat{u}_p \hat{v}_q \right] \neq N \hat{w}_k$$

To recover the formula of  $\hat{w}_k$ , we can apply the **two-third rule**, which forcing all modes  $|k| > N/3$  of  $\hat{u}_k$  and  $\hat{v}_k$  to have zero amplitude.

# FINITE ELEMENT METHOD

We consider the piece-wise linear basis function for 1D advection problem  $u_t + au_x = 0$ , using hat-function:

$$\phi_j(x_i) = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$



The Galerkin approximation gives

$$\sum_{n=1}^N (\phi_n, \phi_k) \frac{dU_n}{dt} + a \sum_{n=1}^N U_n \int_s \frac{d\phi_n}{dx} \phi_k = 0, \text{ for } k = 1, \dots, N$$

where the inner products are

$$(\phi_j, \phi_j) = \int_s \phi_j \phi_j = \frac{2}{3} \Delta x, \quad (\phi_{j-1}, \phi_j) = (\phi_{j+1}, \phi_j) = \int_s \phi_j \phi_{j+1} = \frac{\Delta x}{6},$$

and

$$-\int_s \frac{d\phi_{j-1}}{dx} \phi_j dx = \int_s \frac{d\phi_{j+1}}{dx} \phi_j dx = \frac{1}{2}$$

Since all other integrals of the expansion functions or their derivatives are zero, we have

$$\frac{d}{dt} \left( \frac{U_{j+1} + 4U_j + U_{j-1}}{6} \right) + a \left( \frac{U_{j+1} - U_{j-1}}{2\Delta x} \right) = 0$$

which is identical to the fourth-order compact finite-difference scheme.