

(Ch. 1)

## Lecture 2: Consistency & Stability

Lax Equivalence Theorem states that consistency & stability = convergence; this has been the cornerstone of numerical methods for differential equations.

Consistency To understand consistency, let's look at the PDE

$$(2.1) \quad \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = 0 \quad (\text{assume appropriate ICs \& BCs})$$

If we approximate the derivative using a general order approximation, we get:

$$(2.2) \quad \sum_{n=-N}^{+N} \frac{\alpha_n u_i^{n+1}}{\Delta t} + \sum_{n=-N}^{+N} B_n u_i^n = 0$$

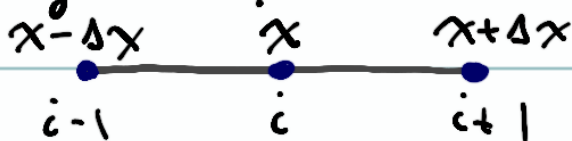
which yields an order  $O(\Delta t^{2N}, \Delta x^{2N})$  method.

But let's simplify the discussion and assume  $N=1$

to get:

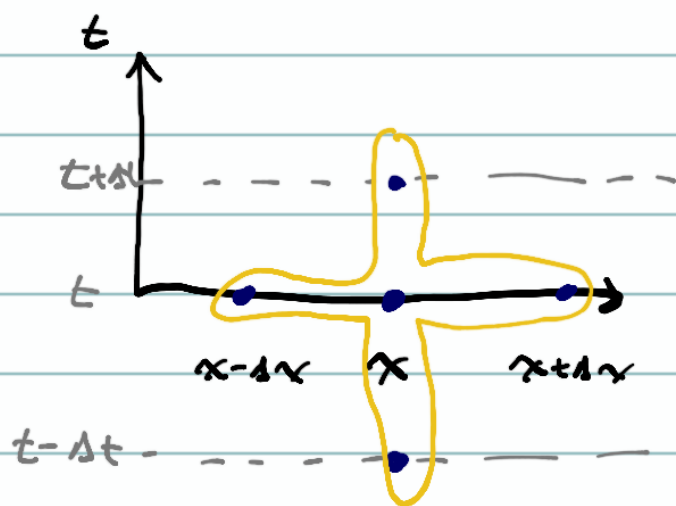
$$(2.3) \quad \frac{\alpha_{-1} u_i^{n+1} + \alpha_0 u_i^n + \alpha_1 u_i^{n+1}}{\Delta t} + \frac{u B_{-1} u_{i-1}^n + B_0 u_i^n + B_1 u_{i+1}^n}{\Delta x} = 0$$

Taking a Taylor Series expansion about  $(x, t)$  s.t.



is the spatial stencil

which is similar for both time & space. The full stencil in the  $x-t$  diagram is:



Applying the Taylor Series yields:

$$f_i^n = f(x, t) \equiv f$$

$$f_i^{n+1} \equiv f(x, t + \Delta t) = f + \Delta t f_t + \frac{\Delta t^2}{2} f_{tt} + o(\Delta t^3)$$

$$(2.4) \quad f_i^{n-1} \equiv f(x, t - \Delta t) = f - \Delta t f_t + \frac{\Delta t^2}{2} f_{tt} - o(\Delta t^3)$$

$$f_{i+1}^n \equiv f(x + \Delta x, t) = f + \Delta x f_x + \frac{\Delta x^2}{2} f_{xx} + o(\Delta x^3)$$

$$f_{i-1}^n \equiv f(x - \Delta x, t) = f - \Delta x f_x + \frac{\Delta x^2}{2} f_{xx} - o(\Delta x^3)$$

For simplicity, pick  $\alpha_{-1} = 0$  (let's focus on the spatial discretization). Subbing (2.4) into (2.3) yields

$$(2.5) \quad \frac{\alpha_1}{\Delta t} \left( f + \Delta t f_t + \frac{\Delta t^2}{2} f_{tt} + o(\Delta t^3) \right) + \frac{\alpha_0}{\Delta t} f$$

$$+ \frac{u \beta_1}{\Delta x} \left( f + \Delta x f_x + \frac{\Delta x^2}{2} f_{xx} + o(\Delta x^3) \right) + \frac{u \beta_0}{\Delta x} f$$

$$+ \frac{u \beta_{-1}}{\Delta x} \left( f - \Delta x f_x + \frac{\Delta x^2}{2} f_{xx} - o(\Delta x^3) \right) = 0$$

From the time-derivative terms we see that:

(TC1)  $\delta$ :  $\frac{\alpha_1}{\Delta t} + \frac{\alpha_0}{\Delta t} = 0 \rightarrow$  to eliminate  $\delta$  term

(TC2)  $\frac{\partial \delta}{\partial t}$ :  $\frac{\alpha_1}{\Delta t} \Delta t = 1 \rightarrow$  to keep  $\delta_t$  term

$\therefore \alpha_1 = 1$  &  $\alpha_0 = -1$  & we omit  $\alpha_{-1} = 0$

From the space-derivative terms we note:

(SC1)  $\delta$ :  $\beta_1 + \beta_0 + \beta_{-1} = 0 \rightarrow$  to eliminate  $\delta$  term

(SC2)  $\frac{\partial \delta}{\partial x}$ :  $\beta_1 - \beta_{-1} = 1 \rightarrow$  to keep  $u \frac{\partial \delta}{\partial x}$  term

2 eqs. & 3 unknowns is over-determined so we do the following: choose  $\beta_0$  is a free parameter. Pick  $\beta_0 = 0$

$$\begin{pmatrix} \beta_1 + \beta_{-1} = -\beta_0 \\ \beta_1 - \beta_{-1} = 1 \end{pmatrix} \quad \text{adding rows 1 \& 2 gives:} \\ 2\beta_1 = 1 - \beta_0 \quad \text{(SC3)}$$

Subbing in the coefficients into Eq. (2.3) yields:

$$\delta_t + u \delta_x + \underbrace{u(\beta_1 + \beta_{-1}) \frac{\Delta x}{2} \delta_{xx}}_{\text{truncation error term}} + o(\Delta t, \Delta x^2) = 0$$

Let's 3 possible choices for our free parameter  $\beta_0$ .

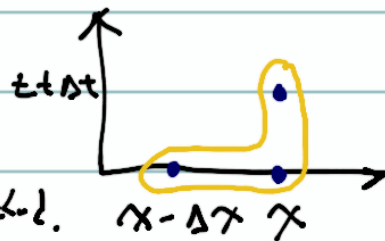
**Case 1** Let  $\beta_0 = 1$  which, from SC3,  $\beta_1 = 0$   
 & from SC2  $\beta_{-1} = -1$  that yields (from Eq. (2.3))  

$$\frac{f_i^{n+1} - f_i^n}{\Delta t} + u \frac{f_i^n - f_{i-1}^n}{\Delta x} + u(0-1) \frac{\Delta x}{2} f_{xx} + O(\Delta t, \Delta x^3) = 0$$

or

$$\frac{f_i^{n+1} - f_i^n}{\Delta t} + u \frac{f_i^n - f_{i-1}^n}{\Delta x} + O(\Delta t, \Delta x) = 0$$

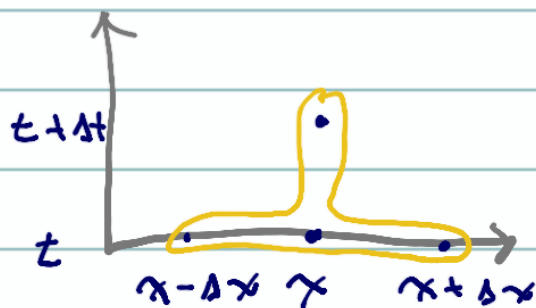
With the corresponding stencil which for  $u > 0$  is an upwinding method.



**Case 2** Let  $\beta_0 = 0$  which yields  $\beta_1 = \frac{1}{2}$  (SC3)  
 &  $\beta_{-1} = -\frac{1}{2}$  (SC2) to get:

$$\frac{f_i^{n+1} - f_i^n}{\Delta t} + u \frac{f_{i+1}^n - f_{i-1}^n}{2\Delta x} + O(\Delta t, \Delta x^2) = 0$$

w/ the stencil

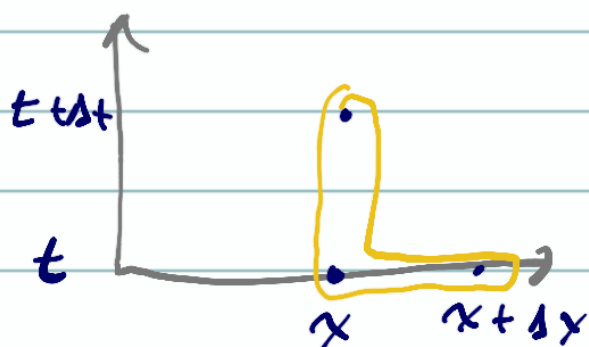


which is a centered method in space.

**Case 3** Let  $B_0 = -1$  which gives  $B_1 = 1$  (SC1)  
 &  $B_{-1} = 0$  (SC2) to get:

$$\frac{f_i^{n+1} - f_i^n}{\Delta t} + u \frac{f_{i+1}^n - f_i^n}{\Delta x} + O(\Delta t, \Delta x) = 0$$

w/ Stencil



which for  $u > 0$  is not an upwind method but rather a downwind method

Recap of Consistency We saw that for the 1D wave equation:

(2.1)  $\frac{\partial f}{\partial t} + u \frac{\partial f}{\partial x} = 0$  we can discretize as

(2.2)  $\sum_{n=-K}^{+K} \frac{\alpha_n f_i^{n+1}}{\Delta t} + u \sum_{n=-N}^{+N} \frac{B_n f_{i+n}^n}{\Delta x} = 0$

where "i" denotes a position in space & "n" in time.

Eq. (2.2) defines an  $O(\leq 2K, \leq 2N)$  method.

For, e.g.,  $N=1$  we have  $n=-1, 0, +1$  or  $2N+1$  points which gives an order  $2N$  approximation

Applying a Taylor Series, we obtained:

$$\delta t + u \rho x + u (B_1 + B_{-1}) \frac{\Delta x}{2} \delta \pi x + o(\Delta t, \Delta x^3) = 0$$

## Key Points

1. as  $(\Delta x, \Delta t) \rightarrow 0$  we recover the original continuous PDE  $\rightarrow$  consistent
  2. For  $B_1 + B_{-1} \neq 0$  the method is  $O(\Delta t, \Delta x)$
  3. For  $B_1 + B_{-1} = 0$  the method is  $O(\Delta t, \Delta x^2)$
- } order  $\leq 2N$ ,  $\leq 2N$

Stability We will use Von Neumann stability analysis & matrix stability to check for the stability of methods. In this section we shall only discuss Von Neumann analysis.

In Fourier analysis we write the variable  $\phi_j$  as  $\phi_j^n \equiv \phi(x_j, t^n) = \sum_{l=-N}^{+N} \tilde{\phi}_l^{(n)} e^{i k_l j \Delta x}$

with the following def:

$f_n^{(n)}$  → multiplier w/ exponent  $n$   
 $e^{i n \omega \Delta x}$  → frequencies with  $i = \sqrt{-1}$

He  $\rightarrow$  here number

$U_L = U_L \Delta x \rightarrow$  phase angle

For  $\nu \rightarrow 0$  we get low frequency (long waves) & for

$\omega \rightarrow \pi$  we get high frequency (short waves)

$N$  denotes the number of points in the domain.

To simplify our discussion, let us only consider a single Fourier mode & write it as follows:

$$q_j^n = \tilde{q}(n) e^{ij\omega}$$

Since  $\tilde{q}$  is an amplitude then we must insist that  $|\tilde{q}(n)| < \infty$ . For this to be satisfied requires that  $|\tilde{q}| \leq 1$ .

Let's now apply this to Cases 1, 2, & 3.

**Case 1** For the upwind method  $\frac{q_j^{n+1} - q_j^n}{\Delta t} + u \frac{q_j^n - q_{j-1}^n}{\Delta x} = 0$

We get:

(2.6)  $q_j^{n+1} = q_j^n - \underbrace{\frac{u \Delta t}{\Delta x}}_{\text{Courant Number}} (q_j^n - q_{j-1}^n)$  & let  $\sigma = \frac{u \Delta t}{\Delta x}$

Writing in terms of Fourier modes yields:

(2.7) 
$$\begin{aligned} q_j^n &= \tilde{q}(n) e^{ij\omega} \\ q_j^{n+1} &= \tilde{q}(n+1) e^{ij\omega} \\ q_{j-1}^n &= \tilde{q}(n) e^{i(j-1)\omega} \end{aligned}$$

Subbing (2.7) into (2.6) yields:



$$(2.8) \quad \tilde{f}^{(n+1)} e^{i\omega} = \tilde{f}^{(n)} e^{i\omega} - \sigma \left[ \tilde{f}^{(n)} e^{i\omega} - \tilde{f}^{(n)} e^{i(j-1)\omega} \right]$$

Dividing by  $\tilde{f}^{(n)} e^{i\omega}$  gives:

$$\tilde{f} = 1 - \sigma (1 - e^{-i\omega}) \rightarrow \text{Euler's formula: } e^{\pm i\omega} = \cos \omega \pm i \sin \omega \text{ gives:}$$

$$\tilde{f} = 1 - \sigma [1 - (\cos \omega - i \sin \omega)] = 1 - \sigma + \sigma \cos \omega - i \sigma \sin \omega$$

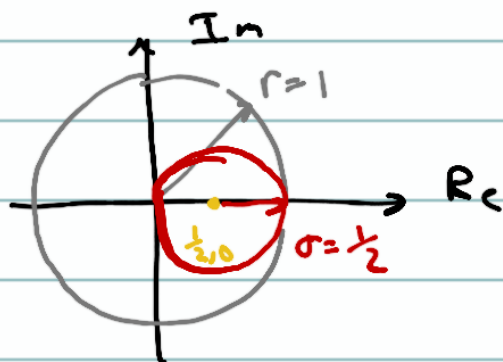
where  $\text{Re}(\tilde{f}) = (1 - \sigma) + \sigma \cos \omega$  &  $\text{Im}(\tilde{f}) = -\sigma \sin \omega$  that defines a circle of radius  $\sigma$  w/ center  $(1 - \sigma, 0)$  in the complex plane. Recall that we need

$|\tilde{f}| \leq 1$  which gives:

$$|\tilde{f}|^2 = \tilde{f} \tilde{f}^* = [(1 - \sigma) + \sigma \cos \omega]^2 + (\sigma \sin \omega)^2 \leq 1$$

Note that  $|\tilde{f}| = (\tilde{f} \tilde{f}^*)^{1/2}$  defines a circle of radius  $\sigma$  w/ center  $(1 - \sigma, 0)$  which will only occur for  $\sigma \leq 1$  (Conditionally Stable)

Here is a schematic:



Note that for  $\sigma = 1$ , the circle starts at  $(0, 0)$

**Case 2** For the centered method:  $\frac{f_j^{n+1} - f_j^n}{\Delta t} + u \frac{f_{j+1}^n - f_{j-1}^n}{2\Delta x} = 0$   
we require the following Fourier representation



$$\begin{aligned}\hat{f}_j^n &= \tilde{f}^{(n)} e^{ij\varphi} \\ \hat{f}_{j+1}^{n+1} &= \tilde{f}^{(n+1)} e^{i(j+1)\varphi} \\ \hat{f}_j^{n+1} &= \tilde{f}^{(n)} e^{i(j+1)\varphi} \\ \hat{f}_{j+1}^n &= \tilde{f}^{(n)} e^{i(j-1)\varphi} \\ \hat{f}_{j-1}^n &= \tilde{f}^{(n)} e^{i(j-1)\varphi}\end{aligned}$$

which we sub into the differencing stencil to get:

$$\tilde{f}^{(n+1)} e^{ij\varphi} = \tilde{f}^{(n)} e^{ij\varphi} - \frac{u \Delta t}{2\Delta x} \left[ \tilde{f}^{(n)} e^{i(j+1)\varphi} - \tilde{f}^{(n)} e^{i(j-1)\varphi} \right]$$

Dividing by  $\hat{f}_j^n = \tilde{f}^{(n)} e^{ij\varphi}$  gives:

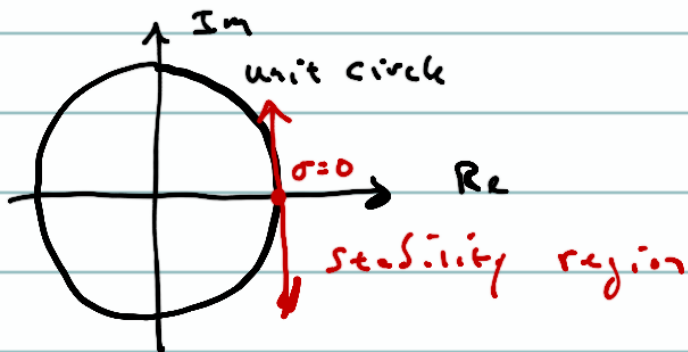
$$\tilde{f} = 1 - \frac{\sigma}{2} [e^{i\varphi} - e^{-i\varphi}] \rightarrow \text{using Euler's formula}$$

$$\tilde{f} = 1 - \frac{\sigma}{2} [(\cancel{\cos\varphi} + i\sin\varphi) - (\cancel{\cos\varphi} - i\sin\varphi)]$$

$$= 1 - i\sigma \sin\varphi$$

$$\operatorname{Re}(\tilde{f}) = 1 \quad \operatorname{Im}(\tilde{f}) = -\sigma \sin\varphi$$

$$\& \quad \tilde{f} \tilde{f}^* = 1 + \sigma^2 \sin^2\varphi \leq 1 \quad \text{only for } \sigma = 0$$



**Case 3** For the downwind method

$$\frac{\tilde{f}_i^{n+1} - \tilde{f}_i^n}{\Delta t} + u \frac{\tilde{f}_{i+1}^n - \tilde{f}_i^n}{\Delta x} = 0$$

We get in Fourier space:

$$\tilde{f}^{(n+1)} e^{i j \psi} = \tilde{f}^{(n)} e^{i j \psi} - \frac{u \Delta t}{\Delta x} \left[ \tilde{f}^{(n)} e^{i(j+1)\psi} - \tilde{f}^{(n)} e^{i j \psi} \right]$$

Dividing by  $\tilde{f}^{(n)} e^{i j \psi}$  gives:

$$\tilde{f} = 1 - \sigma \left[ e^{i \psi} - 1 \right] = 1 - \sigma \left[ \cos \psi + i \sin \psi - 1 \right]$$

$$\tilde{f} = 1 + \sigma - \sigma \cos \psi - i \sigma \sin \psi$$

$$\operatorname{Re}(\tilde{f}) = 1 + \sigma - \sigma \cos \psi \quad \operatorname{Im}(\tilde{f}) = -\sigma \sin \psi$$

$$\tilde{f} \tilde{f}^* = (1 + \sigma - \sigma \cos \psi)^2 + (\sigma \sin \psi)^2 \leq 1$$

which forms a circle of radius  $\sigma$  centered at  $(1 + \sigma, 0)$ .

Here's a Schematic:

