

Convergence analysis of finite difference methods for PDEs

Consider the following initial/boundary value problem for a system of linear PDEs

$$\begin{cases} \frac{\partial \mathbf{U}(\mathbf{x}, t)}{\partial t} = \mathbf{L}(\mathbf{x}, t)\mathbf{U}(\mathbf{x}, t) + \mathbf{f}(\mathbf{x}, t) & t \geq 0 & \mathbf{x} \in \Omega \\ \mathbf{S}\mathbf{U}(\mathbf{x}, t) = 0 & \mathbf{x} \in \partial\Omega \\ \mathbf{U}(\mathbf{x}, 0) = \mathbf{U}_0(\mathbf{x}) \end{cases} \quad (1)$$

Here $\mathbf{U}(\mathbf{x}, t)$ denotes a vector field defined in a compact domain $\Omega \subseteq \mathbb{R}^d$, $\partial\Omega$ is the boundary of Ω , \mathbf{L} is a linear operator that can depend on $\mathbf{x} = (x_1, \dots, x_d)$ and t , and \mathbf{S} is a (linear/affine) boundary operator enforcing Dirichlet, Neumann, or mixed boundary conditions. We assume that the IBVP (1) is well-posed, i.e., that the solution exists and is unique in a certain space of functions. Let us provide a few simple examples of PDEs that can be written in the form (1)

- **Liouville equation:** Consider a dynamical system

$$\begin{cases} \frac{d\mathbf{y}}{dt} = \mathbf{F}(\mathbf{y}, t) \\ \mathbf{y}(0) = \mathbf{y}_0 \end{cases} \quad (2)$$

evolving from a random initial state \mathbf{y}_0 with probability density function $p_0(\mathbf{y})$. The PDE governing the evolution equation of the probability density function (PDF) of $\mathbf{y}(t)$ is

$$\frac{\partial p(\mathbf{y}, t)}{\partial t} + \nabla \cdot [\mathbf{F}(\mathbf{y}, t)p(\mathbf{y}, t)] = 0. \quad (3)$$

Clearly, this PDE can be written in the form $\partial p / \partial t = L(\mathbf{x}, t)p$, where

$$L(\mathbf{x}, t)p = -\nabla \cdot \mathbf{F}(\mathbf{x}, t)p - \mathbf{F}(\mathbf{x}, t) \cdot \nabla p \quad (4)$$

is a first-order differential operator that depends on the phase variables \mathbf{y} as well as on time.

- **Fokker-Planck equation:** The Fokker-Planck equation describes the evolution of the probability density function of the state vector solving the stochastic differential equation (SDE)

$$d\mathbf{y} = \boldsymbol{\mu}(\mathbf{y}, t)dt + \sigma d\mathbf{W}(t). \quad (5)$$

Here, $\mathbf{y}(t) \in \mathbb{R}^n$, $\boldsymbol{\mu}(\mathbf{y}, t)$ is the n -dimensional drift, σ is a constant diffusion coefficient and $\mathbf{W}(t)$ is an n -dimensional Wiener random process. The Fokker-Planck equation that corresponds to (5) has the form

$$\frac{\partial p(\mathbf{y}, t)}{\partial t} = -\nabla \cdot [\boldsymbol{\mu}(\mathbf{y}, t)p(\mathbf{y}, t)] + \frac{\sigma^2}{2} \nabla^2 p(\mathbf{y}, t), \quad p(\mathbf{y}, 0) = p_0(\mathbf{y}), \quad (6)$$

where $p_0(\mathbf{y})$ is the PDF of the initial state $\mathbf{y}(0)$.

- **Wave equation:** Consider the wave equation

$$\frac{\partial^2 \psi(\mathbf{x}, t)}{\partial t^2} = c^2 \nabla^2 \psi(\mathbf{x}, t) \quad (7)$$

and the equivalent system of two first-order PDEs as

$$\begin{cases} \frac{\partial \psi(\mathbf{x}, t)}{\partial t} = \eta(\mathbf{x}, t), \\ \frac{\partial \eta(\mathbf{x}, t)}{\partial t} = c^2 \nabla^2 \psi(\mathbf{x}, t). \end{cases} \quad (8)$$

Clearly, the system (8) can be written in the form (1) as

$$\frac{\partial}{\partial t} \underbrace{\begin{bmatrix} \psi(\mathbf{x}, t) \\ \eta(\mathbf{x}, t) \end{bmatrix}}_{\mathbf{U}(\mathbf{x}, t)} = \underbrace{\begin{bmatrix} 0 & 1 \\ c^2 \nabla^2 & 0 \end{bmatrix}}_{\mathbf{L}(\mathbf{x}, t)} \underbrace{\begin{bmatrix} \psi(\mathbf{x}, t) \\ \eta(\mathbf{x}, t) \end{bmatrix}}_{\mathbf{U}(\mathbf{x}, t)} \quad (9)$$

Note that the linear operator $\mathbf{L}(\mathbf{x}, t)$ in this case does not depend on \mathbf{x} and t .

Lax-Richtmyer stability theory

In this section we provide necessary and sufficient conditions for convergence of finite-difference schemes to approximate the solution of the IBVP (1). For simplicity here we consider the case where the linear operator $\mathbf{L}(\mathbf{x}, t)$ in (1) is time-independent, although all consideration in the present discussion apply as well when \mathbf{L} is time-dependent. The fully discrete finite-difference form of the IBVP (1) with time-independent linear operator \mathbf{L} can always be written in the form

$$\mathbf{u}^{k+1} = \mathbf{B}\mathbf{u}^k + \mathbf{b}^k, \quad (10)$$

where \mathbf{u}^k is the vector representing the approximation of the solution $\mathbf{U}(\mathbf{x}, t)$ at all spatial grid points¹ and time t_k . More generally, \mathbf{u}^k represents a vector representing the solution at all grid points and multiple time instants (see the AB2 method described below).

The matrix \mathbf{B} usually depends on Δt , Δx_1 , Δx_2 , etc., and also on the spatial discretization of functions and operators appearing in $\mathbf{L}(\mathbf{x})$. On the other hand, \mathbf{b}^k takes care of external forcing terms and/or boundary conditions. The vector \mathbf{b}^k may also depend of Δt , Δx_1 , Δx_2 , etc. Hence, we shall write

$$\mathbf{B}(\Delta t, \Delta x_1, \Delta x_2, \dots), \quad \mathbf{b}^k(\Delta t, \Delta x_1, \Delta x_2, \dots). \quad (11)$$

Example: Consider the one-dimensional heat-equation

$$\frac{\partial U}{\partial t} = \alpha \frac{\partial^2 U}{\partial x^2}, \quad (12)$$

with zero Dirichlet boundary conditions. Discretize the second derivative $\partial^2 U / \partial x^2$ by, e.g., second-order centered finite differences. This yields the semi-discrete form

$$\frac{d\mathbf{u}(t)}{dt} = \alpha \mathbf{D}_{\text{FD}}^2 \mathbf{u}. \quad (13)$$

We have seen that we can discretize (13) in time using many different schemes, e.g.,

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \mathbf{u}^k \quad (\text{Euler forward}), \quad (14)$$

$$\left(\mathbf{I} - \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right) \mathbf{u}^{k+1} = \left(\mathbf{I} + \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right) \mathbf{u}^k \quad (\text{Crank-Nicolson}). \quad (15)$$

These schemes can be written in the form (10) provided we define

$$\mathbf{B} = \mathbf{I} + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \quad (\text{Euler forward}), \quad (16)$$

$$\mathbf{B} = \left(\mathbf{I} - \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right)^{-1} \left(\mathbf{I} + \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right) \quad (\text{Crank-Nicolson}). \quad (17)$$

¹The numerical solution \mathbf{u}^k in (10) can be arranged as a vector, a matrix, or a multi-dimensional array. Correspondingly, \mathbf{B} can be a matrix, a tensor or a more general linear operator in the space in which \mathbf{u}^k is defined.

Similarly, if we discretize (12) in time with the two-step Adams-Bashforth method we obtain

$$\mathbf{u}^{k+2} = \mathbf{u}^{k+1} + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \left(\frac{3}{2} \mathbf{u}^{k+1} - \frac{1}{2} \mathbf{u}^k \right). \quad (18)$$

We can always write a two-step method as a one-step method in a higher-dimensional space. To this end, define

$$\mathbf{v}^{k+1} = \mathbf{u}^k \quad (19)$$

and rewrite (18) as

$$\begin{cases} \mathbf{u}^{k+2} = \mathbf{u}^{k+1} + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \left(\frac{3}{2} \mathbf{u}^{k+1} - \frac{1}{2} \mathbf{v}^{k+1} \right) \\ \mathbf{v}^{k+2} = \mathbf{u}^{k+1} \end{cases} \quad (20)$$

i.e.,

$$\mathbf{z}^{k+2} = \mathbf{B} \mathbf{z}^{k+1} \quad (21)$$

where

$$\mathbf{z}^{k+2} = \begin{bmatrix} \mathbf{u}^{k+2} \\ \mathbf{v}^{k+2} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{I} + \frac{3}{2} \alpha \Delta t \mathbf{D}_{\text{FD}}^2 & -\frac{1}{2} \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \\ \mathbf{I} & \mathbf{0} \end{bmatrix}. \quad (22)$$

Note that (21) is in the form (10). However, in this case the vector of unknowns \mathbf{z}^k is not just the solution \mathbf{u}^k at time t_k but rather a concatenation of the solution at time t_k and t_{k-1} .

Definition 1 (Lax-Richtmyer stability). The finite difference scheme (10) is said to be stable if there exists a constant C_T independent of k , Δt , Δx_1 , Δx_2 , etc. such that

$$\|\mathbf{B}^k\| \leq C_T \quad \text{for all } k \text{ such that } k \Delta t \leq T. \quad (23)$$

Here, $\|\cdot\|$ denotes any matrix norm induced by a vector norm. In other words, we require that the matrix powers \mathbf{B} are uniformly bounded² by a constant C_T for all $k \leq T/\Delta t$, where the integration period T is fixed. Note that Lax-Richtmyer stability requirement is very similar to the condition we proved for zero-stability of numerical methods for ODEs, where we bounded the the matrix powers of the matrix that propagates the error (see Proposition 1 in Appendix of the course note 4).

Remark: The definition of Lax-Richtmyer stability follows from the following considerations. Consider a perturbation of the fully discrete scheme (10) of the form

$$\mathbf{z}^{k+1} = \mathbf{B} \mathbf{z}^k + \mathbf{b}^k + \boldsymbol{\delta}^k, \quad \mathbf{z}^0 = \mathbf{u}^0 + \boldsymbol{\delta}^0, \quad (25)$$

where $\boldsymbol{\delta}^k$ is assumed to be small, i.e., $\|\boldsymbol{\delta}^k\| \leq \epsilon$. Define the error

$$\mathbf{e}^k = \mathbf{z}^k - \mathbf{u}^k, \quad (26)$$

where \mathbf{u}^k is obtained by iterating (10). It is straightforward to show (see the proof of Theorem 1) that

$$\mathbf{e}^k \leq \|\mathbf{B}^k\| \|\mathbf{e}^0\| + \Delta t \sum_{j=1}^k \|\mathbf{B}^{k-j}\| \|\boldsymbol{\delta}^{j-1}\|. \quad (27)$$

²Uniformly bounded means that the bound

$$\|\mathbf{B}(\Delta t, \Delta x_1, \Delta x_2, \dots)^k\| \leq C_T \quad (24)$$

holds for every Δt , Δx_1 , Δx_2 , etc., and every $k \leq T/\Delta t$ (fixed T).

Hence, if the Lax-Richtmyer stability condition (23) is satisfied we have that

$$\left\| \mathbf{z}^k - \mathbf{u}^k \right\| \leq C_T \|\delta^0\| + C_T T \max_{j=1, \dots, k} \|\delta^{j-1}\| \leq \epsilon C_T (1 + T) \quad \text{for all } k \leq \frac{T}{\Delta t}. \quad (28)$$

In this form, we see Lax-Richtmyer stability (23) can be equivalently formulated as follows: the difference between the solution of the perturbed scheme (25) and the unperturbed one (10) stays bounded for all $k \leq T/\Delta t$, and goes to zero as we send the perturbation amplitude to zero.

When studying convergence of the finite-difference scheme (10) we are interested, in particular, in the behavior of $\|\mathbf{B}^k\|$ when Δt and Δx_i are sent to zero.

Theorem 1 (Lax-Richtmyer equivalence theorem [1]). Given a well-posed initial-boundary value problem (1) and a consistent³ finite-difference approximation (10), stability is a necessary and sufficient condition for convergence. In other words,

$$\boxed{\text{(consistency + stability)} \Leftrightarrow \text{convergence.}} \quad (29)$$

Proof. For simplicity we consider the case where time integration is defined by a one-step scheme although all consideration in the present proof apply also for multistep schemes. A substitution of the exact solution $\mathbf{U}(\mathbf{x}, t)$ of the IBVP (1) into the fully discrete scheme (10) yields the local truncation error (LTE) $\boldsymbol{\tau}^k$, defined by the equation

$$\mathbf{U}^{k+1} = \mathbf{B}\mathbf{U}^k + \mathbf{b}^k + \Delta t \boldsymbol{\tau}^k. \quad (30)$$

Here, we denoted by

$$\mathbf{U}^k = \begin{bmatrix} U(\mathbf{x}_1, t_k) \\ U(\mathbf{x}_2, t_k) \\ \vdots \\ U(\mathbf{x}_N, t_k) \end{bmatrix} \quad (31)$$

where N denotes the total number of spatial grid points⁴. As an example, consider the local truncation error for the forward-in-time centered-in-space finite difference discretization of the heat equation

$$\tau_j^k = \frac{U_j^{k+1} - U_j^k}{\Delta t} - \frac{U_{j+1}^k - 2U_j^k + U_{j-1}^k}{\Delta x^2} \Rightarrow U_j^{k+1} = U_j^k + \frac{\Delta t}{\Delta x^2} (U_{j+1}^k - 2U_j^k + U_{j-1}^k) + \Delta t \tau_j^k. \quad (33)$$

Subtracting (10) from (30) yields

$$\mathbf{e}^{k+1} = \mathbf{B}\mathbf{e}^k + \Delta t \boldsymbol{\tau}^k, \quad (34)$$

where

$$\mathbf{e}^k = \mathbf{U}^k - \mathbf{u}^k \quad (\text{error}) \quad (35)$$

³Recall that a finite-difference approximation is said to be consistent if the local truncation error goes to zero as we send Δt and Δx_i ($i = 1, \dots, d$) to zero.

⁴For example, in 3D we have

$$\mathbf{x}_i = (x_{l(i)}, y_{m(i)}, z_{q(i)}). \quad (32)$$

i.e., we have $N = n^3$ points, where n is the number of points in each variable x , y or z . Note that for multi-step time integration methods \mathbf{U}^k represents the vector of unknown at different times. In other words, everything we saying applies also to multi-step time integration schemes, provided we replace \mathbf{U}^k with $\mathbf{Z}^k = [\mathbf{U}^k \ \mathbf{U}^{k-1} \ \mathbf{U}^{k-2} \ \dots]^T$, i.e., a vector collecting the solution vector at times t_k, t_{k-1} , etc. (see, Eqs. (20)-(22)).

The recursion (34) can be iterated back to the error e^0 . To this end,

$$\begin{aligned}
e^k &= B e^{k-1} + \Delta t \tau^{k-1} \\
&= B \left(B e^{k-2} + \Delta t \tau^{k-2} \right) + \Delta t \tau^{k-1} \\
&= B^2 e^{k-2} + \Delta t B \tau^{k-2} + \Delta t \tau^{k-1} \\
&\vdots \\
&= B^k e^0 + \Delta t \sum_{j=1}^k B^{k-j} \tau^{j-1}.
\end{aligned} \tag{36}$$

At this point we take any vector norm of e^k (and corresponding induced matrix norm), and use the stability assumption (23) to obtain

$$\begin{aligned}
\|e^k\| &= \left\| B^k e^0 + \Delta t \sum_{j=1}^k B^{k-j} \tau^{j-1} \right\| \\
&\leq \|B^k\| \|e^0\| + \Delta t \sum_{j=1}^k \|B^{k-j}\| \|\tau^{j-1}\| \\
&\leq C_T \|e^0\| + C_T \underbrace{k \Delta t}_{\leq T} \max_{j=1, \dots, k} \|\tau^{j-1}\| \\
&\leq C_T \|e^0\| + T C_T \max_{j=1, \dots, k} \|\tau^{j-1}\|,
\end{aligned} \tag{37}$$

where T is period of integration, and C_T is the the uniform bound in (23). The upper bound in (37) goes to zero if the method is consistent, i.e., if

$$\max_{j=1, \dots, k} \|\tau^{j-1}\| \rightarrow 0 \quad \text{for} \quad \Delta t, \Delta x_1, \Delta x_2, \dots \rightarrow 0, \tag{38}$$

and if the error at initial time $\|e^0\|$ is either zero or goes to zero as we send Δt and $\Delta x_1, \Delta x_2$, etc., to zero. This proves that consistency plus Lax-Richtmyer stability implies convergence. \square

Necessary and sufficient conditions for stability. Recall that for any matrix B , any matrix norm $\|\cdot\|$, and any integer k we have

$$\|B^k\| \leq \|B\|^k. \tag{39}$$

Therefore, to prove stability of (10) it is sufficient to show that for sufficiently small Δt

$$\|B\| \leq 1 + \beta \Delta t \quad \text{for some } \beta \in \mathbb{R}. \tag{40}$$

In fact,

$$\|B^k\| \leq \|B\|^k \leq (1 + \beta \Delta t)^k \leq e^{k \Delta t \beta} \leq e^{T \beta}. \tag{41}$$

The last two inequalities follow from the basic inequality

$$\log(1+x) \leq x \quad \Leftrightarrow \quad \log\left[(1+x)^k\right] = k \log(1+x) \leq kx \quad \Leftrightarrow \quad (1+x)^k \leq e^{kx}. \tag{42}$$

Of course, if the matrix B is a contraction in some norm, i.e., if

$$\|B\| \leq 1, \tag{43}$$

then the finite-difference scheme (10) is Lax-Richtmyer stable. Next, recall that the spectral radius of a matrix

$$\rho(\mathbf{B}) = \max \{|\lambda_1|, \dots, |\lambda_n|\}. \quad (44)$$

represents a lower bound for any natural matrix norm⁵, i.e.,

$$\rho(\mathbf{B}) \leq \|\mathbf{B}\| \Rightarrow \rho(\mathbf{B})^k \leq \|\mathbf{B}^k\| \leq C_T \Leftrightarrow \rho(\mathbf{B})^k \leq C_T. \quad (46)$$

From this equation it follows that

$$\rho(\mathbf{B}) \leq 1 + \beta\Delta t \quad (47)$$

is necessary for stability. Hence, the modulus of the eigenvalues of the matrix \mathbf{B} must always be smaller than $1 + \beta\Delta t$. If the matrix \mathbf{B} is normal (i.e. $\mathbf{B}\mathbf{B}^T = \mathbf{B}^T\mathbf{B}$) then the 2-norm coincides with the spectral radius, i.e.

$$\rho(\mathbf{B}) = \|\mathbf{B}\|_2 \quad (48)$$

and (47) is *necessary and sufficient for stability*.

Remark: The norm of the matrix \mathbf{B} depends on the discretization scheme, on Δx_i ($i = 1, 2, \dots$), and on Δt . Hence, condition (40) effectively sets a condition on all these variables for the scheme to be stable. Hereafter we provide some examples.

An example of Lax-Richtmyer stability analysis

Consider the one-dimensional heat equation (12) with zero Dirichlet boundary conditions, and the matrix \mathbf{B} corresponding to the forward-in-time centered-in-space finite-difference discretization (16). The matrix \mathbf{B} is symmetric, and therefore the 2-norm coincides with the spectral radius $\rho(\mathbf{B})$. This yields⁶,

$$\|\mathbf{B}\|_2 = \max_{i=1, \dots, N} |\alpha\Delta t\lambda_i + 1| \quad (49)$$

where λ_i are the the eigenvalues of the differentiation matrix \mathbf{D}_{FD}^2 . At this point we recall that all the eigenvalues λ_i are negative, and that

$$\min_{i=1, \dots, N} \lambda_i = \lambda_N = \frac{2}{\Delta x^2} (\cos(N\pi\Delta x) - 1) \geq -\frac{4}{\Delta x^2}. \quad (50)$$

Hence,

$$\|\mathbf{B}\|_2 \leq \left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right|. \quad (51)$$

Recalling equation (40), we conclude that a necessary and sufficient condition for stability is

$$\left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \leq 1 + \beta\Delta t. \quad (52)$$

This equation defines a *stability region* in the $(\Delta t, \Delta x)$ -plane for each β (see Figure 23). Such a stability region can be computed analytically, although the computation is a bit cumbersome, except for the case $\beta = 0$. In the latter case we have

$$\left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \leq 1 \Leftrightarrow \frac{\Delta t}{\Delta x^2} \leq \frac{1}{2\alpha} \quad (\text{Lax-Richtmyer stability condition}). \quad (53)$$

⁵A natural matrix norm is a matrix norm induced by vector norm as

$$\|\mathbf{B}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{B}\mathbf{x}\|}{\|\mathbf{x}\|}. \quad (45)$$

⁶The eigenvalues of the matrix $\mathbf{B} = \mathbf{I} + \alpha\Delta t\mathbf{D}_{\text{FD}}^2$ are $1 + \alpha\Delta t\lambda_i$, where λ_i are the eigenvalues of \mathbf{D}_{FD}^2 .

Condition (53) shows that Δt and Δx cannot be sent to zero at arbitrary rates. Indeed, we must have $\Delta t \sim \kappa \Delta x^2$ for (53) to hold in the limits $\Delta t, \Delta x \rightarrow 0$.

Lax-Richtmyer stability analysis can be applied to both implicit and explicit temporal integration schemes. However, for implicit schemes the matrix \mathbf{B} involves the inverse of some other matrix (see, e.g., equation (17)). This makes the stability analysis of implicit schemes not straightforward nor practical using the matrix \mathbf{B} . We shall see hereafter that this issue can be mitigated (at least for linear PDEs) by using discrete Fourier series (for periodic boundary conditions), which are at the basis of the *Von-Neumann stability analysis*.

Convergence analysis for nonlinear PDEs. The fully discrete finite-difference formulation of a one-dimensional nonlinear PDE can be written as

$$\sum_{k=0}^q \alpha_k u_j^{k+q} = \Delta t \Phi_j \left(\mathbf{u}^{k+q}, \dots, \mathbf{u}^k, \Delta t, \Delta x \right). \quad (54)$$

For instance, the second-order central finite-difference discretization of the Kuramoto-Sivashinsky equation with Euler forward time stepping can be written as

$$u_j^{k+1} - u_j^k = \Delta t \left(-u_j^k \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} - \frac{u_{j-1}^k - 2u_j^k + u_{j+1}^k}{\Delta x^2} - \frac{u_{j-2}^k - 3u_{j-1}^k + 6u_j^k - 3u_{j+1}^k + u_{j+2}^k}{\Delta x^4} \right), \quad (55)$$

To study convergence of this scheme for Δt and Δx going to zero, we can use methods similar to the those we used in the convergence analysis of numerical schemes for ODEs, in particular the convergence proof in the appendix of course note 4. To use such a proof in the context of finite-difference (or other) approximations of PDEs, we need to make sure that the Lipschitz constant of Φ in (54) can be bounded by some constant when we send Δt and Δx to zero at an appropriate rate. Under this assumption, it is rather straightforward to show that the method is convergent (provided the method it is consistent). To this end, just follow the proof in the appendix of the course note 4.

Von-Neumann stability theory

Stability analysis of finite-difference schemes can be simplified substantially if the PDE is defined in a periodic spatial domain. The key idea is to use discrete Fourier series applied to the finite-difference discretization of the PDE and determine under which conditions on Δt , Δx_1 , Δx_2 , etc., the scheme is “stable”. The idea is to write the finite difference form of the PDE as a recurrence in Fourier space, i.e., as

$$\mathbf{c}^k = \mathbf{G} \mathbf{c}^{k-1} + \mathbf{v}^{k-1}, \quad (56)$$

where \mathbf{c}^k are the Fourier coefficient of the finite-difference numerical solution, and \mathbf{v}^{k-1} are Fourier coefficient of some external forcing term. With this form available, we can apply the same reasoning of Theorem 1 to conclude that the scheme is convergent if and only if it is stable and consistent. Stability here refers to the norm of the matrix \mathbf{G} which defines the discrete Fourier transform of the finite-difference scheme. As we shall see hereafter there is a one-to-one correspondence between the matrix \mathbf{G} and the matrix \mathbf{B} in the previous section. In particular, \mathbf{G} and \mathbf{F} are unitarily equivalent via the Fourier transformation matrix \mathbf{F} , i.e.,

$$\mathbf{G} = \mathbf{F}^{-1} \mathbf{B} \mathbf{F}. \quad (57)$$

One reason for the Fourier series analysis is that it allows us to determine stability conditions for both implicit and explicit schemes is a rather straightforward way.

Example: To illustrate how the discrete Fourier transform method works, consider the prototype IBVP

$$\begin{cases} \frac{\partial U}{\partial t} = \alpha \frac{\partial^2 U}{\partial x^2} & t \geq 0 & x \in [0, L] \\ U(x, 0) = U_0(x) \\ \text{Periodic B.C.} \end{cases} \quad (58)$$

We have seen that (58) can be discretized with second-order finite-differences in space and Euler-forward time integration as

$$u_j^{k+1} = u_j^k + \frac{\alpha \Delta t}{\Delta x^2} \left(u_{j-1}^k - 2u_j^k + u_{j+1}^k \right) \quad j = 0, \dots, N-1, \quad (59)$$

where u_j^k is the approximation of $U(x_j, t_k)$, and

$$x_j = j \frac{L}{N} \quad j = 0, \dots, N. \quad (60)$$

The scheme (59) is supplemented with periodic boundary conditions

$$u_j^k = u_{j+N}^k \quad \text{for all } j \in \mathbb{Z}. \quad (61)$$

and the initial condition

$$u_j^0 = U_0(x_j) \quad j = 0, \dots, N-1. \quad (62)$$

Let us now expand the numerical solution u_j^k in a discrete Fourier series⁷

$$u_j^k = \frac{1}{N} \sum_{p=0}^{N-1} c_p^k e^{ipj\xi}, \quad \text{where} \quad \xi = \frac{2\pi \Delta x}{L} \quad \text{and} \quad c_p^k \in \mathbb{C}. \quad (67)$$

A substitution on (67) into (59) yields

$$\begin{aligned} \sum_{p=0}^{N-1} c_p^{k+1} e^{ipj\xi} &= \sum_{p=0}^{N-1} c_p^k e^{ipj\xi} \left[1 + \frac{\alpha \Delta t}{\Delta x^2} \left(e^{-ip\xi} - 2 + e^{ip\xi} \right) \right] \\ &= \sum_{p=0}^{N-1} c_p^k e^{ipj\xi} \left[1 + \frac{\alpha \Delta t}{\Delta x^2} (2 \cos(p\xi) - 2) \right], \end{aligned} \quad (68)$$

⁷As is well known, the solution to (58) can be expanded in a Fourier series as

$$U(x, t) = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k x / L} \quad (63)$$

Evaluating $U(x, t)$ on the grid (60) yields the discrete Fourier series

$$U(x_j, t) = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k x_j / L} = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k j / N} \quad j = 0, \dots, N-1. \quad (64)$$

Moreover,

$$U(x_j, t) = \sum_{k=-\infty}^{\infty} C_k(t) e^{2ikj\pi/N} = \sum_{k=0}^{N-1} \sum_{p=-\infty}^{\infty} C_{k+pN}(t) e^{2\pi i(k+pN)j/N} = \sum_{k=0}^{N-1} e^{2\pi i k j / N} \underbrace{\sum_{p=-\infty}^{\infty} C_{k+pN}(t)}_{c_k(t)/N}. \quad (65)$$

This can be written as

$$U(x_j, t) = \frac{1}{N} \sum_{k=0}^{N-1} c_k(t) e^{ikj\xi}, \quad \xi = \frac{2\pi \Delta x}{L} = \frac{2\pi}{N} \quad (66)$$

i.e.,

$$c_p^{k+1} = c_p^k \underbrace{\left[1 + \frac{2\alpha\Delta t}{\Delta x^2} \left(\cos\left(2\pi p \frac{\Delta x}{L}\right) - 1 \right) \right]}_{\text{"amplification factor" } G_p(\Delta t, \Delta x)}. \quad (69)$$

Upon definition of

$$\mathbf{c}^k = \begin{bmatrix} c_0^k \\ c_0^k \\ \vdots \\ c_{N-1}^k \end{bmatrix}, \quad \mathbf{G}(\Delta t, \Delta x) = \begin{bmatrix} G_0(\Delta t, \Delta x) & 0 & \cdots & 0 \\ 0 & G_1(\Delta t, \Delta x) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & G_{N-1}(\Delta t, \Delta x) \end{bmatrix} \quad (70)$$

we can write (69) as

$$\mathbf{c}^{k+1} = \mathbf{G}(\Delta t, \Delta x) \mathbf{c}^k. \quad (71)$$

From this equation we see that the matrix $\mathbf{G}(\Delta t, \Delta x)$ in (71) plays the same role in Fourier space as the matrix \mathbf{B} in (10) does in physical space. Of course, any finite-difference scheme for an arbitrary linear PDE can be written in the form (71). This prompts the following definition.

Definition 2 (Von-Neumann stability). The finite difference scheme is said to be stable if there exists a constant H_T independent of k , Δt , Δx_1 , Δx_2 , etc. such that

$$\|\mathbf{G}^k\| \leq H_T \quad \text{for all } k \text{ such that } k\Delta t \leq T, \quad (72)$$

where \mathbf{G} is the Von-Neumann matrix appearing in (71) or (56).

Equivalence between Von-Neumann and Lax-Richtmyer stability conditions. In this section we show that the definitions of Von-Neumann stability (Eq. (72)) and Lax-Richtmyer stability (Eq. (23)) are equivalent for PDEs with periodic boundary conditions.

To this end, we notice that if we compute the inverse Fourier transform of (71) we obtain

$$\mathbf{u}^{k+1} = \mathbf{F} \mathbf{G} \mathbf{F}^{-1} \mathbf{u}^k, \quad (73)$$

where \mathbf{F} is the Fourier transform matrix such that

$$\mathbf{u}^k = \frac{1}{\sqrt{N}} \mathbf{F} \mathbf{c}^k, \quad \mathbf{F} = \frac{1}{\sqrt{N}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & e^{i\xi} & e^{2i\xi} & \cdots & e^{i(N-1)\xi} \\ 1 & e^{2i\xi} & e^{4i\xi} & \cdots & e^{2i(N-1)\xi} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & e^{(N-1)i\xi} & e^{2(N-1)\xi} & \cdots & e^{i(N-1)^2\xi} \end{bmatrix} \quad (74)$$

A comparison between (73) and (10) shows that

$$\mathbf{B} = \mathbf{F} \mathbf{G} \mathbf{F}^{-1}. \quad (75)$$

This equation justifies why stability can be equivalently studied in Fourier space by studying the norm of \mathbf{G}^k . In fact,

$$\|\mathbf{B}^k\| = \|\mathbf{F} \mathbf{G}^k \mathbf{F}^{-1}\|. \quad (76)$$

The Fourier transform matrix \mathbf{F} plays no role in the stability properties of the scheme. In fact, the DFT matrix \mathbf{F} is *unitary*, i.e.,

$$\mathbf{F} \mathbf{F}^* = \mathbf{F}^* \mathbf{F} = \mathbf{I} \quad (\mathbf{F}^* \text{ is the conjugate transpose of } \mathbf{F}). \quad (77)$$

This can be verified directly using the definition (74). We know that the singular values and the eigenvalues of a matrix do not change under unitary transformations. This implies that \mathbf{B}^k and \mathbf{G}^k have exactly the same eigenvalues and the same singular values. Hence the matrices \mathbf{B}^k and \mathbf{G}^k have the same spectral radius and the same 2-norm

$$\left\| \mathbf{B}^k \right\|_2 = \left\| \mathbf{G}^k \right\|_2. \quad (78)$$

Example: The matrix \mathbf{G} that corresponds to the finite difference scheme

Necessary and sufficient conditions for Von-Neumann stability

Let us recall that the spectral radius of a matrix \mathbf{G} , i.e.,

$$\rho(\mathbf{G}) = \max_i |\lambda_i| \quad (79)$$

where λ_i are the eigenvalues of \mathbf{G} , is a lower bound for any natural matrix norm⁸ of \mathbf{G} , i.e.,

$$\rho(\mathbf{G}) \leq \|\mathbf{G}\| \quad \text{for every natural matrix norm } \|\cdot\|. \quad (83)$$

Moreover, the spectral radius of the matrix power \mathbf{G}^k is equal to $\rho(\mathbf{G})^k$ (recall that the eigenvalues of \mathbf{G}^k are λ_i^k). By using the stability condition (72) we obtain

$$\rho(\mathbf{G})^k \leq \left\| \mathbf{G}^k \right\| \leq H_T \quad (84)$$

i.e.,

$$\rho(\mathbf{G})^k \leq H_T. \quad (85)$$

As before, this implies that for sufficiently small Δt the spectral radius of \mathbf{G} must satisfy (see Eq. (40))

$$\rho(\mathbf{G}) \leq 1 + \gamma \Delta t \quad \gamma \in \mathbb{R}. \quad (86)$$

This is a *necessary condition*, not a sufficient condition. In fact, it is possible that $\rho(\mathbf{G})^k \leq H_T$ even though $\left\| \mathbf{G}^k \right\|$ grows unboundedly as we send Δt , Δx_1 , Δx_2 , etc., to zero. In other words,

$$\rho(\mathbf{G})^k \leq H_T \quad \text{does not imply} \quad \left\| \mathbf{G}^k \right\| \leq H_T. \quad (87)$$

However, if the matrix \mathbf{G} is normal, i.e., if $\mathbf{G}\mathbf{G}^* = \mathbf{G}^*\mathbf{G}$ (where $*$ denotes the conjugate transpose) then it is easy to show that Von-Neumann stability condition (86) is sufficient.

Lemma 1. The Von-Neumann stability condition (86) is sufficient if the amplification matrix \mathbf{G} is normal.

⁸All natural matrix norms, i.e., matrix norms defined in terms of vector norms as

$$\|\mathbf{A}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \quad (80)$$

are sub-multiplicative. This means that they

$$\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad (81)$$

for all matrices \mathbf{A} and \mathbf{B} . To prove (83) we leverage the sub-multiplicative property of natural matrix norm. To this end, consider one eigenvalue λ_i of \mathbf{G} and the corresponding eigenvector \mathbf{v} . Construct the matrix $\mathbf{V} = [\mathbf{v} \ \cdots \ \mathbf{v}]$, and note that

$$\|\mathbf{G}\mathbf{V}\| = |\lambda_i| \|\mathbf{V}\| \leq \|\mathbf{G}\| \|\mathbf{V}\| \quad \Rightarrow \quad \|\mathbf{G}\| \geq \max_i |\lambda_i| = \rho(\mathbf{G}). \quad (82)$$

the matrix

Proof. The spectral radius of normal matrices is equal to the matrix 2-norm

$$\rho(\mathbf{G}) = \sqrt{\rho(\mathbf{G}\mathbf{G}^*)} = \|\mathbf{G}\|_2. \quad (88)$$

This allows us to write (84) as

$$\rho(\mathbf{G})^k = \left\| \mathbf{G}^k \right\|_2 \leq H_T. \quad (89)$$

Hence, for normal matrices \mathbf{G} we have that (86) implies

$$\left\| \mathbf{G}^k \right\|_2 \leq H_T \quad (90)$$

i.e., that the scheme is stable. □

For normal matrices \mathbf{G} a simplified condition that grants us stability of the scheme is

$$\rho(\mathbf{G}) \leq 1. \quad (91)$$

Examples of Von-Neumann stability analysis

The Fourier series of the solution of linear PDEs with constant coefficients can be always decoupled into a system of equations involving one Fourier mode at a time. Hence, to determine the amplification factors of the Fourier coefficients it is sufficient to consider only mode. In practice, we can simply substitute

$$u_j^k = c_p^k e^{ijp\xi} \quad \text{where} \quad \xi = \frac{2\pi\Delta x}{L} \quad (92)$$

into the numerical scheme and compute the amplification factors for the p -th mode. Let us show how to perform this calculation for second-order centered finite-difference discretization of the heat equation with Crank-Nicolson time-integration.

Heat equation with Euler-forward time integration. We have seen that the Fourier transform of finite-difference scheme (59) yields the diagonal matrix of amplification factors defined in (70). The diagonal entries of \mathbf{G} are the eigenvalues of \mathbf{G} . Hence, the spectral radius of \mathbf{G} is

$$\rho(\mathbf{G}) = \max_{p=0,\dots,N-1} \left| 1 + \frac{2\alpha\Delta t}{\Delta x^2} \left(\cos\left(2\pi p \frac{\Delta x}{L}\right) - 1 \right) \right| \leq \left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right|. \quad (93)$$

By using the Von-Neumann condition (86) we conclude that the scheme (59) is stable if and only if

$$\left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \leq 1 + \gamma\Delta t. \quad (94)$$

This is exactly the same condition we obtained in (52) (see the discussion thereafter).

Heat equation with Crank-Nicolson time integration. Consider the fully discrete finite-difference scheme

$$u_j^{k+1} - \frac{\alpha\Delta t}{2\Delta x^2} \left(u_{j-1}^{k+1} - 2u_j^{k+1} + u_{j+1}^{k+1} \right) = u_j^k + \frac{\alpha\Delta t}{2\Delta x^2} \left(u_{j-1}^k - 2u_j^k + u_{j+1}^k \right). \quad (95)$$

Substitute (92) into (95) to obtain

$$c_p^{k+1} \left[1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right] = c_p^k \left[1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right], \quad (96)$$

i.e.,

$$c_p^{k+1} = \frac{1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)}{1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)} c_p^k. \quad (97)$$

Again, the amplification matrix \mathbf{G} is diagonal with spectral radius

$$\rho(\mathbf{G}) = \max_{p=0, \dots, N-1} \left| \frac{1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)}{1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)} \right|. \quad (98)$$

At this point we notice that $\cos(p\xi) - 1 \leq 0$ for any p . This implies that

$$\left| 1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right| \leq \left| 1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right| \quad (99)$$

and

$$\rho(\mathbf{G}) \leq 1. \quad (100)$$

Recalling the Von-Neumann stability condition (86) we conclude that the second-order centered finite-difference scheme with Crank-Nicolson time integration is unconditionally stable. Moreover, the scheme is consistent, and therefore convergent.

References

- [1] P. D. Lax and R. D. Richtmyer. Survey of the stability of linear finite-difference equations. *Communications on Pure and Applied Mathematics*, 9:267–293, 1956.