

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, Robin or mixed boundary conditions. We assume that the IBVP (1) is well-posed, i.e., that it admits a unique solution. 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 joint probability density function 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

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

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

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

$$\frac{\partial^2 \psi(\mathbf{x}, t)}{\partial t^2} = c^2 \nabla^2 \psi(\mathbf{x}, t) \quad (5)$$

and the equivalent system of two first-order PDEs as

$$\begin{cases} \frac{\partial \psi(\mathbf{x}, t)}{\partial t} = \eta(\mathbf{x}, t), \\ \frac{\partial \eta(\mathbf{x}, t)}{\partial t} = c^2 \nabla^2 \psi(\mathbf{x}, t). \end{cases} \quad (6)$$

Clearly, the system (6) can be written in the form (1) as

$$\frac{\partial}{\partial t} \underbrace{\begin{bmatrix} \psi(\mathbf{x}, t) \\ \eta(\mathbf{x}, t) \end{bmatrix}}_{\mathbf{U}(\mathbf{x}, t)} = \underbrace{\begin{bmatrix} 0 & 1 \\ c^2 \nabla^2 & 0 \end{bmatrix}}_{\mathbf{L}(\mathbf{x}, t)} \underbrace{\begin{bmatrix} \psi(\mathbf{x}, t) \\ \eta(\mathbf{x}, t) \end{bmatrix}}_{\mathbf{U}(\mathbf{x}, t)} \quad (7)$$

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). In the interest of simplicity 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

$$\mathbf{u}^{k+1} = \mathbf{B}\mathbf{u}^k + \mathbf{b}^k, \quad (8)$$

where \mathbf{u}^k is the vector representing the approximation of the solution $U(x, t)$ at *all* grid points¹ and time t_k , or (more generally) a vector representing 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 the functions and operators appearing in $\mathbf{L}(\mathbf{x})$, while \mathbf{b}^k takes care of external forcing terms and/or the boundary conditions. The vector \mathbf{b}^k may also depend of Δt , Δx_1 , Δx_2 , etc.

Example: Consider the one-dimensional heat-equation

$$\frac{\partial U}{\partial t} = \alpha \frac{\partial^2 U}{\partial x^2}, \quad (9)$$

with 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

$$\frac{d\mathbf{u}(t)}{dt} = \alpha \mathbf{D}_{\text{FD}}^2 \mathbf{u}. \quad (10)$$

We have seen that we can discretize (10) in time using many different schemes, e.g.,

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \mathbf{u}^k \quad (\text{Euler forward}), \quad (11)$$

$$\left(\mathbf{I} - \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right) \mathbf{u}^{k+1} = \left(\mathbf{I} + \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right) \mathbf{u}^k \quad (\text{Crank-Nicolson}). \quad (12)$$

These schemes can be written in the form (8) provided we define

$$\mathbf{B} = \mathbf{I} + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \quad (\text{Euler forward}), \quad (13)$$

$$\mathbf{B} = \left(\mathbf{I} - \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right)^{-1} \left(\mathbf{I} + \frac{\alpha \Delta t}{2} \mathbf{D}_{\text{FD}}^2 \right) \quad (\text{Crank-Nicolson}). \quad (14)$$

Similarly, if we discretize (9) in time with the two-step Adams-Bashforth method we obtain

$$\mathbf{u}^{k+2} = \mathbf{u}^{k+1} + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \left(\frac{3}{2} \mathbf{u}^{k+1} - \frac{1}{2} \mathbf{u}^k \right). \quad (15)$$

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

$$\mathbf{v}^{k+1} = \mathbf{u}^k \quad (16)$$

¹The numerical solution \mathbf{u}^k in (8) 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.

and rewrite (15) as

$$\begin{cases} \mathbf{u}^{k+2} = \mathbf{u}^{k+1} + \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \left(\frac{3}{2} \mathbf{u}^{k+1} - \frac{1}{2} \mathbf{v}^{k+1} \right) \\ \mathbf{v}^{k+2} = \mathbf{u}^{k+1} \end{cases} \quad (17)$$

i.e.,

$$\mathbf{z}^{k+2} = \mathbf{B} \mathbf{z}^{k+1} \quad (18)$$

where

$$\mathbf{z}^{k+2} = \begin{bmatrix} \mathbf{u}^{k+2} \\ \mathbf{v}^{k+2} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{I} + \frac{3}{2} \alpha \Delta t \mathbf{D}_{\text{FD}}^2 & -\frac{1}{2} \alpha \Delta t \mathbf{D}_{\text{FD}}^2 \\ \mathbf{I} & \mathbf{0} \end{bmatrix}. \quad (19)$$

Note that (18) is in the form (8). 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 (8) is *stable* if there exists a constant C_T independent of k , Δt , Δx_1 , Δx_2 etc such that

$$\|\mathbf{B}^k\| \leq C_T \quad \text{for all } k \text{ such that } k\Delta t \leq T. \quad (20)$$

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 and is chosen arbitrarily Δt .

When studying convergence of (8) 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 properly posed initial-boundary value problem (1) and a consistent³ finite-difference approximation (8), stability is a necessary and sufficient condition for convergence.*

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 as well for multistep schemes. A substitution of the exact solution $\mathbf{U}(\mathbf{x}, t)$ of the IBVP (1) into the fully discrete scheme (8) yields the local truncation error (LTE) $\boldsymbol{\tau}^k$, defined by the equation

$$\mathbf{U}^{k+1} = \mathbf{B} \mathbf{U}^k + \mathbf{b}^k + \Delta t \boldsymbol{\tau}^k. \quad (22)$$

Here, we denoted by

$$\mathbf{U}^k = \begin{bmatrix} U(\mathbf{x}_1, t_k) \\ U(\mathbf{x}_2, t_k) \\ \vdots \\ U(\mathbf{x}_N, t_k) \end{bmatrix} \quad (23)$$

where N denotes the total number of spatial grid points. For example, in 3D we have

$$\mathbf{x}_i = (x_{l(i)}, y_{m(i)}, z_{n(i)}). \quad (24)$$

²Uniformly bounded means that the bound

$$\|\mathbf{B}(\Delta t, \Delta x_1, \Delta x_2, \dots)^k\| \leq C_T \quad (21)$$

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

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.

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 we just need to replace \mathbf{U}^k by $\mathbf{Z}^k = [U^k, U^{k-1}, U^{k-2}, \dots]^T$, i.e., a vector collecting the solution vector at times t_k, t_{k-1} , etc. (see, e.g., Eqs. (17)-(19)). Subtracting (8) from (22) yields

$$\mathbf{e}^{k+1} = \mathbf{B}\mathbf{e}^k + \Delta t\boldsymbol{\tau}^k, \quad (25)$$

where

$$\mathbf{e}^k = \mathbf{U}^k - \mathbf{u}^k \quad (\text{error}) \quad (26)$$

The recursion (25) can be iterated back to the error \mathbf{e}^0 . To this end,

$$\begin{aligned} \mathbf{e}^k &= \mathbf{B}\mathbf{e}^{k-1} + \Delta t\boldsymbol{\tau}^{k-1} \\ &= \mathbf{B}(\mathbf{B}\mathbf{e}^{k-2} + \Delta t\boldsymbol{\tau}^{k-2}) + \Delta t\boldsymbol{\tau}^{k-1} \\ &= \mathbf{B}^2\mathbf{e}^{k-2} + \Delta t\mathbf{B}\boldsymbol{\tau}^{k-2} + \Delta t\boldsymbol{\tau}^{k-1} \\ &\vdots \\ &= \mathbf{B}^k\mathbf{e}^0 + \Delta t \sum_{j=1}^k \mathbf{B}^{k-j}\boldsymbol{\tau}^{j-1}. \end{aligned} \quad (27)$$

At this point we take any vector norm of \mathbf{e}^k (and corresponding induced matrix norm), and use the stability assumption (20) to obtain

$$\begin{aligned} \|\mathbf{e}^k\| &= \left\| \mathbf{B}^k\mathbf{e}^0 + \Delta t \sum_{j=1}^k \mathbf{B}^{k-j}\boldsymbol{\tau}^{j-1} \right\| \\ &\leq \|\mathbf{B}^k\| \|\mathbf{e}^0\| + \Delta t \sum_{j=1}^k \|\mathbf{B}^{k-j}\| \|\boldsymbol{\tau}^{j-1}\| \\ &\leq C_T \|\mathbf{e}^0\| + k\Delta t C_T \max_{j=1, \dots, k} \|\boldsymbol{\tau}^{j-1}\| \\ &\leq C_T \|\mathbf{e}^0\| + TC_T \max_{j=1, \dots, k} \|\boldsymbol{\tau}^{j-1}\|, \end{aligned} \quad (28)$$

where T is period of integration, and C_T is the the uniform bound in (20). The upper bound in (28) goes to zero if the method is consistent, i.e., if

$$\max_{j=1, \dots, k} \|\boldsymbol{\tau}^{j-1}\| \rightarrow 0 \quad \text{for} \quad \Delta t, \Delta x_i \rightarrow 0, \quad (29)$$

and if the error at initial time $\|\mathbf{e}^0\|$ 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

At this point a few remarks are in order.

- **Sufficient condition for stability:** Recall that for any matrix norm and any $k \in \mathbb{N}$

$$\|\mathbf{B}^k\| \leq \|\mathbf{B}\|^k. \quad (30)$$

Therefore, to prove stability it is sufficient to show that for sufficiently small Δt

$$\|\mathbf{B}\| \leq 1 + \beta\Delta t \quad \text{for some } \beta \in \mathbb{R}. \quad (31)$$

In fact⁴,

$$\left\| \mathbf{B}^k \right\| \leq \|\mathbf{B}\|^k \leq (1 + \beta\Delta t)^k \leq e^{k\Delta t\beta} \leq e^{T\beta}. \quad (33)$$

- **Necessary and sufficient conditions for stability:** The spectral radius of a matrix is a lower bound for any matrix sub-multiplicative matrix norm. This implies that

$$\rho(\mathbf{B}) \leq \|\mathbf{B}\| \quad \Rightarrow \quad \rho(\mathbf{B})^k \leq \left\| \mathbf{B}^k \right\| \leq C_T \quad \Leftrightarrow \quad \rho(\mathbf{B})^k \leq C_T. \quad (34)$$

From this equation it follows that

$$\rho(\mathbf{B}) \leq 1 + \beta\Delta t \quad (35)$$

is necessary for stability. 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.

$$\rho(\mathbf{B}) = \|\mathbf{B}\|_2 \quad (36)$$

and (35) is *necessary and sufficient for stability*.

Stability analysis of forward-in-time centered-in-space scheme for the heat equation: It is important to remark that the stability condition may depend on the way we send Δt and Δx_i to zero. To show this, consider the one-dimensional heat equation (9) with zero Dirichlet boundary conditions, and the matrix \mathbf{B} corresponding to the centered-in-space forward-in-time finite-difference discretization (13). The matrix \mathbf{B} is symmetric, and therefore the 2-norm coincides with the spectral radius. This yields⁵,

$$\|\mathbf{B}\|_2 = \max_{i=1,\dots,N} |\alpha\Delta t\lambda_i + 1|. \quad (37)$$

At this point we recall that, for large N (number of spatial grid points), we have

$$\min_{i=1,\dots,N} \lambda_i = \lambda_N = \frac{2}{\Delta x^2} (\cos(N\pi\Delta x) - 1) \simeq -\frac{4}{\Delta x^2}. \quad (38)$$

Hence,

$$\|\mathbf{B}\|_2 \simeq \left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \quad (39)$$

Recalling equation (31), we conclude that a necessary and sufficient condition for stability is

$$\left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \leq 1 + \beta\Delta t. \quad (40)$$

This equation defines a *stability region* in the $(\Delta t, \Delta x)$ -plane for each β (see Figure 20). Such a stability region can be computed analytically, although the computation is a bit cumbersome (except for the case $\beta = 0$). In fact, define

$$\eta = \frac{\Delta t}{\Delta x^2}. \quad (41)$$

The inequality (40) can be split into the following two inequalities

$$\begin{cases} 4\alpha\eta - \beta\eta\Delta x^2 \leq 2 & \Leftrightarrow & \eta(4\alpha - \beta\Delta x^2) \leq 2 & \text{for } \Delta x^2 \leq 4\alpha\Delta t, \\ -4\alpha\eta - \beta\eta\Delta x^2 \leq 0 & \Leftrightarrow & 4\alpha + \beta\Delta x^2 \geq 0 & \text{for } \Delta x^2 \geq 4\alpha\Delta t. \end{cases} \quad (42)$$

⁴The inequalities (33) follow from the basic inequality

$$\log(1+x) \leq x \quad \Leftrightarrow \quad \log(1+x)^k = k \log(1+x) \leq kx \quad \Leftrightarrow \quad (1+x)^k \leq e^{kx} \quad (32)$$

⁵The 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 .

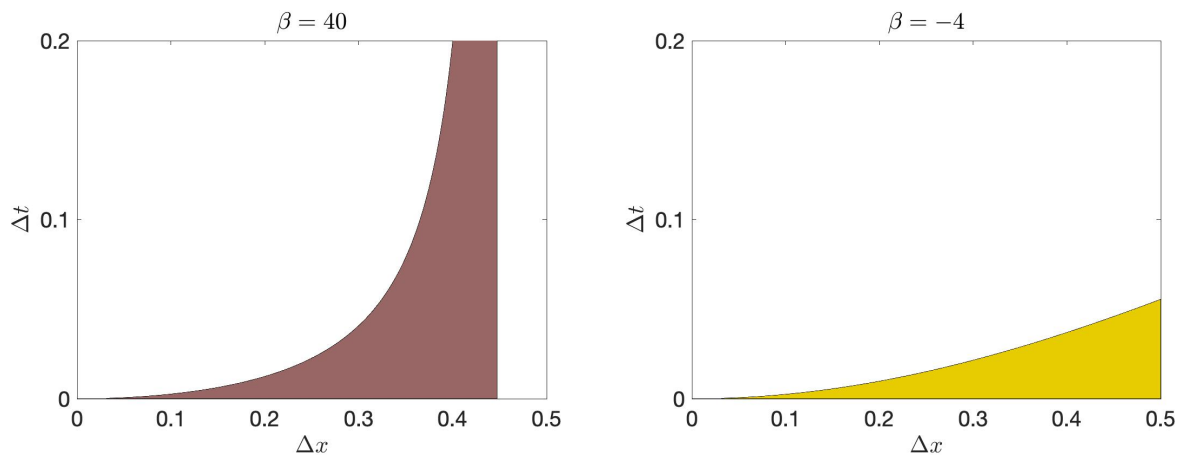


Figure 1: Lax-Richtmyer stability regions for the forward-in-time (Euler) centered-in-space (second-order) discretization of the heat equation (9) with $\alpha = 2$ and zero Dirichlet boundary conditions. The regions of stability are computed numerically using (40). Note the vertical asymptote for $\beta = 40$ at $\Delta x = 2\sqrt{\alpha/\beta} = 0.4472$ (see Eq. (43))

The first one can be written as

$$\frac{\Delta t}{\Delta x^2} \leq \frac{2}{4\alpha - \beta\Delta x^2}. \quad (43)$$

For $\beta > 0$ this yields the additional condition $\Delta x \leq 2\sqrt{\alpha/\beta}$ (see Figure 20). On the other hand, for $\beta = 0$ everything simplifies substantially. In particular, the second inequality in (42) yields the trivial condition $\alpha \geq 0$, while the first inequality yields

$$\frac{\Delta t}{\Delta x^2} \leq \frac{1}{2\alpha}. \quad (44)$$

The condition (44) also 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 (44) (or (43)) to hold in the limits $\Delta t, \Delta x \rightarrow 0$.

Lax-Richtmyer stability analysis applies 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 (14)). This makes the stability analysis of implicit schemes not straightforward nor practical using the matrix \mathbf{B} . We will see hereafter that this issue can be mitigated (at least for linear PDEs) by using discrete Fourier series.

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

$$\sum_{k=0}^q \alpha_k u_j^{k+q} = \Delta t \Phi_j \left(\mathbf{u}^{k+q}, \dots, \mathbf{u}^k, \Delta t, \Delta x \right). \quad (45)$$

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

$$u_j^{k+1} - u_j^k = \Delta t \left(-u_j^k \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} - \frac{u_{j-1}^k - 2u_j^k + u_{j+1}^k}{\Delta x^2} - \frac{u_{j-2}^k - 3u_{j-1}^k + 6u_j^k - 3u_{j+1}^k + u_{j+2}^k}{\Delta x^4} \right), \quad (46)$$

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 course note 4. To use such a proof in the context of finite-difference approximations of PDEs, we need to make sure that the Lipschitz constant of Φ in (45) 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 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 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. One reason for the Fourier series analysis is that it allows us to determine stability conditions for both implicit and explicit schemes in a rather straightforward way. To illustrate the method, let us consider the prototype IBVP

$$\begin{cases} \frac{\partial U(x, t)}{\partial t} = \alpha \frac{\partial^2 U}{\partial x^2} & t \geq 0 \quad x \in [0, L] \\ U(x, 0) = U_0(x) \\ \text{Periodic B.C.} \end{cases} \quad (47)$$

We have seen that (47) can be discretized with second-order finite-differences in space and Euler-forward time integration as

$$u_j^{k+1} = u_j^k + \frac{\alpha \Delta t}{\Delta x^2} (u_{j-1}^k - 2u_j^k + u_{j+1}^k) \quad j = 0, \dots, N-1, \quad (48)$$

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

$$x_j = j \frac{L}{N} \quad j = 0, \dots, N. \quad (49)$$

The scheme (48) is supplemented with periodic boundary conditions

$$u_j^k = u_{j+N}^k \quad \text{for all } j \in \mathbb{Z}. \quad (50)$$

and the initial condition

$$u_j^0 = U_0(x_j) \quad j = 0, \dots, N-1. \quad (51)$$

Let us now expand the numerical solution u_j^k in a discrete Fourier series⁶

⁶As is well known, the solution to (47) can be expanded in a Fourier series as

$$U(x, t) = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k x / L} \quad (52)$$

Evaluating $U(x, t)$ on the grid (49) yields the discrete Fourier series

$$U(x_j, t) = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k x_j / L} = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k j / N} \quad j = 0, \dots, N-1. \quad (53)$$

Moreover,

$$U(x_j, t) = \sum_{k=-\infty}^{\infty} C_k(t) e^{2\pi i k j / N} = \sum_{k=0}^{N-1} \sum_{p=-\infty}^{\infty} C_{k+pN}(t) e^{2\pi i (k+pN) j / N} = \sum_{k=0}^{N-1} e^{2\pi i k j / N} \underbrace{\sum_{p=-\infty}^{\infty} C_{k+pN}(t)}_{c_k(t)/N}. \quad (54)$$

This can be written as

$$U(x_j, t) = \frac{1}{N} \sum_{k=0}^{N-1} c_k(t) e^{2\pi i k j / N} = \frac{\Delta x}{L} \sum_{k=0}^{N-1} c_k(t) e^{i k j \xi}, \quad \xi = \frac{2\pi \Delta x}{L}. \quad (55)$$

$$u_j^k = \frac{1}{N} \sum_{p=0}^{N-1} c_p^k e^{ipj\xi}, \quad \text{where} \quad \xi = \frac{2\pi\Delta x}{L}, \quad (56)$$

and substitute it into (48) to obtain

$$\begin{aligned} \sum_{p=0}^{N-1} c_p^{k+1} e^{ipj\xi} &= \sum_{p=0}^{N-1} c_p^k e^{ipj\xi} \left[1 + \frac{\alpha\Delta t}{\Delta x^2} (e^{-ip\xi} - 2 + e^{ip\xi}) \right] \\ &= \sum_{p=0}^{N-1} c_p^k e^{ipj\xi} \left[1 + \frac{\alpha\Delta t}{\Delta x^2} (2 \cos(p\xi) - 2) \right], \end{aligned} \quad (57)$$

i.e.,

$$c_p^{k+1} = c_p^k \underbrace{\left[1 + \frac{2\alpha\Delta t}{\Delta x^2} \left(\cos\left(2\pi p \frac{\Delta x}{L}\right) - 1 \right) \right]}_{\text{amplification factor } G_p(\Delta t, \Delta x)}. \quad (58)$$

Upon definition of

$$\mathbf{c}^k = \begin{bmatrix} c_0^k \\ c_1^k \\ \vdots \\ c_{N-1}^k \end{bmatrix}, \quad \mathbf{G}(\Delta t, \Delta x) = \begin{bmatrix} G_0(\Delta t, \Delta x) & 0 & \cdots & 0 \\ 0 & G_1(\Delta t, \Delta x) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & G_{N-1}(\Delta t, \Delta x) \end{bmatrix} \quad (59)$$

we can write (58) as

$$\mathbf{c}^{k+1} = \mathbf{G}(\Delta t, \Delta x) \mathbf{c}^k. \quad (60)$$

The matrix $\mathbf{G}(\Delta t, \Delta x)$ in (60) plays the same role in Fourier space as the matrix \mathbf{B} in (8) does in physical space. In other words, for the scheme (48) to be stable we must have

$$\|\mathbf{G}^k\| \leq H_T \quad \text{for all } k \text{ such that } k\Delta t \leq T. \quad (61)$$

where H_T is a constant that does not depend on Δx or on Δt . In (61) $\|\cdot\|$ denotes any matrix norm induced by a vector norm.

Remark: Clearly, if we compute the inverse Fourier transform of (60) we obtain

$$\mathbf{u}^{k+1} = \mathbf{F} \mathbf{G} \mathbf{F}^{-1} \mathbf{u}^k, \quad (62)$$

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

$$\mathbf{u}^k = \mathbf{F} \mathbf{c}^k, \quad \mathbf{F} = \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & e^{i\xi} & e^{2i\xi} & \cdots & e^{i(N-1)\xi} \\ 1 & e^{2i\xi} & e^{4i\xi} & \cdots & e^{2i(N-1)\xi} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & e^{(N-1)i\xi} & e^{2(N-1)\xi} & \cdots & e^{i(N-1)^2\xi} \end{bmatrix} \quad (63)$$

A comparison between (62) and (8) shows that

$$\mathbf{B} = \mathbf{F} \mathbf{G} \mathbf{F}^{-1}. \quad (64)$$

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

$$\|\mathbf{B}^k\| = \|\mathbf{F}\mathbf{G}^k\mathbf{F}^{-1}\|. \quad (65)$$

The Fourier transform matrix \mathbf{F} plays no role in the stability properties of the scheme.

Necessary and sufficient conditions for Von-Neumann stability. Let us recall that the spectral radius of a matrix \mathbf{G} , i.e.,

$$\rho(\mathbf{G}) = \max_i |\lambda_i| \quad (66)$$

where λ_i are the eigenvalues of \mathbf{G} , is a lower bound for any (sub-multiplicative) matrix norm⁷ of \mathbf{G} , i.e.,

$$\rho(\mathbf{G}) \leq \|\mathbf{G}\| \quad \text{for every sub-multiplicative matrix norm } \|\cdot\|. \quad (69)$$

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 (61) we obtain

$$\rho(\mathbf{G})^k \leq \|\mathbf{G}^k\| \leq H_T \quad (70)$$

i.e.,

$$\rho(\mathbf{G})^k \leq H_T. \quad (71)$$

As before, this implies that for sufficiently small Δt the spectral radius of \mathbf{G} must satisfy (see Eq. (31))

$$\rho(\mathbf{G}) \leq 1 + \gamma\Delta t \quad (72)$$

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

$$\rho(\mathbf{G})^k \leq H_T \quad \text{does not imply} \quad \|\mathbf{G}^k\| \leq H_T. \quad (73)$$

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 (72) is sufficient.

Lemma 1. The Von-Neumann stability condition (72) is sufficient if the amplification matrix \mathbf{G} is normal.

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

$$\rho(\mathbf{G}) = \sqrt{\rho(\mathbf{G}\mathbf{G}^*)} = \|\mathbf{G}\|_2. \quad (74)$$

This allows us to write (70) as

$$\rho(\mathbf{G})^k = \|\mathbf{G}^k\|_2 \leq H_T. \quad (75)$$

Hence, for normal matrices \mathbf{G} we have that (72) implies

$$\|\mathbf{G}^k\|_2 \leq H_T \quad (76)$$

⁷A sub-multiplicative matrix norm is a norm satisfying

$$\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad (67)$$

for all matrices \mathbf{A} and \mathbf{B} . All matrix norms induced by vector norms are sub-multiplicative. To prove (69) it is sufficient to 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

$$\|\mathbf{G}\mathbf{V}\| = |\lambda_i| \|\mathbf{V}\| \leq \|\mathbf{G}\| \|\mathbf{V}\| \quad \Rightarrow \quad \|\mathbf{G}\| \geq \max_i |\lambda_i| = \rho(\mathbf{G}). \quad (68)$$

the matrix

i.e., that the scheme is stable. Recall that stability in one norm implies stability in any other norm. \square

Fast computation of the amplification factors. 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 one wave number. In practice, we can simply substitute

$$u_j^k = c_p^k e^{ijp\xi} \quad \text{where} \quad \xi = \frac{2\pi\Delta x}{L} \quad (77)$$

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.

- **Von-Neumann stability analysis of the heat equation (Euler-forward time integration).** We have seen that the Fourier transform of finite-difference scheme (48) yields the diagonal matrix of amplification factors defined in (59). The diagonal entries of \mathbf{G} are the eigenvalues of \mathbf{G} . Hence, the spectral radius of \mathbf{G} is

$$\rho(\mathbf{G}) = \max_{p=0,\dots,N-1} \left| 1 + \frac{2\alpha\Delta t}{\Delta x^2} \left(\cos\left(2\pi p \frac{\Delta x}{L}\right) - 1 \right) \right| \simeq \left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \quad (\text{for large } N). \quad (78)$$

By using the Von-Neumann condition (72) we conclude that the scheme (48) is stable if and only if

$$\left| \frac{4\alpha\Delta t}{\Delta x^2} - 1 \right| \leq 1 + \gamma\Delta t. \quad (79)$$

This is exactly the same condition we obtained in (40) (see the discussion thereafter).

- **Von-Neumann stability analysis of the heat equation (Crank-Nicolson time integration).** Consider the fully discrete finite-difference scheme

$$u_j^{k+1} - \frac{\alpha\Delta t}{2\Delta x^2} (u_{j-1}^{k+1} - 2u_j^{k+1} + u_{j+1}^{k+1}) = u_j^k + \frac{\alpha\Delta t}{2\Delta x^2} (u_{j-1}^k - 2u_j^k + u_{j+1}^k). \quad (80)$$

Substitute (79) into (80) to obtain

$$c_p^{k+1} \left[1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right] = c_p^k \left[1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right], \quad (81)$$

i.e.,

$$c_p^{k+1} = \frac{1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)}{1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)} c_p^k. \quad (82)$$

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

$$\rho(\mathbf{G}) = \max_{p=0,\dots,N-1} \left| \frac{1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)}{1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1)} \right|. \quad (83)$$

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

$$\left| 1 + \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right| \leq \left| 1 - \frac{\alpha\Delta t}{\Delta x^2} (\cos(p\xi) - 1) \right| \quad (84)$$

and

$$\rho(\mathbf{G}) \leq 1. \quad (85)$$

Recalling that the Von-Neumann stability condition (72) 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.

Clearly, by following the same steps that lead us to (31) we see that (61) yields the following *sufficient condition* for stability

$$\|\mathbf{G}\| \leq 1 + \delta\Delta t \quad \text{as} \quad \Delta t \rightarrow 0, \quad (86)$$

where $\|\mathbf{G}\|$ denotes any matrix norm compatible with a vector norm.

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.