

## HW04

### 2.3.1.

We consider the  $\theta$ -scheme written in operator form

$$(I - \theta k D_0) v^{n+1} = (I + (1 - \theta) k D_0) v^n,$$

where  $k = \Delta t$  and  $D_0$  is the usual second-difference operator (discrete Laplacian) in space. Write the amplification operator

$$G = (I - \theta k D_0)^{-1} (I + (1 - \theta) k D_0),$$

so that  $v^{n+1} = G v^n$ . To prove unconditional stability for  $\theta \geq \frac{1}{2}$  it suffices to show the spectral radius  $\rho(G) \leq 1$  (or equivalently that every Fourier mode or eigenmode is not amplified in modulus).

Assume periodic boundary conditions so that  $D_0$  is diagonalizable with an orthonormal basis of eigenvectors and real eigenvalues. For the standard centered second-difference operator  $D_0$  the eigenvalues are real and nonpositive; denote any eigenvalue of  $k D_0$  by  $\lambda \leq 0$ . On the corresponding eigenvector the scalar amplification factor is

$$g(\lambda) = \frac{1 + (1 - \theta)\lambda}{1 - \theta\lambda}.$$

Put  $\lambda = -\alpha$  with  $\alpha \geq 0$ . Then

$$g(\lambda) = \frac{1 - (1 - \theta)\alpha}{1 + \theta\alpha}.$$

We need to show  $|g(\lambda)| \leq 1$  for all  $\alpha \geq 0$  when  $\theta \geq \frac{1}{2}$ .

Observe first that for any  $\alpha \geq 0$ ,

$$|1 - (1 - \theta)\alpha| \leq 1 + (1 - \theta)\alpha,$$

and when  $\theta \geq \frac{1}{2}$  we have  $1 - \theta \leq \theta$ , hence

$$1 + (1 - \theta)\alpha \leq 1 + \theta\alpha.$$

Combining these inequalities gives

$$|1 - (1 - \theta)\alpha| \leq 1 + \theta\alpha.$$

Therefore

$$|g(\lambda)| = \frac{|1 - (1 - \theta)\alpha|}{1 + \theta\alpha} \leq \frac{1 + \theta\alpha}{1 + \theta\alpha} = 1.$$

Since this holds for every eigenvalue  $\lambda$  of  $k D_0$ , the spectral radius  $\rho(G) \leq 1$ . Thus the  $\theta$ -scheme is unconditionally stable for  $\theta \geq \frac{1}{2}$ .  $\square$

### 2.4.1.

Let  $u(x, t)$  be sufficiently smooth in a neighbourhood of a fixed point  $(x_*, t_*)$ . Write  $\xi = x - x_*$ ,  $\tau = t - t_*$ . The two-variable Taylor expansion about  $(x_*, t_*)$  up to second order is

$$\begin{aligned} u(x_* + \xi, t_* + \tau) &= u(x_*, t_*) + u_x(x_*, t_*)\xi + u_t(x_*, t_*)\tau \\ &\quad + \frac{1}{2}u_{xx}(x_*, t_*)\xi^2 + u_{xt}(x_*, t_*)\xi\tau + \frac{1}{2}u_{tt}(x_*, t_*)\tau^2 + R(\xi, \tau), \end{aligned}$$

with the remainder  $R(\xi, \tau) = O(|\xi|^3 + |\xi|^2|\tau| + |\xi||\tau|^2 + |\tau|^3)$ .

When deriving the order of accuracy of a finite difference formula we substitute the exact solution values at the stencil points into the difference formula and compare with the continuous operator. If the stencil points are at offsets  $\xi_i = O(\Delta x)$ ,  $\tau_j = O(\Delta t)$  relative to  $(x_*, t_*)$ , then each term in the Taylor expansion can be grouped by powers of  $\Delta x, \Delta t$ ; the remainder becomes  $O(\Delta x^p, \Delta t^q)$  accordingly. Hence the leading-order truncation terms (and therefore the formal order) depend only on the relative offsets  $\xi_i, \tau_j$  (i.e. their orders in  $\Delta x, \Delta t$ ) and on the derivatives of  $u$  at  $(x_*, t_*)$ , not on whether  $(x_*, t_*)$  coincides with a grid point. Therefore one may choose the expansion centre  $(x_*, t_*)$  arbitrarily (provided  $u$  is smooth there); in particular it need not be a grid point.  $\square$

#### 2.4.2.

We specialise to the one-dimensional linear advection equation

$$u_t + au_x = 0, \quad a \in \mathbb{R},$$

and show that both the leap-frog and the Crank-Nicolson (CN) schemes have second order accuracy in time and second order in space, i.e. they are (2, 2) schemes. Then we compare their numerical behaviour.

**(i) Leap-frog.** The standard leap-frog discretisation for the advection equation (space by centred difference, time by centred difference) is

$$\frac{u_j^{n+1} - u_j^{n-1}}{2\Delta t} + a \frac{u_{j+1}^n - u_{j-1}^n}{2\Delta x} = 0.$$

Substitute the exact solution  $u(x_j, t_n)$  and expand in Taylor series about  $(x_j, t_n)$ .

Time expansion:

$$u(x_j, t_{n\pm 1}) = u + (\pm \Delta t)u_t + \frac{(\Delta t)^2}{2}u_{tt} \pm \frac{(\Delta t)^3}{6}u_{ttt} + O(\Delta t^4),$$

hence

$$\frac{u(x_j, t_{n+1}) - u(x_j, t_{n-1})}{2\Delta t} = u_t + \frac{(\Delta t)^2}{6}u_{ttt} + O(\Delta t^4).$$

Space expansions (centre-difference for  $u_x$ ):

$$\frac{u(x_{j+1}, t_n) - u(x_{j-1}, t_n)}{2\Delta x} = u_x + \frac{(\Delta x)^2}{6}u_{xxx} + O(\Delta x^4).$$

Substituting into the scheme and using the PDE  $u_t = -au_x$  to cancel leading terms gives the local truncation error

$$\text{LTE}_{\text{LF}} = \frac{(\Delta t)^2}{6}u_{ttt} + a \cdot \frac{(\Delta x)^2}{6}u_{xxx} + O(\Delta t^4, \Delta x^4).$$

Thus the leading time error is  $O(\Delta t^2)$  and the leading space error is  $O(\Delta x^2)$ . Therefore the leap-frog scheme is (2, 2).

**(ii) Crank-Nicolson.** Although CN is usually used for diffusion problems, one may form a CN-type centred-in-time scheme for advection by applying the trapezoidal rule in time to the spatial difference:

$$\frac{u_j^{n+1} - u_j^n}{\Delta t} + a \frac{1}{2} \left( \frac{u_{j+1}^{n+1} - u_{j-1}^{n+1}}{2\Delta x} + \frac{u_{j+1}^n - u_{j-1}^n}{2\Delta x} \right) = 0.$$

Expand  $u(x_j, t_{n+1})$  in time about  $t_n$ :

$$u(x_j, t_{n+1}) = u + \Delta t u_t + \frac{(\Delta t)^2}{2}u_{tt} + \frac{(\Delta t)^3}{6}u_{ttt} + O(\Delta t^4),$$

and similarly expand the spatial centred differences. Using the PDE to eliminate  $u_t$ , the trapezoidal (midpoint) time discretisation yields a time truncation error  $O(\Delta t^2)$ . With centred spatial differences giving  $O(\Delta x^2)$  spatial truncation, the CN scheme is therefore (2, 2).

**(iii) Why two (2, 2) schemes can behave differently.** Although both schemes have the same formal orders, they differ in:

1. **Truncation error constants.** The leading coefficients multiplying  $\Delta t^2$  and  $\Delta x^2$  differ; a smaller constant yields a smaller error for the same grid.
2. **Dissipation and dispersion.** Leap–frog is non-dissipative (no amplitude damping for Fourier modes) but dispersive (phase errors). CN (as a trapezoidal method) typically damps high-frequency modes when combined with discrete spatial dissipation terms, so it behaves differently on under-resolved scales.
3. **Computational modes and filtering.** Leap–frog supports a computational (parasitic) mode that alternates in time (the scheme couples even and odd time levels); in practice one often needs a filter (or start-up step) to suppress this. CN is implicit and does not produce the same alternating parasitic mode.
4. **Stability domain.** The von–Neumann stability (and thus practical allowable time step) differs — this affects which scheme is more accurate in practice for a given  $\Delta t, \Delta x$ .

Consequently, with the same nominal (2, 2) order, one scheme can produce noticeably smaller errors or more robust behaviour than the other depending on the problem and the chosen mesh parameters.  $\square$

### Supplement 1.

We construct a conservative finite–volume style scheme of temporal order 1 and spatial order 3 for

$$u_t + au_x = 0, \quad a = \text{const.}$$

**Derivation (integral form + quadratic reconstruction).** Integrate the conservation law over the cell  $I_j = [x_{j-1/2}, x_{j+1/2}]$  and over one time step  $[t^n, t^{n+1}]$ :

$$\bar{u}_j^{n+1} = \bar{u}_j^n - \frac{\Delta t}{\Delta x} \left( F_{j+1/2}^n - F_{j-1/2}^n \right),$$

where  $\bar{u}_j^n$  denotes the cell average of  $u$  on  $I_j$  at  $t^n$ , and the time–averaged numerical flux

$$F_{j+1/2}^n \approx \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} a u(x_{j+1/2}, t) dt.$$

Using forward Euler in time (first-order), replace the time averaged flux by the flux at the left endpoint  $t^n$ :

$$F_{j+1/2}^n \approx a u(x_{j+1/2}, t^n) \quad (\text{first-order in time}).$$

To obtain third-order accuracy in space we reconstruct a quadratic interpolant through three neighbouring pointwise values. Using point values  $u_{j-1}^n, u_j^n, u_{j+1}^n$  (assume the solution is smooth so point values and cell averages differ by higher-order terms), the Lagrange quadratic interpolant evaluated at the right face  $x_{j+1/2}$  gives

$$u_{j+1/2}^- = -\frac{1}{8}u_{j-1}^n + \frac{3}{4}u_j^n + \frac{3}{8}u_{j+1}^n,$$

which is third-order accurate as an approximation of the face value from the left. (This value follows from Lagrange interpolation through the three points  $x_{j-1}, x_j, x_{j+1}$  and evaluating at  $x_{j+1/2}$ .)

Using upwinding for the flux: if  $a > 0$  then

$$F_{j+1/2}^n = a u_{j+1/2}^-,$$

if  $a < 0$  use the right-biased reconstruction  $u_{j+1/2}^+$  analogously.

Therefore the scheme is

$$\bar{u}_j^{n+1} = \bar{u}_j^n - \frac{\Delta t}{\Delta x} a \left( u_{j+1/2}^- - u_{j-1/2}^- \right), \quad u_{j+1/2}^- = -\frac{1}{8}u_{j-1}^n + \frac{3}{4}u_j^n + \frac{3}{8}u_{j+1}^n.$$

This scheme is first order in time (forward Euler) and third order in space (quadratic reconstruction), thus it is a (1, 3) scheme. If higher temporal accuracy is needed one may replace forward Euler by a higher-order ODE integrator (e.g. strong-stability-preserving Runge–Kutta) preserving the spatial third-order accuracy.

**Remark.** The expression for  $u_{j+1/2}^-$  above is obtained by writing the quadratic interpolant  $p(x)$  through the three points  $(x_{j-1}, u_{j-1}), (x_j, u_j), (x_{j+1}, u_{j+1})$  and evaluating  $p(x_{j+1/2})$ . One can also derive equivalent third-order upwind-biased finite-difference formulas for the derivative  $u_x$  and construct a method-of-lines (1, 3) ODE discretisation with forward Euler in time.

## Supplement 2.

We again consider  $u_t + au_x = 0$ . Below we derive (a) the PDE domain of dependence for the exact solution, (b) the discrete domain of dependence of the leap–frog numerical solution, and (c) the CFL condition that ensures the continuous domain is contained in the discrete one (consistency with causality) and coincides with the von–Neumann stability restriction.

**(a) Continuous domain of dependence.** The characteristic equations are  $\frac{dx}{dt} = a, \frac{du}{dt} = 0$ . A characteristic through  $(x, t)$  reaches the initial line  $t = 0$  at

$$x_0 = x - at.$$

Hence the exact solution at  $(x, t)$  depends only on the initial datum at  $x_0$ ; in other words the continuous domain of dependence for the point  $(x, t)$  is the single point  $x_0$  on the initial line (or, if initial data were given on an interval, the intersection of the characteristic with that interval). For finite-support initial data the continuous physical domain of dependence at time  $t$  is the characteristic line segment  $x - as$  for  $0 \leq s \leq t$ .

**(b) Discrete domain of dependence for leap–frog.** The leap–frog update uses the values at the two previous time levels and the two neighbouring spatial points. Written again,

$$u_j^{n+1} = u_j^{n-1} - \mu (u_{j+1}^n - u_{j-1}^n), \quad \mu = \frac{a\Delta t}{2\Delta x}.$$

From this stencil one sees that  $u_j^{n+1}$  depends directly on  $u_{j-1}^n, u_{j+1}^n$  and on  $u_j^{n-1}$ . By induction on  $n$ , any grid value  $u_j^n$  depends only on initial grid values  $u_k^0$  with indices  $k$  satisfying  $|k - j| \leq n$  (i.e. at most one grid index per time level can be added/subtracted). Thus the discrete domain of dependence after  $n$  time steps is the index set  $\{j - n, \dots, j + n\}$ . In physical coordinates this corresponds to the interval

$$[x_j - n\Delta x, x_j + n\Delta x] = [x_j - \frac{t}{\Delta t} \Delta x, x_j + \frac{t}{\Delta t} \Delta x].$$

**(c) CFL condition.** For the numerical scheme to be consistent with the PDE causality, the continuous characteristic foot  $x_0 = x_j - at$  (or the continuous domain of dependence) should lie inside the discrete domain of dependence. After one time step ( $t = \Delta t$ ) this requirement becomes

$$|a|\Delta t \leq \Delta x,$$

or equivalently

$$|\mu_{\text{CFL}}| := \left| \frac{a\Delta t}{\Delta x} \right| \leq 1.$$

This is the usual CFL condition in this context: the physical information (characteristic displacement  $a\Delta t$ ) should not travel more than one grid cell per time step.

**Consistency with von-Neumann stability.**

Von-Neumann analysis of the leap-frog scheme (substituting  $u_j^n = \xi^n e^{ikx_j}$ ) yields the relation

$$\frac{\xi - \xi^{-1}}{2} + i\mu \sin \theta = 0, \quad \theta = k\Delta x, \quad \mu = \frac{a\Delta t}{\Delta x}.$$

This implies  $\xi = e^{\pm i\omega\Delta t}$  with

$$\sin(\omega\Delta t) = -\mu \sin \theta.$$

For real  $\omega$  (no spurious amplifying modes) one requires  $|\mu \sin \theta| \leq 1$  for all  $\theta$ . The maximum of  $|\sin \theta|$  is 1, so a sufficient (and sharp) condition is  $|\mu| \leq 1$ , i.e. the CFL condition above. Thus the causality requirement that continuous domain of dependence lie within the discrete one coincides with the von-Neumann stability restriction  $|a\Delta t/\Delta x| \leq 1$ .

**Conclusion.** The continuous characteristic for  $u_t + au_x = 0$  is  $x - x_0 = at$ ; the leap-frog discrete domain after  $n$  steps is  $\{j - n, \dots, j + n\}$  (physical width  $2n\Delta x$ ); consistency and von-Neumann stability require the CFL bound  $|a|\Delta t/\Delta x \leq 1$ .  $\square$

*If you want, I can:*

- replace the reconstruction in Supplement 1 by the explicit algebraic derivation of the interpolation coefficients (showing the Lagrange polynomials step by step),
- or replace the forward-Euler time discretisation by a specific SSP-RK(2 or 3) integrator to obtain a (2, 3) or (3, 3) scheme,
- or add a short numerical example (one-period sine wave) with error constants to illustrate the differences between leap-frog and CN.