

List of topics in this lecture

- 2D IBVP of the heat equation, implementation of 2D FTCS and 2D BTCS
- Crank-Nicolson method, alternating direction implicit method (ADI)
- Numerical solution of hyperbolic PDEs, upwind method, downwind method, FTCS method, Lax-Friedrichs method, local truncation errors and stability of these methods, characteristics propagation, boundary conditions in IBVP

Review

Lax equivalence theorem: Consistency + Stability = Convergence

Method of lines (MOL), numerical stability of ODE solvers vs PDE solvers

ODE solvers	PDE solvers
An ODE system may be stiff but the <i>stiffness is fixed</i> as numerical resolution is refined. As Δt is refined, eventually all methods will be well behaved for small Δt . We use L-stable ODE solvers to avoid tiny Δt .	After MOL (method of lines) discretization, a PDE becomes an ODE system. <i>The stiffness of the system increases</i> as numerical resolution is refined. As $(\Delta x, \Delta t)$ is refined, a constraint on $\Delta t/(\Delta x)^2$ is needed to ensure the numerical stability of <u>explicit methods</u> . We use L-stable ODE solvers to get out of this constraint.

Procedure of von Neumann stability analysis

- Substitute $u_i^n = \rho^n \exp(\sqrt{-1} \xi i \Delta x)$ into the numerical method to calculate ρ .
- Determine the stability as follows.

If $|\rho(\xi, \Delta t)| \leq 1 + C \Delta t$ for small Δt and all ξ is valid for r within a certain range, then the method is stable for r in that range.

If $|\rho(\xi, \Delta t)| \leq 1 + C \Delta t$ for small Δt and all ξ is valid for all r , then the method is unconditionally stable (stable for all r).

If $|\rho(\xi, \Delta t)| \leq 1 + C \Delta t$ for small Δt and all ξ is valid for none of r , then the method is unconditionally unstable (unstable for all r).

Two-dimensional IBVP of the heat equation

$$\begin{cases} u_t = u_{xx} + u_{yy}, & (x, y) \in \Omega \\ u(x, y, 0) = f(x, y), & (x, y) \in \Omega \\ u(x, y, t) = g(x, y, t), & (x, y) \in \partial\Omega \end{cases} \quad \Omega = [0, L_x] \times [0, L_y]$$

FTCS method:

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

$$u_{i,j}^0 = f(x_i, y_j)$$

$$u_{0,j}^n = g(0, y_j, t_n), \quad u_{(N_x+1),j}^n = g(L_x, y_j, t_n)$$

$$u_{i,0}^n = g(x_i, 0, t_n), \quad u_{i,(N_y+1)}^n = g(x_i, L_y, t_n)$$

Matlab matrix $\{u^n(j, i), 1 \leq j \leq N_y + 2, 1 \leq i \leq N_x + 2\}$ stores $\{u_{i-1,j-1}^n\}$, values of u at all points, including boundary points. In each time step, we do two things:

- Update $\{u^n(j, i)\}$ at boundary points
- Calculate $\{u^{n+1}(j, i)\}$ at internal points

Remark: We don't need to write $\{u^n(j, i)\}$ as a long vector.

BTCS method

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

$$u_{i,j}^0 = f(x_i, y_j)$$

$$u_{0,j}^{n+1} = g(0, y_j, t_{n+1}), \quad u_{(N_x+1),j}^{n+1} = g(L_x, y_j, t_{n+1})$$

$$u_{i,0}^{n+1} = g(x_i, 0, t_{n+1}), \quad u_{i,(N_y+1)}^{n+1} = g(x_i, L_y, t_{n+1})$$

Matlab matrix $\{u^{n+1}(j, i), 1 \leq j \leq N_y, 1 \leq i \leq N_x\}$ stores $\{u_{i,j}^{n+1}\}$, values of u at internal points only, excluding boundary points.

BTCS in matrix-vector form

$$u^{n+1} = u^n + \Delta t A_x u^{n+1} + \Delta t A_y u^{n+1} + \Delta t b^{n+1}$$

$$\Rightarrow (I - \Delta t A_x - \Delta t A_y) u^{n+1} = u^n + \Delta t b^{n+1}$$

$$\Rightarrow u^{n+1} = (I - \Delta t A_x - \Delta t A_y)^{-1} (u^n + \Delta t b^{n+1}) \quad \text{solution of a linear system}$$

Remark:

- We need to write $\{u^{n+1}(j,i)\}$ as a long vector in the linear system.
- We need to construct matrices A_x and A_y and vector b^{n+1} in the linear system.

Mapping between 2D array and 1D vector

In the linear system, we need to write $\{u^{n+1}(j,i)\}$ as a long vector. We list all internal points by scanning each column of $\{u^{n+1}(j,i)\}$, storing $\{u_{i,j}^{n+1}\}$.

We write out the mapping between the Matlab 2D index and 1D index.

$$\underbrace{(i, j)}_{\text{2D index}} \xleftrightarrow{\text{Mapping}} \underbrace{k}_{\text{1D index}}$$

$$(i, j) \rightarrow k: \quad k(i, j) = (i-1)N_y + j$$

$$k \rightarrow (i, j): \quad \begin{cases} i(k) = \text{floor}((k-1)/N_y) + 1 \\ j(k) = (k-1) - N_y \text{floor}((k-1)/N_y) + 1 \end{cases}$$

The 1D Matlab vector u is

$$\{u_{1D}^n(k) = u_{2D}^n(j(k), i(k)), \quad k = 1, 2, \dots, N_{xy}\}, \quad N_{xy} \equiv N_x N_y$$

In Matlab, `u1D = reshape(u2D, Nx*Ny, 1);`

Constructing matrix A_x corresponding to $D_x^2 \{u_{i,j}\} = \frac{1}{(\Delta x)^2} [u_{i+1,j} - 2u_{i,j} + u_{i-1,j}]$

Row k of A_x has at most 3 non-zero entries.

- map k to $(i(k), j(k))$
- two neighboring points in finite difference are $(i(k)-1, j(k))$, $(i(k)+1, j(k))$
- if $(i(k)-1, j(k))$ is still an internal point, map it to $k_{\text{Left}} \equiv k(i(k)-1, j(k))$
- if $(i(k)+1, j(k))$ is still an internal point, map it to $k_{\text{Right}} \equiv k(i(k)+1, j(k))$
- $A_x(k, k) = -2/(\Delta x)^2$, $A_x(k, k_{\text{Left}}) = 1/(\Delta x)^2$, $A_x(k, k_{\text{Right}}) = 1/(\Delta x)^2$

For 2D problems, matrix A_x needs to be stored as a sparse matrix.

Example:

$$N_x = N_y = 400 \quad \Rightarrow \quad N_{xy} \equiv N_x N_y = 160000$$

a 160000×160000 matrix is too big to store as a dense matrix.

Constructing matrix A_y corresponding to $D_y^2 \{u_{i,j}\} = \frac{1}{(\Delta x)^2} [u_{i,j+1} - 2u_{i,j} + u_{i,j-1}]$

Row k of A_y has at most 3 non-zero entries.

- map k to $(i(k), j(k))$
- two neighboring points in finite difference are $(i(k), j(k)-1)$, $(i(k), j(k)+1)$
- if $(i(k), j(k)-1)$ is still an internal point, map it to $k_{\text{Down}} \equiv k(i(k), j(k)-1)$
- if $(i(k), j(k)+1)$ is still an internal point, map it to $k_{\text{Up}} \equiv k(i(k), j(k)+1)$
- $A_y(k, k) = -2/(\Delta x)^2$, $A_y(k, k_{\text{Down}}) = 1/(\Delta x)^2$, $A_y(k, k_{\text{Up}}) = 1/(\Delta x)^2$

Calculating vector b^{n+1} containing the effect of boundary conditions

Start with matrix version $\{b_{i,j}\} = 0$, stored in Matlab matrix $\{b(j, i)\} = 0$.

Left boundary affects $\{b_{1,j}\}$, stored in $\{b(j, 1)\}$

$$b(j, 1) = b(j, 1) + g_L(j\Delta y, (n+1)\Delta t)/(\Delta x)^2$$

Right boundary affects $\{b_{N_x,j}\}$, stored in $\{b(j, N_x)\}$

$$b(j, N_x) = b(j, N_x) + g_R(j\Delta y, (n+1)\Delta t)/(\Delta x)^2$$

Bottom boundary affects $\{b_{i,1}\}$, stored in $\{b(1, i)\}$

$$b(1, i) = b(1, i) + g_B(i\Delta x, (n+1)\Delta t)/(\Delta y)^2$$

Top boundary affects $\{b_{i,N_y}\}$, stored in $\{b(j, N_x)\}$

$$b(N_y, i) = b(N_y, i) + g_T(i\Delta x, (n+1)\Delta t)/(\Delta y)^2$$

Map 2D array $\{b(j, i)\}$ to 1D vector $\{b(k)\}$

$$\left\{ b_{1D}^{n+1}(k) = u_{2D}(j(k), i(k)), k = 1, 2, \dots, N_{xy} \right\}, \quad N_{xy} \equiv N_x N_y$$

In Matlab, $b1D = \text{reshape}(b2D, N_x*N_y, 1);$

Calculating vector u^{n+1} containing u at internal points at time t_{n+1} .

Once we obtain matrices A_x and A_y and vector b^{n+1} , we calculate vector u^{n+1}

$$u^{n+1} = (I - \Delta t A_x - \Delta t A_y) \backslash (u^n + \Delta t b^{n+1})$$

Then the process is repeated in each time step.

Remarks:

- Matrices A_x and A_y only need to be calculated once and then are used in all time steps. Vector b^{n+1} needs to be calculated in each time step.
- We keep evolving 1D vector u forward in time until we need to output its 2D array version. The 2D array u^n gives a better description of function $u(x, y, t_n)$.

$$\left\{ u_{2D}^n(j, i) = u_{1D}^n((i-1)N_y + j), \quad j = 1, 2, \dots, N_y, \quad i = 1, 2, \dots, N_x \right\}$$

In Matlab, `u2D = reshape(u1D, Ny, Nx);`

Next study the Crank-Nicolson method for solving the 2D heat equation, which has the second order accuracy both in space and in time.

Crank-Nicolson method

$$u^{n+1} = u^n + \Delta t A_x \left(\frac{u^{n+1} + u^n}{2} \right) + \Delta t A_y \left(\frac{u^{n+1} + u^n}{2} \right) + \Delta t \left(\frac{b^{n+1} + b^n}{2} \right)$$

To calculate u^{n+1} from u^n , we need to solve a sparse linear system

$$\underbrace{\left(I - \frac{\Delta t}{2} A_x - \frac{\Delta t}{2} A_y \right)}_{\equiv B} u^{n+1} = \underbrace{\left(I + \frac{\Delta t}{2} A_x + \frac{\Delta t}{2} A_y \right)}_{\equiv B_2} u^n + \underbrace{\Delta t \left(\frac{b^{n+1} + b^n}{2} \right)}_{\equiv \tilde{b}^{n+1/2}}$$

$$\Rightarrow u^{n+1} = B \setminus (B_2 u^n + \Delta t \tilde{b}^{n+1/2})$$

On a numerical grid of 400×400 , matrix B is 160000×160000 , a huge matrix.

In the Crank-Nicolson method, the linear system is sparse. Nevertheless, it is still a huge linear system and its sparse structure is not easy to accommodate (as we saw in the implementation of BTCS). We like to replace the huge linear system with a collection of small tridiagonal linear systems. We now discuss such an approach.

Alternating direction implicit method (ADI)

Strategy:

- divide each time step into two half steps,
- in each half step, use implicit difference along only one direction (x or y), and
- alternate the implicit direction in two half steps.

$$\begin{cases} u^{n+1/2} = u^n + \underbrace{\frac{\Delta t}{2} A_x u^{n+1/2}}_{\text{Implicit}} + \underbrace{\frac{\Delta t}{2} A_y u^n}_{\text{Explicit}} + \frac{\Delta t}{2} (b_x^{n+1/2} + b_y^n) \\ u^{n+1} = u^{n+1/2} + \underbrace{\frac{\Delta t}{2} A_x u^{n+1/2}}_{\text{Explicit}} + \underbrace{\frac{\Delta t}{2} A_y u^{n+1}}_{\text{Implicit}} + \frac{\Delta t}{2} (b_x^{n+1/2} + b_y^{n+1}) \end{cases}$$

where

A_x : matrix corresponding to difference operator D_x^2

A_y : matrix corresponding to difference operator D_y^2

$b_x^{n+1/2}$: vector containing the effects of boundary conditions in $D_x^2 u^{n+1/2}$,

b_y^{n+1} : vector containing the effects of boundary conditions in $D_y^2 u^{n+1}$,

The first half step $u^n \rightarrow u^{n+1/2}$

It is 1st order in time; it is implicit in x only; it is explicit in y.

In this half step, we only need to solve a collection of small tridiagonal linear systems, each corresponding to the discretization of D_x^2 at a fixed y.

The second half step $u^{n+1/2} \rightarrow u^{n+1}$

It is 1st order in time; it is implicit in y only; it is explicit in x.

In this half step, we only need to solve a collection of small tridiagonal linear systems, each corresponding to the discretization of D_y^2 at a fixed x.

The combined step $u^n \rightarrow u^{n+1}$

It is 2nd order in time and 2nd order in space!

This can be seen by writing out the combined step.

$$u^{n+1} = u^n + \Delta t A_x u^{n+1/2} + \Delta t A_y \frac{u^{n+1} + u^n}{2} + \Delta t \left(b_x^{n+1/2} + \frac{b_y^n + b_y^{n+1}}{2} \right)$$

Note: every term on the RHS is evaluated either directly or effectively at $t_{n+1/2}$.

Stability of numerical methods for solving the 2D heat equation

FTCS method:

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

It is stable if and only if $\frac{\Delta t}{(\Delta x)^2} + \frac{\Delta t}{(\Delta y)^2} \leq \frac{1}{2}$.

Tool: 2D von Neumann stability analysis (see Appendix A).

BTCS method:

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

It is unconditionally stable.

Alternating direction implicit method (ADI)

$$\begin{cases} u^{n+1/2} = u^n + \underbrace{\frac{\Delta t}{2} D_x^2 u^{n+1/2}}_{\text{Implicit}} + \underbrace{\frac{\Delta t}{2} D_y^2 u^n}_{\text{Explicit}} \\ u^{n+1} = u^{n+1/2} + \underbrace{\frac{\Delta t}{2} D_y^2 u^{n+1/2}}_{\text{Explicit}} + \underbrace{\frac{\Delta t}{2} D_x^2 u^{n+1}}_{\text{Implicit}} \end{cases}$$

It is unconditionally stable (see Appendix B).

Numerical solution of hyperbolic PDEs

We first review a simple initial value problem.

IVP of $u_t + au_x = 0$

$$\begin{cases} u_t + au_x = 0, & a = \text{const} \\ u(x, 0) = u_0(x) \end{cases}$$

The initial condition in parametric form

$$\begin{cases} x_0(s) = s \\ t_0(s) = 0 \\ u(x_0(s), t_0(s)) = u_0(s) \end{cases} \quad s \text{ is a parameter}$$

Characteristics:

$$\begin{cases} \frac{dt}{d\tau} = 1, & t(0) = t_0(s) = 0 \\ \frac{dx}{d\tau} = a, & x(0) = x_0(s) = s \\ \frac{du}{d\tau} = 0, & u(0) = u_0(s) \end{cases} \quad \begin{array}{l} s \text{ is a parameter;} \\ \tau \text{ is the time along C-line.} \end{array}$$

Solution in parametric form, as a function of (s, τ) :

$$\begin{cases} t(s, \tau) = \tau \\ x(s, \tau) = s + a\tau \\ u(s, \tau) = u_0(s) \end{cases}$$

We like to express u as a function of (x, t) .

The inverse mapping of $(t(s, \tau), x(s, \tau))$:

$$\begin{cases} \tau = t \\ s = x - at \end{cases}$$

Solution u , as a function of (x, t) :

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

Propagation of solution:

$a > 0$, $u(x, t) = u_0(x - at)$ is a traveling wave to the right.

(Draw the characteristics and wave propagation in the x - t plane)

$a < 0$, $u(x, t) = u_0(x - at)$ is a traveling wave to the left.

(Draw the characteristics and wave propagation in the x - t plane)

In all real applications, we solve the PDE in a finite domain.

IBVP of $u_t + au_x = 0$

For $a > 0$ (wave propagating to the right)

$$\begin{cases} u_t + au_x = 0, & 0 < x < L \\ u(x, 0) = u_0(x), & 0 < x < L \\ u(0, t) = g(t) \end{cases}$$

Note: A boundary condition is imposed only at $x = 0$, not at $x = L$.

For $a < 0$ (wave propagating to the left)

$$\begin{cases} u_t + au_x = 0, & 0 < x < L \\ u(x, 0) = u_0(x), & 0 < x < L \\ u(L, t) = g(t) \end{cases}$$

Note: A boundary condition is imposed only at $x = L$, not at $x = 0$.

Remark:

- A boundary condition is imposed only at the boundary where the characteristics go into the computational domain.

- No boundary condition is imposed at the boundary where the characteristics go out of the computational domain.

We consider the case of $a > 0$ and study several numerical methods.

The upwind method

Exact solution: $u(x, t) = u_0(x - at)$

Information is coming from the left. We go to the left side to get information.

$$u_x(x_i, t_n) \approx \frac{u_i^n - u_{i-1}^n}{\Delta x}, \quad u_t(x_i, t_n) \approx \frac{u_i^{n+1} - u_i^n}{\Delta t}$$

$$u_t = -au_x \implies \frac{u_i^{n+1} - u_i^n}{\Delta t} \approx -a \frac{u_i^n - u_{i-1}^n}{\Delta x}$$

The upwind method:

$$u_i^{n+1} = u_i^n - ar(u_i^n - u_{i-1}^n), \quad r = \frac{\Delta t}{\Delta x}, \quad a > 0$$

This is called the upwind method since we go upstream of characteristics propagation to seek information for calculating u^{n+1} .

Notice the **new definition of r** for numerical solution of hyperbolic PDEs.

Local truncation error of upwind method:

$$e_i^n(\Delta x, \Delta t) = \Delta t \cdot O(\Delta t + \Delta x)$$

It is 1st order in time and 1st order in space.

Stability of upwind method (an alternative analysis):

We write the upwind method as

$$u_i^{n+1} = (1 - ar)u_i^n + ar u_{i-1}^n$$

- The case of $0 \leq (ar) \leq 1$

$$|u_i^{n+1}| \leq (1 - ar)|u_i^n| + ar|u_{i-1}^n| \leq (1 - ar)\sup_k |u_k^n| + (ar)\sup_k |u_k^n| = \sup_k |u_k^n|$$

Since this is valid at any i , we have

$$\sup_k |u_k^{n+1}| \leq \sup_k |u_k^n|$$

\implies The method is stable for $(ar) \leq 1$.

- The case of $(ar) > 1$

We try $u_i^n = \rho^n (-1)^i$. Substituting it into the upwind method, we get

$$\rho = (1 - ar) + ar(-1) = 1 - 2ar$$

$$\Rightarrow |\rho| = |1 - 2ar| = 2ar - 1 > 1$$

$$\Rightarrow \text{The method is unstable for } (ar) > 1.$$

The analysis above also shows that the method is unstable for $a < 0$.

$$|\rho| = |1 - 2ar| = 1 + (-2ar) > 1 \quad \text{for } a < 0$$

Therefore, the upwind method is stable if and only if $0 \leq (ar) \leq 1$.

von Neumann stability analysis (the universal tool)

We substitute $u_i^n = \rho^n e^{\sqrt{-1}\xi i \Delta x}$ into $u_i^{n+1} = (1 - ar)u_i^n + (ar)u_{i-1}^n$.

$$\rho^{n+1} e^{\sqrt{-1}\xi i \Delta x} = (1 - ar)\rho^n e^{\sqrt{-1}\xi i \Delta x} + (ar)\rho^n e^{\sqrt{-1}\xi (i-1) \Delta x}$$

$$\Rightarrow \rho = (1 - ar) + (ar)e^{-\sqrt{-1}\xi \Delta x}$$

$$\Rightarrow \rho = (1 - ar) + (ar)\cos(\xi \Delta x) - \sqrt{-1}(ar)\sin(\xi \Delta x)$$

$$\begin{aligned} \Rightarrow |\rho|^2 &= (1 - ar)^2 + 2(1 - ar)(ar)\cos(\xi \Delta x) + (ar)^2 \\ &= \left((1 - ar)^2 + 2(1 - ar)(ar) + (ar)^2 \right) + 2(1 - ar)(ar)(\cos(\xi \Delta x) - 1) \\ &= 1 - 4(1 - ar)(ar)\sin^2\left(\frac{\xi \Delta x}{2}\right) \end{aligned}$$

- The case of $0 \leq (ar) \leq 1$

$$|\rho|^2 \leq 1 \quad \text{for all } \xi$$

- The case of $(ar) > 1$

$$\text{At } \xi \Delta x / 2 = \pi / 2, \sin(\xi \Delta x / 2) = 1 \text{ and } |\rho|^2 = 1 + 4(ar - 1)(ar) > 1$$

- The case of $(ar) < 0$

$$\text{At } \xi \Delta x / 2 = \pi / 2, \sin(\xi \Delta x / 2) = 1 \text{ and } |\rho|^2 = 1 + 4(1 - ar)(-ar) > 1$$

In summary, the upwind method is stable if and only if $0 \leq (ar) \leq 1$.

The downwind method:

$$u_i^{n+1} = u_i^n - ar(u_{i+1}^n - u_i^n), \quad r = \frac{\Delta t}{\Delta x}, \quad a > 0$$

This is called the downwind method since we go downstream of characteristics propagation to seek information for calculating u^{n+1} .

Local truncation error of downwind method:

$$e_i^n(\Delta x, \Delta t) = \Delta t \cdot O(\Delta t + \Delta x)$$

It is 1st order in time and 1st order in space.

Stability of downwind method:

We write the downwind method as

$$u_i^{n+1} = (1 + ar)u_i^n - ar u_{i+1}^n$$

We carry out the von Neumann stability analysis

...

$$|\rho|^2 = 1 + 4(1 + ar)(ar) \sin^2\left(\frac{\xi \Delta x}{2}\right)$$

...

In summary, the downwind method is unstable for $a > 0$.

The FTCS method

We use central difference in the spatial dimension.

$$u_i^{n+1} = u_i^n - \frac{ar}{2}(u_{i+1}^n - u_{i-1}^n), \quad r = \frac{\Delta t}{\Delta x}$$

Local truncation error of FTCS method:

$$e_i^n(\Delta x, \Delta t) = \Delta t \cdot O(\Delta t + (\Delta x)^2)$$

It is 1st order in time and 2nd order in space.

von Neumann stability analysis of FTCS method:

We substitute $u_i^n = \rho^n e^{\sqrt{-1}\xi i \Delta x}$ into $u_i^{n+1} = u_i^n - \frac{ar}{2}(u_{i+1}^n - u_{i-1}^n)$.

$$\rho^{n+1} e^{\sqrt{-1}\xi i \Delta x} = \rho^n e^{\sqrt{-1}\xi i \Delta x} - \frac{ar}{2} \left[e^{\sqrt{-1}\xi (i+1) \Delta x} - e^{\sqrt{-1}\xi (i-1) \Delta x} \right]$$

$$\Rightarrow \rho = 1 - \frac{ar}{2} \left[e^{\sqrt{-1}\xi \Delta x} - e^{-\sqrt{-1}\xi \Delta x} \right]$$

$$= 1 - \sqrt{-1}(ar) \sin(\xi \Delta x)$$

$$\Rightarrow |\rho|^2 = 1 + (ar)^2 \sin^2(\xi \Delta x)$$

At $\xi \Delta x = \pi/2$, $\sin(\xi \Delta x) = 1$ and

$$|\rho|^2 = 1 + (ar)^2 > 1$$

==> The FTCS method is unstable for any (ar) .

In summary, the FTCS method is unconditionally unstable.

Next we introduce the Lax-Friedrichs method, which is a stable variation of FTCS.

Lax-Friedrichs method

$$u_i^{n+1} = \frac{u_{i+1}^n + u_{i-1}^n}{2} - \frac{ar}{2}(u_{i+1}^n - u_{i-1}^n), \quad r = \frac{\Delta t}{\Delta x}$$

Note that the u_i^n in FTCS is replaced by $(u_{i+1}^n + u_{i-1}^n)/2$ in Lax-Friedrichs.

Local truncation error of Lax-Friedrichs method:

$$e_i^n(\Delta x, \Delta t) = \Delta t \cdot O(\Delta t + \Delta x) \quad \text{for fixed } r$$

It is 1st order in time and 1st order in space.

Stability of Lax-Friedrichs method (an alternative analysis)

We write the Lax-Friedrichs method as

$$u_i^{n+1} = \frac{1}{2}(1+ar)u_{i+1}^n + \frac{1}{2}(1-ar)u_{i-1}^n$$

- For $|ar| \leq 1$,

$$\begin{aligned} |u_i^{n+1}| &\leq \frac{1}{2}(1+ar)|u_{i+1}^n| + \frac{1}{2}(1-ar)|u_{i-1}^n| \\ &\leq \frac{1}{2}(1+ar)\sup_k |u_k^n| + \frac{1}{2}(1-ar)\sup_k |u_k^n| = \sup_k |u_k^n| \end{aligned}$$

Since this is valid at any i , we have

$$\sup_k |u_k^{n+1}| \leq \sup_k |u_k^n|$$

==> The method is stable for $|ar| \leq 1$.

- For $|ar| > 1$, we try $u_i^n = \rho^n (\sqrt{-1})^i$.

Substituting it into the Lax-Friedrichs method, we get

$$\rho = \frac{1}{2}(1+ar)\sqrt{-1} + \frac{1}{2}(1-ar)(\sqrt{-1})^{-1} = \frac{1}{2}(1+ar)\sqrt{-1} - \frac{1}{2}(1-ar)\sqrt{-1} = (ar)\sqrt{-1}$$

==> $|\rho| = |ar| > 1$

==> The method is unstable for $|ar| > 1$.

In summary, the Lax-Friedrichs method is stable if and only if $|a r| \leq 1$.

Remarks:

Lax-Friedrichs method has the properties below

- It is consistent.
- It is stable for $|a r| \leq 1$.
- It has a unified form for both the case of $a > 0$ and the case of $a < 0$.

But Lax-Friedrichs is only 1st order in time and in space.

We want a method having these properties with 2nd order in time and space.

Appendix A Stability of the FTCS method for solving the 2D heat equation

2D FTCS method:

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

We substitute $u_{i,j}^n = \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)}$ into the 2D FTCS method.

$$\begin{aligned} \rho^{n+1} e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)} &= \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)} + \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)} \frac{\Delta t}{(\Delta x)^2} (e^{\sqrt{-1}\xi \Delta x} - 2 + e^{-\sqrt{-1}\xi \Delta x}) \\ &\quad + \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)} \frac{\Delta t}{(\Delta y)^2} (e^{\sqrt{-1}\eta \Delta y} - 2 + e^{-\sqrt{-1}\eta \Delta y}) \\ \implies \rho &= 1 + \frac{\Delta t}{(\Delta x)^2} (e^{\sqrt{-1}\xi \Delta x} - 2 + e^{-\sqrt{-1}\xi \Delta x}) + \frac{\Delta t}{(\Delta y)^2} (e^{\sqrt{-1}\eta \Delta y} - 2 + e^{-\sqrt{-1}\eta \Delta y}) \\ \implies \rho &= 1 - 4 \frac{\Delta t}{(\Delta x)^2} \sin^2\left(\frac{\xi \Delta x}{2}\right) - 4 \frac{\Delta t}{(\Delta y)^2} \sin^2\left(\frac{\eta \Delta y}{2}\right) \end{aligned}$$

Assuming the ratios $\Delta t/(\Delta x)^2$ and $\Delta t/(\Delta y)^2$ remain constant as $(\Delta x, \Delta y, \Delta t) \rightarrow 0$.

The 2D FTCS method is stable if and only if $\rho > -1$ for all (ξ, η)

$$\text{if and only if } 4 \frac{\Delta t}{(\Delta x)^2} \sin^2\left(\frac{\xi \Delta x}{2}\right) + 4 \frac{\Delta t}{(\Delta y)^2} \sin^2\left(\frac{\eta \Delta y}{2}\right) \leq 2 \text{ for all } (\xi, \eta)$$

$$\text{if and only if } \frac{\Delta t}{(\Delta x)^2} + \frac{\Delta t}{(\Delta y)^2} \leq \frac{1}{2}.$$

Appendix B Stability of the ADI method for solving the 2D heat equation

Alternating direction implicit method (ADI)

$$\begin{cases} u^{n+1/2} = u^n + \underbrace{\frac{\Delta t}{2} D_x^2 u^{n+1/2}}_{\text{Implicit}} + \underbrace{\frac{\Delta t}{2} D_y^2 u^n}_{\text{Explicit}} \\ u^{n+1} = u^{n+1/2} + \underbrace{\frac{\Delta t}{2} D_y^2 u^{n+1/2}}_{\text{Explicit}} + \underbrace{\frac{\Delta t}{2} D_x^2 u^{n+1}}_{\text{Implicit}} \end{cases}$$

We substitute $u_{i,j}^n = \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)}$ and $u_{i,j}^{n+1/2} = \omega \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)}$ into the first half.

$$\begin{aligned} \omega &= 1 + \omega \frac{\Delta t}{(\Delta x)^2} \left(e^{\sqrt{-1} \xi \Delta x} - 2 + e^{-\sqrt{-1} \xi \Delta x} \right) + \frac{\Delta t}{(\Delta y)^2} \left(e^{\sqrt{-1} \eta \Delta y} - 2 + e^{-\sqrt{-1} \eta \Delta y} \right) \\ \Rightarrow \quad \omega \left(1 + 4 \frac{\Delta t}{(\Delta x)^2} \sin^2 \left(\frac{\xi \Delta x}{2} \right) \right) &= \left(1 - 4 \frac{\Delta t}{(\Delta y)^2} \sin^2 \left(\frac{\eta \Delta x}{2} \right) \right) \end{aligned} \quad (\text{E01A})$$

We substitute $u_{i,j}^{n+1/2} = \omega \rho^n e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)}$ and $u_{i,j}^{n+1} = \rho^{n+1} e^{\sqrt{-1}(\xi i \Delta x + \eta j \Delta y)}$ into the second half.

$$\begin{aligned} \rho &= \omega + \omega \frac{\Delta t}{(\Delta x)^2} \left(e^{\sqrt{-1} \xi \Delta x} - 2 + e^{-\sqrt{-1} \xi \Delta x} \right) + \rho \frac{\Delta t}{(\Delta y)^2} \left(e^{\sqrt{-1} \eta \Delta y} - 2 + e^{-\sqrt{-1} \eta \Delta y} \right) \\ \Rightarrow \quad \rho \left(1 + 4 \frac{\Delta t}{(\Delta y)^2} \sin^2 \left(\frac{\eta \Delta x}{2} \right) \right) &= \omega \left(1 - 4 \frac{\Delta t}{(\Delta x)^2} \sin^2 \left(\frac{\xi \Delta x}{2} \right) \right) \end{aligned} \quad (\text{E01B})$$

Combining (E01A) and (E01B), we obtain

$$\rho = \frac{\left(1 - 4 \frac{\Delta t}{(\Delta x)^2} \sin^2 \left(\frac{\xi \Delta x}{2} \right) \right)}{\left(1 + 4 \frac{\Delta t}{(\Delta x)^2} \sin^2 \left(\frac{\xi \Delta x}{2} \right) \right)} \cdot \frac{\left(1 - 4 \frac{\Delta t}{(\Delta y)^2} \sin^2 \left(\frac{\eta \Delta x}{2} \right) \right)}{\left(1 + 4 \frac{\Delta t}{(\Delta y)^2} \sin^2 \left(\frac{\eta \Delta x}{2} \right) \right)}$$

It is straightforward to show that $|\rho| \leq 1$ for all (ξ, η) .