

## §3.2. Initial Value Problems

Hyperbolic and parabolic PDEs are **initial value problems** or **marching problems** (a term introduced by Richardson).

## §3.2. Initial Value Problems

Hyperbolic and parabolic PDEs are **initial value problems** or **marching problems** (a term introduced by Richardson).

The solution is obtained by using the known initial values and marching or advancing in time.

## §3.2. Initial Value Problems

Hyperbolic and parabolic PDEs are **initial value problems** or **marching problems** (a term introduced by Richardson).

The solution is obtained by using the known initial values and marching or advancing in time.

If boundary values are necessary, they are called **mixed initial-boundary value problems**.

## §3.2. Initial Value Problems

Hyperbolic and parabolic PDEs are **initial value problems** or **marching problems** (a term introduced by Richardson).

The solution is obtained by using the known initial values and marching or advancing in time.

If boundary values are necessary, they are called **mixed initial-boundary value problems**.

The simplest **prototypes** of these initial value problems are:

## §3.2. Initial Value Problems

Hyperbolic and parabolic PDEs are **initial value problems** or **marching problems** (a term introduced by Richardson).

The solution is obtained by using the known initial values and marching or advancing in time.

If boundary values are necessary, they are called **mixed initial-boundary value problems**.

The simplest **prototypes** of these initial value problems are:

- The **advection equation** (with solution  $u(x, t) = u(x - ct, 0)$ )

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

which is a hyperbolic equation.

## §3.2. Initial Value Problems

Hyperbolic and parabolic PDEs are **initial value problems** or **marching problems** (a term introduced by Richardson).

The solution is obtained by using the known initial values and marching or advancing in time.

If boundary values are necessary, they are called **mixed initial-boundary value problems**.

The simplest **prototypes** of these initial value problems are:

- The **advection equation** (with solution  $u(x, t) = u(x - ct, 0)$ )

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

which is a hyperbolic equation.

- The **diffusion equation**,

$$\frac{\partial u}{\partial t} = \sigma \frac{\partial^2 u}{\partial x^2}$$

which is a parabolic equation.

# The Finite Difference Method

We take *discrete values* for  $x$  and  $t$ :  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

# The Finite Difference Method

We take *discrete values* for  $x$  and  $t$ :  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

The solution of the finite difference equation is also defined at the discrete points  $(x_j, t_n) = (j\Delta x, n\Delta t)$ :

$$U_j^n = U(j\Delta x, n\Delta t) = U(x_j, t_n).$$



# The Finite Difference Method

We take *discrete values* for  $x$  and  $t$ :  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

The solution of the finite difference equation is also defined at the discrete points  $(x_j, t_n) = (j\Delta x, n\Delta t)$ :

$$U_j^n = U(j\Delta x, n\Delta t) = U(x_j, t_n).$$

That is, we use a **small**  $u$  to denote the solution of the PDE (*continuous*) and a **capital**  $U$  to denote the solution of the finite difference equation (FDE, a *discrete solution*).

# The Finite Difference Method

We take *discrete values* for  $x$  and  $t$ :  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

The solution of the finite difference equation is also defined at the discrete points  $(x_j, t_n) = (j\Delta x, n\Delta t)$ :

$$U_j^n = U(j\Delta x, n\Delta t) = U(x_j, t_n).$$

That is, we use a **small**  $u$  to denote the solution of the PDE (*continuous*) and a **capital**  $U$  to denote the solution of the finite difference equation (FDE, a *discrete solution*).

Consider again the advection equation

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

# The Finite Difference Method

We take *discrete values* for  $x$  and  $t$ :  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

The solution of the finite difference equation is also defined at the discrete points  $(x_j, t_n) = (j\Delta x, n\Delta t)$ :

$$U_j^n = U(j\Delta x, n\Delta t) = U(x_j, t_n).$$

That is, we use a **small**  $u$  to denote the solution of the PDE (*continuous*) and a **capital**  $U$  to denote the solution of the finite difference equation (FDE, a *discrete solution*).

Consider again the advection equation

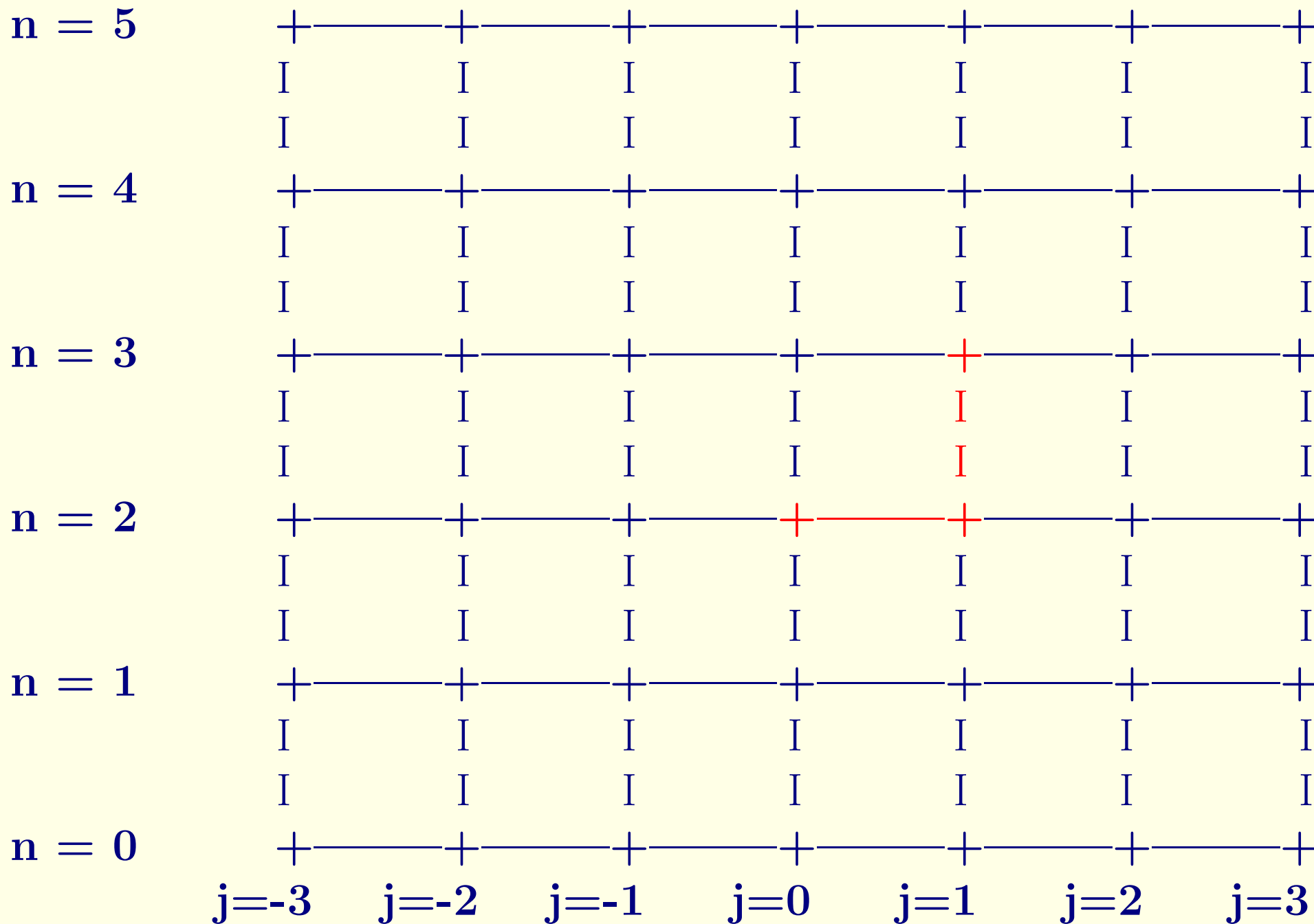
$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

Suppose we choose to approximate this **PDE** with the **FDE**

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

# Space-Time Grid:

Space axis horizontal  
Time axis vertical



To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

Note that both differences are **non-centered** with respect to the point  $(x_j, t_n) = (j\Delta x, n\Delta t)$ .

To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

Note that both differences are **non-centered** with respect to the point  $(x_j, t_n) = (j\Delta x, n\Delta t)$ .

We can ask two fundamental questions:



To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

Note that both differences are **non-centered** with respect to the point  $(x_j, t_n) = (j\Delta x, n\Delta t)$ .

We can ask two fundamental questions:

- [1] Is the FDE **consistent** with the PDE?

To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

Note that both differences are **non-centered** with respect to the point  $(x_j, t_n) = (j\Delta x, n\Delta t)$ .

We can ask two fundamental questions:

- [1] Is the FDE **consistent** with the PDE?
- [2] For a given time  $t > 0$ , will the solution of the FDE **converge** to that of the PDE as  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ ?

To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

Note that both differences are **non-centered** with respect to the point  $(x_j, t_n) = (j\Delta x, n\Delta t)$ .

We can ask two fundamental questions:

- [1] Is the FDE **consistent** with the PDE?
- [2] For a given time  $t > 0$ , will the solution of the FDE **converge** to that of the PDE as  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ ?

We will clarify these questions below.

★ ★ ★

To repeat:

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$$

This is called an **upstream scheme** (we are assuming  $c > 0$ ).

Note that both differences are **non-centered** with respect to the point  $(x_j, t_n) = (j\Delta x, n\Delta t)$ .

We can ask two fundamental questions:

- [1] Is the FDE **consistent** with the PDE?
- [2] For a given time  $t > 0$ , will the solution of the FDE **converge** to that of the PDE as  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ ?

We will clarify these questions below.

★ ★ ★

**Warning:** Sometimes superscript  $n$  denotes a *power*; sometimes it is just an index. Be careful!

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

**Definition:** The difference between the PDE and the FDE is called the discretization error or **local truncation error**.

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

**Definition:** The difference between the PDE and the FDE is called the discretization error or **local truncation error**.

Consistency is rather simple to verify:



# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

**Definition:** The difference between the PDE and the FDE is called the discretization error or **local truncation error**.

Consistency is rather simple to verify:

- Substitute  $u$  instead of  $U$  in the FDE.

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

**Definition:** The difference between the PDE and the FDE is called the discretization error or **local truncation error**.

Consistency is rather simple to verify:

- Substitute  $u$  instead of  $U$  in the FDE.
- Evaluate all terms using a Taylor series expansion centered on the point  $(x_j, t_n)$ .

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

**Definition:** The difference between the PDE and the FDE is called the discretization error or **local truncation error**.

Consistency is rather simple to verify:

- Substitute  $u$  instead of  $U$  in the FDE.
- Evaluate all terms using a Taylor series expansion centered on the point  $(x_j, t_n)$ .
- Subtract the PDE from the FDE.

# Truncation Errors and Consistency

The FDE is **consistent** with the PDE if, in the limit  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , the FDE coincides with the PDE.

Obviously, this is a basic requirement that the FDE should fulfill if its solutions are going to be good approximations of the solutions of the PDE.

**Definition:** The difference between the PDE and the FDE is called the discretization error or **local truncation error**.

Consistency is rather simple to verify:

- Substitute  $u$  instead of  $U$  in the FDE.
- Evaluate all terms using a Taylor series expansion centered on the point  $(x_j, t_n)$ .
- Subtract the PDE from the FDE.

If the difference (local truncation error) goes to zero as  $\Delta x \rightarrow 0, \Delta t \rightarrow 0$ , then the FDE is consistent with the PDE.

# Example

Let us verify the consistency of the upstream scheme for the advection equation.

# Example

Let us verify the consistency of the upstream scheme for the advection equation.

First, consider the Taylor series expansion:

$$\left. \begin{aligned} u_j^{n+1} &= \left( u + u_t \Delta t + \frac{1}{2} u_{tt} \Delta t^2 + \dots \right)_j^n \\ u_{j-1}^n &= \left( u - u_x \Delta x + \frac{1}{2} u_{xx} \Delta x^2 - \dots \right)_j^n \end{aligned} \right\}$$

# Example

Let us verify the consistency of the upstream scheme for the advection equation.

First, consider the Taylor series expansion:

$$\left. \begin{aligned} u_j^{n+1} &= \left( u + u_t \Delta t + \frac{1}{2} u_{tt} \Delta t^2 + \dots \right)_j^n \\ u_{j-1}^n &= \left( u - u_x \Delta x + \frac{1}{2} u_{xx} \Delta x^2 - \dots \right)_j^n \end{aligned} \right\}$$

We substitute this in the FDE and obtain

$$\left( u_t + \frac{1}{2} u_{tt} \Delta t + \dots \right)_j^n + c \left( u_x - \frac{1}{2} u_{xx} \Delta x + \dots \right)_j^n \simeq 0$$

# Example

Let us verify the consistency of the upstream scheme for the advection equation.

First, consider the Taylor series expansion:

$$\left. \begin{aligned} u_j^{n+1} &= \left( u + u_t \Delta t + \frac{1}{2} u_{tt} \Delta t^2 + \dots \right)_j^n \\ u_{j-1}^n &= \left( u - u_x \Delta x + \frac{1}{2} u_{xx} \Delta x^2 - \dots \right)_j^n \end{aligned} \right\}$$

We substitute this in the FDE and obtain

$$\left( u_t + \frac{1}{2} u_{tt} \Delta t + \dots \right)_j^n + c \left( u_x - \frac{1}{2} u_{xx} \Delta x + \dots \right)_j^n \simeq 0$$

Subtracting the PDE gives the **local truncation error**:

$$\tau = \left( \frac{u_{tt}}{2} \right) \Delta t - \left( \frac{c u_{xx}}{2} \right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$



Again, the **local truncation error** is:

$$\tau = \left(\frac{u_{tt}}{2}\right) \Delta t - \left(\frac{cu_{xx}}{2}\right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$

Again, the **local truncation error** is:

$$\tau = \left(\frac{u_{tt}}{2}\right) \Delta t - \left(\frac{cu_{xx}}{2}\right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$

Clearly, as  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$ , we have  $\tau \rightarrow 0$ .

Again, the **local truncation error** is:

$$\tau = \left(\frac{u_{tt}}{2}\right) \Delta t - \left(\frac{cu_{xx}}{2}\right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$

Clearly, as  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$ , we have  $\tau \rightarrow 0$ .

Therefore, **the FDE is consistent.**

★ ★ ★

Again, the **local truncation error** is:

$$\tau = \left(\frac{u_{tt}}{2}\right) \Delta t - \left(\frac{cu_{xx}}{2}\right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$

Clearly, as  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$ , we have  $\tau \rightarrow 0$ .

Therefore, **the FDE is consistent**.

★       ★       ★

Note that both the **time** and the **space** truncation errors are of **first order**, because the finite differences are **uncentered** in both space and time.

Again, the **local truncation error** is:

$$\tau = \left(\frac{u_{tt}}{2}\right) \Delta t - \left(\frac{cu_{xx}}{2}\right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$

Clearly, as  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$ , we have  $\tau \rightarrow 0$ .

Therefore, **the FDE is consistent.**

★ ★ ★

Note that both the **time** and the **space** truncation errors are of **first order**, because the finite differences are **uncentered** in both space and time.

Truncation errors for **centered differences** are **second order**.

Therefore, in general, centered differences are more accurate than uncentered differences.

★ ★ ★

Again, the **local truncation error** is:

$$\tau = \left(\frac{u_{tt}}{2}\right) \Delta t - \left(\frac{cu_{xx}}{2}\right) \Delta x + \text{H.O.T.} = O(\Delta t) + O(\Delta x)$$

Clearly, as  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$ , we have  $\tau \rightarrow 0$ .

Therefore, **the FDE is consistent**.

★ ★ ★

Note that both the **time** and the **space** truncation errors are of **first order**, because the finite differences are **uncentered** in both space and time.

Truncation errors for **centered differences** are **second order**.

Therefore, in general, centered differences are more accurate than uncentered differences.

★ ★ ★

**Truncation errors are a crucial factor in determining forecast accuracy in NWP.**

# Convergence and Stability

The second question posed above was whether the solution of the FDE **converges** to the PDE solution.

That is, if we let  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , so that  $j\Delta x \rightarrow x$  and  $n\Delta t \rightarrow t$ , does  $U(j\Delta x, n\Delta t) \rightarrow u(x, t)$ ?

# Convergence and Stability

The second question posed above was whether the solution of the FDE **converges** to the PDE solution.

That is, if we let  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , so that  $j\Delta x \rightarrow x$  and  $n\Delta t \rightarrow t$ , does  $U(j\Delta x, n\Delta t) \rightarrow u(x, t)$ ?

This is clearly important, and can be answered by considering another problem, that of **computational stability**.



# Convergence and Stability

The second question posed above was whether the solution of the FDE **converges** to the PDE solution.

That is, if we let  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , so that  $j\Delta x \rightarrow x$  and  $n\Delta t \rightarrow t$ , does  $U(j\Delta x, n\Delta t) \rightarrow u(x, t)$ ?

This is clearly important, and can be answered by considering another problem, that of **computational stability**.

Consider again the advection equation

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

which has the solution  $u(x, t) = u(x - ct, 0)$ .

# Convergence and Stability

The second question posed above was whether the solution of the FDE **converges** to the PDE solution.

That is, if we let  $\Delta x \rightarrow 0$  and  $\Delta t \rightarrow 0$ , so that  $j\Delta x \rightarrow x$  and  $n\Delta t \rightarrow t$ , does  $U(j\Delta x, n\Delta t) \rightarrow u(x, t)$ ?

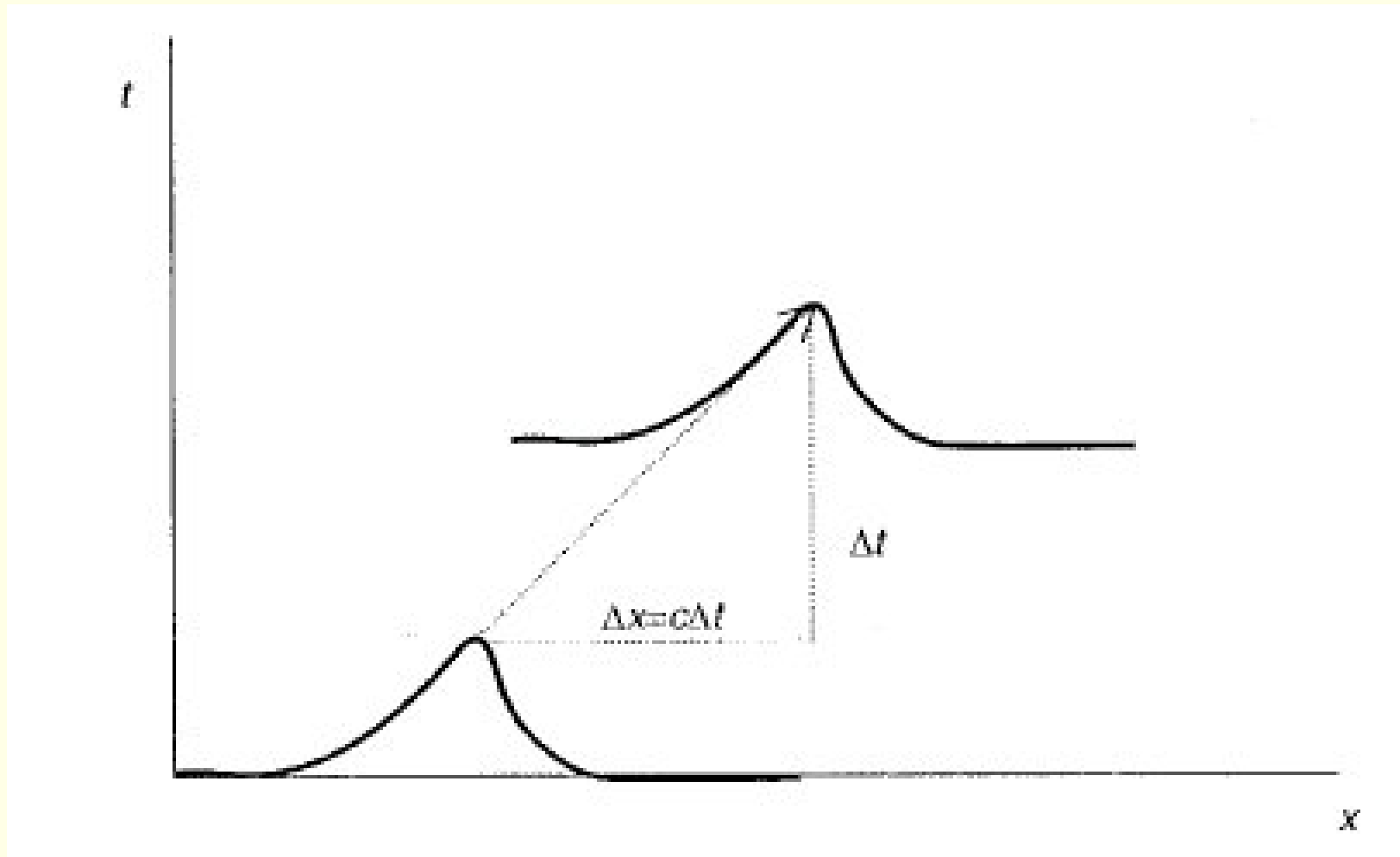
This is clearly important, and can be answered by considering another problem, that of **computational stability**.

Consider again the advection equation

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0,$$

which has the solution  $u(x, t) = u(x - ct, 0)$ .

The shape of the solution  $u(x, 0)$  **translates along the  $x$ -axis** with velocity  $c$  (see Figure below).



Schematic of the **solution** of the advection equation (for  $c > 0$ ).

The FDE for the upstream scheme can be written as

$$U_j^{n+1} = (1 - \mu)U_j^n + \mu U_{j-1}^n$$

where

$$\mu \equiv \frac{c\Delta t}{\Delta x}$$

is the *Courant number* (or *Lewy Number*).

The FDE for the upstream scheme can be written as

$$U_j^{n+1} = (1 - \mu)U_j^n + \mu U_{j-1}^n$$

where

$$\mu \equiv \frac{c\Delta t}{\Delta x}$$

is the *Courant number* (or *Lewy Number*).

Let us suppose that that  $0 \leq \mu \leq 1$ .

The FDE for the upstream scheme can be written as

$$U_j^{n+1} = (1 - \mu)U_j^n + \mu U_{j-1}^n$$

where

$$\mu \equiv \frac{c\Delta t}{\Delta x}$$

is the *Courant number* (or *Lewy Number*).

Let us suppose that that  $0 \leq \mu \leq 1$ .

Then the FDE solution at the new time level  $U_j^{n+1}$  is **interpolated** between the values  $U_j^n$  and  $U_{j-1}^n$ .

The FDE for the upstream scheme can be written as

$$U_j^{n+1} = (1 - \mu)U_j^n + \mu U_{j-1}^n$$

where

$$\mu \equiv \frac{c\Delta t}{\Delta x}$$

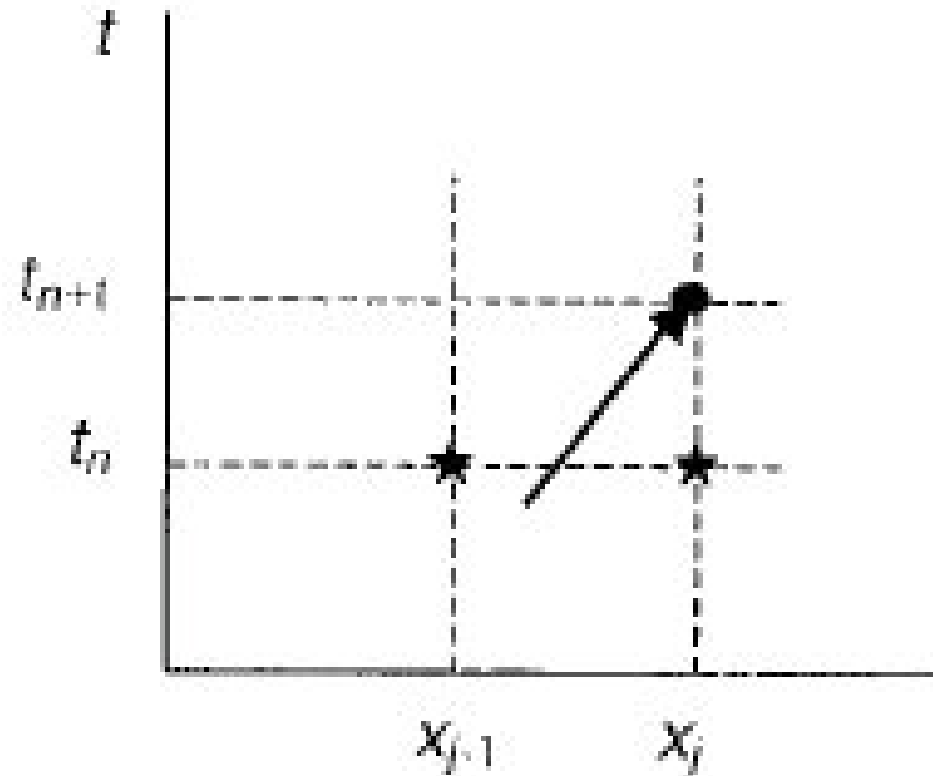
is the *Courant number* (or *Lewy Number*).

Let us suppose that that  $0 \leq \mu \leq 1$ .

Then the FDE solution at the new time level  $U_j^{n+1}$  is **interpolated** between the values  $U_j^n$  and  $U_{j-1}^n$ .

In this case the advection scheme works the way it should, because the true solution lies in between those values.

$$(a) \quad 0 \leq c \leq \frac{\Delta x}{\Delta t}$$



Schematic of the relationship between  $\Delta x$ ,  $\Delta t$  and  $c$  leading to **interpolation** of the solution at time-level  $n + 1$ .

$$0 < \mu \equiv \frac{c\Delta t}{\Delta x} < 1$$



However, suppose this condition is not satisfied, so that

$$\mu = \frac{c\Delta t}{\Delta x} > 1 \quad \text{or else} \quad \mu = \frac{c\Delta t}{\Delta x} < 0.$$

However, **suppose this condition is not satisfied**, so that

$$\mu = \frac{c\Delta t}{\Delta x} > 1 \quad \text{or else} \quad \mu = \frac{c\Delta t}{\Delta x} < 0.$$

Then the parcel arriving at point  $x_j$  at time  $t_{n+1}$  **comes from somewhere outside the interval  $(x_{j-1}, x_j)$  at time  $t_n$ .**

[Recall that  $\partial u/\partial t + c\partial u/\partial x = 0$  is a linear approximation to  $du/dt = 0$ .]

However, **suppose this condition is not satisfied**, so that

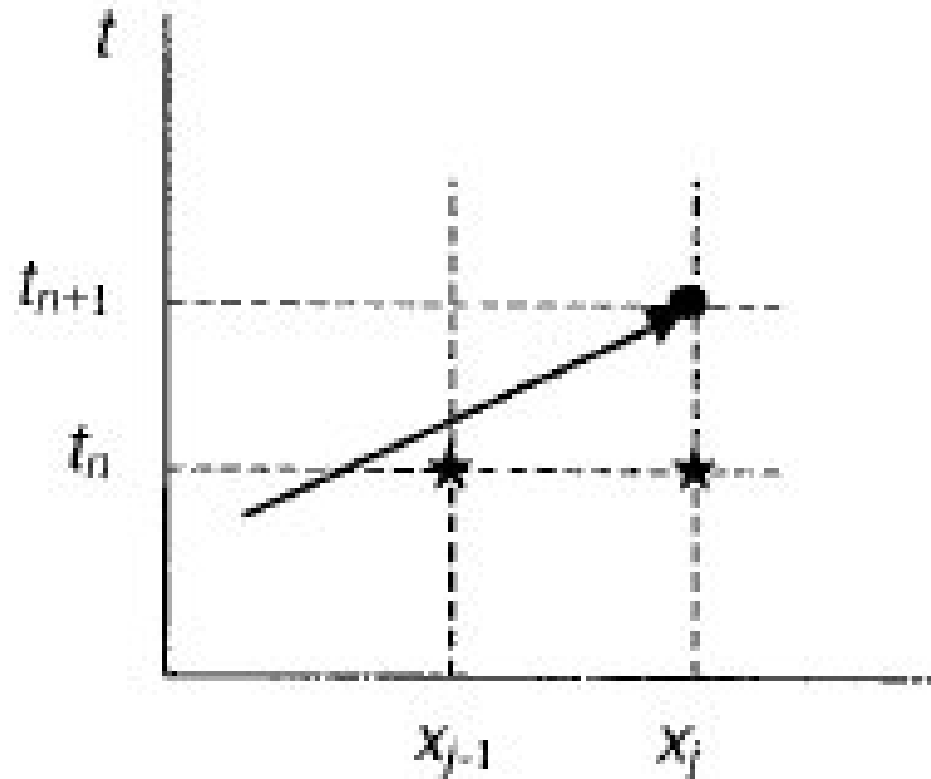
$$\mu = \frac{c\Delta t}{\Delta x} > 1 \quad \text{or else} \quad \mu = \frac{c\Delta t}{\Delta x} < 0.$$

Then the parcel arriving at point  $x_j$  at time  $t_{n+1}$  **comes from somewhere outside the interval  $(x_{j-1}, x_j)$  at time  $t_n$ .**

[Recall that  $\partial u/\partial t + c\partial u/\partial x = 0$  is a linear approximation to  $du/dt = 0$ .]

Thus, the value of  $U_j^{n+1}$  is **extrapolated** from the values  $U_j^n$  and  $U_{j-1}^n$ .

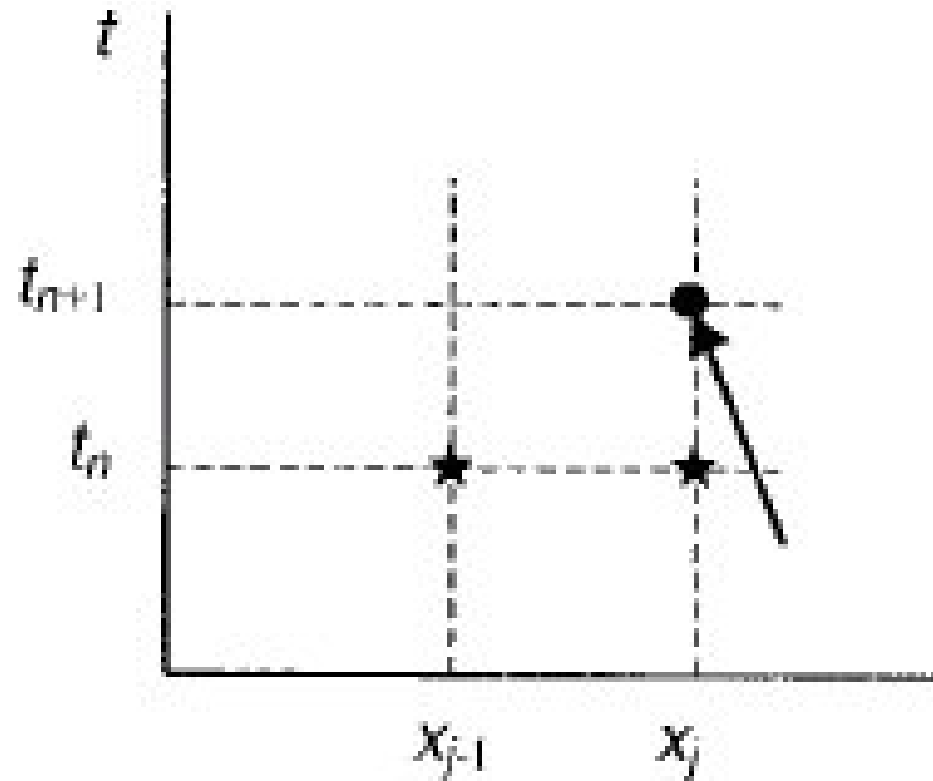
$$(b) 0 \leq \frac{\Delta x}{\Delta t} \leq c$$



Schematic of the relationship between  $\Delta x$ ,  $\Delta t$  and  $c$  leading to **extrapolation** of the solution at time-level  $n + 1$ .

$$\mu \equiv \frac{c\Delta t}{\Delta x} > 1$$

$$(c) \quad c \leq 0 \leq \frac{\Delta x}{\Delta t}$$



Schematic of the relationship between  $\Delta x$ ,  $\Delta t$  and  $c$  leading to **extrapolation** of the solution at time-level  $n + 1$ .

$$\mu \equiv \frac{c\Delta t}{\Delta x} < 0$$

The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

Taking absolute values of the FDE we get

$$|U_j^{n+1}| \leq |U_j^n| |1 - \mu| + |U_{j-1}^n| |\mu|$$

The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

Taking absolute values of the FDE we get

$$|U_j^{n+1}| \leq |U_j^n| |1 - \mu| + |U_{j-1}^n| |\mu|$$

Now defining  $\Upsilon^n = \max_j |U_j^n|$ , we have

$$\Upsilon^{n+1} \leq \{|1 - \mu| + |\mu|\} \Upsilon^n$$

and  $\Upsilon^{n+1} \leq \Upsilon^n$  **if and only if**  $0 \leq \mu \leq 1$ .



The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

Taking absolute values of the FDE we get

$$|U_j^{n+1}| \leq |U_j^n| |1 - \mu| + |U_{j-1}^n| |\mu|$$

Now defining  $\Upsilon^n = \max_j |U_j^n|$ , we have

$$\Upsilon^{n+1} \leq \{|1 - \mu| + |\mu|\} \Upsilon^n$$

and  $\Upsilon^{n+1} \leq \Upsilon^n$  **if and only if**  $0 \leq \mu \leq 1$ .

If the condition  $0 \leq \mu \leq 1$  is **not** satisfied, then **the solution is not bounded** and it grows with  $n$ .

The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

Taking absolute values of the FDE we get

$$|U_j^{n+1}| \leq |U_j^n| |1 - \mu| + |U_{j-1}^n| |\mu|$$

Now defining  $\Upsilon^n = \max_j |U_j^n|$ , we have

$$\Upsilon^{n+1} \leq \{|1 - \mu| + |\mu|\} \Upsilon^n$$

and  $\Upsilon^{n+1} \leq \Upsilon^n$  **if and only if**  $0 \leq \mu \leq 1$ .

If the condition  $0 \leq \mu \leq 1$  is **not** satisfied, then **the solution is not bounded** and it grows with  $n$ .

If we let  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$  with  $\mu = \text{const.}$ , it only makes things worse, because then  $n \rightarrow \infty$ .

The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

Taking absolute values of the FDE we get

$$|U_j^{n+1}| \leq |U_j^n| |1 - \mu| + |U_{j-1}^n| |\mu|$$

Now defining  $\Upsilon^n = \max_j |U_j^n|$ , we have

$$\Upsilon^{n+1} \leq \{|1 - \mu| + |\mu|\} \Upsilon^n$$

and  $\Upsilon^{n+1} \leq \Upsilon^n$  **if and only if**  $0 \leq \mu \leq 1$ .

If the condition  $0 \leq \mu \leq 1$  is **not** satisfied, then **the solution is not bounded** and it grows with  $n$ .

If we let  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$  with  $\mu = \text{const.}$ , it only makes things worse, because then  $n \rightarrow \infty$ .

In practice, if the condition  $0 \leq \mu \leq 1$  is not satisfied, the FDE **blows up** in a few time steps.

★ ★ ★

The **problem with extrapolation** is that the maximum absolute value of the solution  $U_j^n$  increases with each time step.

Taking absolute values of the FDE we get

$$|U_j^{n+1}| \leq |U_j^n| |1 - \mu| + |U_{j-1}^n| |\mu|$$

Now defining  $\Upsilon^n = \max_j |U_j^n|$ , we have

$$\Upsilon^{n+1} \leq \{|1 - \mu| + |\mu|\} \Upsilon^n$$

and  $\Upsilon^{n+1} \leq \Upsilon^n$  **if and only if**  $0 \leq \mu \leq 1$ .

If the condition  $0 \leq \mu \leq 1$  is **not** satisfied, then **the solution is not bounded** and it grows with  $n$ .

If we let  $\Delta t \rightarrow 0$  and  $\Delta x \rightarrow 0$  with  $\mu = \text{const.}$ , it only makes things worse, because then  $n \rightarrow \infty$ .

In practice, if the condition  $0 \leq \mu \leq 1$  is not satisfied, the FDE **blows up** in a few time steps.

★ ★ ★

Do you believe me? See the following exercise.

## Practical Exercise:

Use the simple model SLAM to explore the phenomenon of computational instability.

## Practical Exercise:

Use the simple model SLAM to explore the phenomenon of computational instability.

- Choose  $\Delta t$  small and carry out a complete integration.

## Practical Exercise:

Use the simple model SLAM to explore the phenomenon of computational instability.

- Choose  $\Delta t$  small and carry out a complete integration.
- Increase  $\Delta t$  and watch the model solution “blowing up”.

## Practical Exercise:

Use the simple model SLAM to explore the phenomenon of computational instability.

- Choose  $\Delta t$  small and carry out a complete integration.
- Increase  $\Delta t$  and watch the model solution “blowing up”.
- By numerical experiment, determine approximately the maximum value of  $\Delta t$  which yields stable integrations.



## Practical Exercise:

Use the simple model SLAM to explore the phenomenon of computational instability.

- Choose  $\Delta t$  small and carry out a complete integration.
- Increase  $\Delta t$  and watch the model solution “blowing up”.
- By numerical experiment, determine approximately the maximum value of  $\Delta t$  which yields stable integrations.
- Relate this maximum to the Courant Number.  
What does it imply about the maximum phase-speed of the system?

★ ★ ★

## Practical Exercise:

Use the simple model SLAM to explore the phenomenon of computational instability.

- Choose  $\Delta t$  small and carry out a complete integration.
- Increase  $\Delta t$  and watch the model solution “blowing up”.
- By numerical experiment, determine approximately the maximum value of  $\Delta t$  which yields stable integrations.
- Relate this maximum to the Courant Number.  
What does it imply about the maximum phase-speed of the system?

★ ★ ★

Break here

# Computational Stability

We now define **computational stability**.

# Computational Stability

We now define **computational stability**.

**Definition:** An FDE is **computationally stable** if the solution of the FDE at a fixed time  $t = n\Delta t$  is bounded as  $\Delta t \rightarrow 0$ .

# Computational Stability

We now define **computational stability**.

**Definition:** An FDE is **computationally stable** if the solution of the FDE at a fixed time  $t = n\Delta t$  is bounded as  $\Delta t \rightarrow 0$ .

Note that with  $n\Delta t$  fixed,  $\Delta t \rightarrow 0$  implies  $n \rightarrow \infty$ .

# Computational Stability

We now define **computational stability**.

**Definition:** An FDE is **computationally stable** if the solution of the FDE at a fixed time  $t = n\Delta t$  is bounded as  $\Delta t \rightarrow 0$ .

Note that with  $n\Delta t$  fixed,  $\Delta t \rightarrow 0$  implies  $n \rightarrow \infty$ .

We will derive a condition for stability which involves the **Courant Number**.

# Computational Stability

We now define **computational stability**.

**Definition:** An FDE is **computationally stable** if the solution of the FDE at a fixed time  $t = n\Delta t$  is bounded as  $\Delta t \rightarrow 0$ .

Note that with  $n\Delta t$  fixed,  $\Delta t \rightarrow 0$  implies  $n \rightarrow \infty$ .

We will derive a condition for stability which involves the **Courant Number**.

The condition on the Courant number is usually known as the **Courant-Friedrichs-Lewy criterion** or simply the **CFL condition**.

★ ★ ★

# Computational Stability

We now define **computational stability**.

**Definition:** An FDE is **computationally stable** if the solution of the FDE at a fixed time  $t = n\Delta t$  is bounded as  $\Delta t \rightarrow 0$ .

Note that with  $n\Delta t$  fixed,  $\Delta t \rightarrow 0$  implies  $n \rightarrow \infty$ .

We will derive a condition for stability which involves the **Courant Number**.

The condition on the Courant number is usually known as the **Courant-Friedrichs-Lewy criterion** or simply the **CFL condition**.

★ ★ ★

Recall the story of Courant, Friedrichs and Lewy in Göttingen.



# The Lax-Richtmyer Theorem

We can now state the fundamental Lax–Richtmyer theorem:

# The Lax-Richtmyer Theorem

We can now state the fundamental Lax–Richtmyer theorem:

*Given a properly posed linear initial value problem, and a finite difference scheme that satisfies the consistency condition, then the **stability** of the FDE is the necessary and sufficient condition for **convergence**.*

$\underbrace{\text{Stability} \iff \text{Convergence}}_{\text{For consistent systems}}$

# The Lax-Richtmyer Theorem

We can now state the fundamental Lax–Richtmyer theorem:

*Given a properly posed linear initial value problem, and a finite difference scheme that satisfies the consistency condition, then the **stability** of the FDE is the necessary and sufficient condition for **convergence**.*

$$\underbrace{\text{Stability} \iff \text{Convergence}}_{\text{For consistent systems}}$$

This theorem allows us to establish **convergence** by examining the easier questions of consistency and **stability**.

# The Lax-Richtmyer Theorem

We can now state the fundamental Lax–Richtmyer theorem:

*Given a properly posed linear initial value problem, and a finite difference scheme that satisfies the consistency condition, then the **stability** of the FDE is the necessary and sufficient condition for **convergence**.*

$$\underbrace{\text{Stability} \iff \text{Convergence}}_{\text{For consistent systems}}$$

This theorem allows us to establish **convergence** by examining the easier questions of **consistency** and **stability**.

We are interested in convergence not because we want to let  $\Delta t, \Delta x \rightarrow 0$ , but because we want to make sure that the errors  $[u(j\Delta x, n\Delta t) - U_j^n]$  are acceptably small.

# The Lax-Richtmyer Theorem

We can now state the fundamental Lax–Richtmyer theorem:

*Given a properly posed linear initial value problem, and a finite difference scheme that satisfies the consistency condition, then the **stability** of the FDE is the necessary and sufficient condition for **convergence**.*

$$\underbrace{\text{Stability} \iff \text{Convergence}}_{\text{For consistent systems}}$$

This theorem allows us to establish **convergence** by examining the easier questions of **consistency** and **stability**.

We are interested in convergence not because we want to let  $\Delta t, \Delta x \rightarrow 0$ , but because we want to make sure that the errors  $[u(j\Delta x, n\Delta t) - U_j^n]$  are acceptably small.

**Definition:**  $[u(j\Delta x, n\Delta t) - U_j^n]$  is the **global truncation error**.

**Example:** We use the **criterion of the maximum** method to study the stability condition of the diffusion equation

$$\frac{\partial u}{\partial t} = \sigma \frac{\partial^2 u}{\partial x^2}$$

**Example:** We use the **criterion of the maximum** method to study the stability condition of the **diffusion equation**

$$\frac{\partial u}{\partial t} = \sigma \frac{\partial^2 u}{\partial x^2}$$

A FDE approximation (FTCS scheme) is given by

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} = \sigma \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{\Delta x^2}$$

(verification of **consistency** of this FDE is immediate).

**Example:** We use the **criterion of the maximum** method to study the stability condition of the diffusion equation

$$\frac{\partial u}{\partial t} = \sigma \frac{\partial^2 u}{\partial x^2}$$

A FDE approximation (FTCS scheme) is given by

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} = \sigma \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{\Delta x^2}$$

(verification of **consistency** of this FDE is immediate).

**Note:** Since the differences are **centered in space** but **forward in time**, the truncation error is first order in time and second order in space

$$\tau = O(\Delta t) + O(\Delta x)^2.$$



**Example:** We use the **criterion of the maximum** method to study the stability condition of the diffusion equation

$$\frac{\partial u}{\partial t} = \sigma \frac{\partial^2 u}{\partial x^2}$$

A FDE approximation (FTCS scheme) is given by

$$\frac{U_j^{n+1} - U_j^n}{\Delta t} = \sigma \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{\Delta x^2}$$

(verification of **consistency** of this FDE is immediate).

**Note:** Since the differences are **centered in space** but **forward in time**, the truncation error is first order in time and second order in space

$$\tau = O(\Delta t) + O(\Delta x)^2.$$

We can write the FDE in the form

$$U_j^{n+1} = \mu U_{j+1}^n + (1 - 2\mu)U_j^n + \mu U_{j-1}^n$$

where  $\mu = \sigma \Delta t / \Delta x^2$ .

Again,

$$U_j^{n+1} = \mu U_{j+1}^n + (1 - 2\mu)U_j^n + \mu U_{j-1}^n$$

Again,

$$U_j^{n+1} = \mu U_{j+1}^n + (1 - 2\mu)U_j^n + \mu U_{j-1}^n$$

If we take absolute values, and let  $\Upsilon^n = \max_j |U_j^n|$ , we get

$$\Upsilon^{n+1} \leq \{|\mu| + |1 - 2\mu| + |\mu|\} \Upsilon^n$$

Again,

$$U_j^{n+1} = \mu U_{j+1}^n + (1 - 2\mu)U_j^n + \mu U_{j-1}^n$$

If we take absolute values, and let  $\Upsilon^n = \max_j |U_j^n|$ , we get

$$\Upsilon^{n+1} \leq \{|\mu| + |1 - 2\mu| + |\mu|\} \Upsilon^n$$

Thus, we obtain the condition

$$0 \leq \mu \leq 1/2$$

to insure that the solution remains bounded as  $n \rightarrow \infty$ .

This is the necessary condition for stability of the FDE.

★ ★ ★

Again,

$$U_j^{n+1} = \mu U_{j+1}^n + (1 - 2\mu)U_j^n + \mu U_{j-1}^n$$

If we take absolute values, and let  $\Upsilon^n = \max_j |U_j^n|$ , we get

$$\Upsilon^{n+1} \leq \{|\mu| + |1 - 2\mu| + |\mu|\} \Upsilon^n$$

Thus, we obtain the condition

$$0 \leq \mu \leq 1/2$$

to insure that the solution remains bounded as  $n \rightarrow \infty$ .

This is the necessary condition for stability of the FDE.

★ ★ ★

Unfortunately, the **criterion of the maximum** can only be applied in very few cases.

In most FDEs some coefficients of the equations are negative, and the criterion cannot be applied.

Again,

$$U_j^{n+1} = \mu U_{j+1}^n + (1 - 2\mu)U_j^n + \mu U_{j-1}^n$$

If we take absolute values, and let  $\Upsilon^n = \max_j |U_j^n|$ , we get

$$\Upsilon^{n+1} \leq \{|\mu| + |1 - 2\mu| + |\mu|\} \Upsilon^n$$

Thus, we obtain the condition

$$0 \leq \mu \leq 1/2$$

to insure that the solution remains bounded as  $n \rightarrow \infty$ .

This is the necessary condition for stability of the FDE.

★ ★ ★

Unfortunately, the **criterion of the maximum** can only be applied in very few cases.

In most FDEs some coefficients of the equations are negative, and the criterion cannot be applied.

We need a **more powerful method** of establishing stability.

# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

We assume that we can expand the solution of the FDE in an appropriate set of eigenfunctions.



# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

We assume that we can expand the solution of the FDE in an appropriate set of eigenfunctions.

For simplicity we assume an expansion into **Fourier series**:

$$U(x, t) = \sum_k Z_k(t) e^{ikx}$$

# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

We assume that we can expand the solution of the FDE in an appropriate set of eigenfunctions.

For simplicity we assume an expansion into **Fourier series**:

$$U(x, t) = \sum_k Z_k(t) e^{ikx}$$

The space variable  $x$  and the wavenumber  $k$  can be multi-dimensional, e.g.,  $\mathbf{x} = (x_1, x_2, x_3)$ ,  $\mathbf{k} = (k_1, k_2, k_3)$  but, for simplicity, we will consider the scalar case.

# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

We assume that we can expand the solution of the FDE in an appropriate set of eigenfunctions.

For simplicity we assume an expansion into **Fourier series**:

$$U(x, t) = \sum_k Z_k(t) e^{ikx}$$

The space variable  $x$  and the wavenumber  $k$  can be multi-dimensional, e.g.,  $\mathbf{x} = (x_1, x_2, x_3)$ ,  $\mathbf{k} = (k_1, k_2, k_3)$  but, for simplicity, we will consider the scalar case.

We have  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

We assume that we can expand the solution of the FDE in an appropriate set of eigenfunctions.

For simplicity we assume an expansion into **Fourier series**:

$$U(x, t) = \sum_k Z_k(t) e^{ikx}$$

The space variable  $x$  and the wavenumber  $k$  can be multi-dimensional, e.g.,  $\mathbf{x} = (x_1, x_2, x_3)$ ,  $\mathbf{k} = (k_1, k_2, k_3)$  but, for simplicity, we will consider the scalar case.

We have  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

We define the wavenumber for the Fourier series:  $p = k\Delta x$ .

# The von Neumann Method

Another stability criterion that has much wider application is **the von Neumann stability criterion**.

We assume that we can expand the solution of the FDE in an appropriate set of eigenfunctions.

For simplicity we assume an expansion into **Fourier series**:

$$U(x, t) = \sum_k Z_k(t) e^{ikx}$$

The space variable  $x$  and the wavenumber  $k$  can be multi-dimensional, e.g.,  $\mathbf{x} = (x_1, x_2, x_3)$ ,  $\mathbf{k} = (k_1, k_2, k_3)$  but, for simplicity, we will consider the scalar case.

We have  $x_j = j\Delta x$  and  $t_n = n\Delta t$ .

We define the wavenumber for the Fourier series:  $p = k\Delta x$ .

Then the Fourier expansion is

$$U_j^n = \sum_p Z_p^n e^{ipj} \quad (\text{Note: } kx = kj\Delta x = pj)$$

When we substitute this Fourier expansion into a linear FDE, we obtain a system of equations

$$Z_p^{n+1} = \rho_p Z_p^n$$

When we substitute this Fourier expansion into a linear FDE, we obtain a system of equations

$$Z_p^{n+1} = \rho_p Z_p^n$$

Here  $\rho_p$  is an **amplification factor** that, applied to the  $p$ -th Fourier component of the solution at time  $n\Delta t$ , advances it to the time  $(n+1)\Delta t$ ;  $\rho_p$  depends on  $p, \Delta t$  and  $\Delta x$ .

When we substitute this Fourier expansion into a linear FDE, we obtain a system of equations

$$Z_p^{n+1} = \rho_p Z_p^n$$

Here  $\rho_p$  is an **amplification factor** that, applied to the  $p$ -th Fourier component of the solution at time  $n\Delta t$ , advances it to the time  $(n+1)\Delta t$ ;  $\rho_p$  depends on  $p$ ,  $\Delta t$  and  $\Delta x$ .

If we know the initial conditions

$$U_j^0 = \sum_p Z_p^0 e^{ipj}$$

then the solution of the FDE is (**remember warning about superscripts**)

$$Z_p^n = \rho_p^n Z_p^0$$



When we substitute this Fourier expansion into a linear FDE, we obtain a system of equations

$$Z_p^{n+1} = \rho_p Z_p^n$$

Here  $\rho_p$  is an **amplification factor** that, applied to the  $p$ -th Fourier component of the solution at time  $n\Delta t$ , advances it to the time  $(n+1)\Delta t$ ;  $\rho_p$  depends on  $p$ ,  $\Delta t$  and  $\Delta x$ .

If we know the initial conditions

$$U_j^0 = \sum_p Z_p^0 e^{ipj}$$

then the solution of the FDE is (**remember warning about superscripts**)

$$Z_p^n = \rho_p^n Z_p^0$$

Therefore, stability is guaranteed **if the factor  $\rho^n$  is bounded for all  $p$**  when  $\Delta t \rightarrow 0$  and  $n \rightarrow \infty$ .

When we substitute this Fourier expansion into a linear FDE, we obtain a system of equations

$$Z_p^{n+1} = \rho_p Z_p^n$$

Here  $\rho_p$  is an **amplification factor** that, applied to the  $p$ -th Fourier component of the solution at time  $n\Delta t$ , advances it to the time  $(n+1)\Delta t$ ;  $\rho_p$  depends on  $p$ ,  $\Delta t$  and  $\Delta x$ .

If we know the initial conditions

$$U_j^0 = \sum_p Z_p^0 e^{ipj}$$

then the solution of the FDE is (**remember warning about superscripts**)

$$Z_p^n = \rho_p^n Z_p^0$$

Therefore, stability is guaranteed **if the factor  $\rho^n$  is bounded for all  $p$**  when  $\Delta t \rightarrow 0$  and  $n \rightarrow \infty$ .

So, we must have  $|\rho_p|^n < M$  for all  $p$  as  $n \rightarrow \infty$ .

# Aside: Spectral Radius

For the multi-dimensional case, the modulus of  $\rho$  is replaced by the **norm of the matrix**  $G$  and the stability condition becomes  $\|G^n\| < M$  for all  $n$ , as  $n \rightarrow \infty$ .

# Aside: Spectral Radius

For the multi-dimensional case, the modulus of  $\rho$  is replaced by the **norm of the matrix**  $G$  and the stability condition becomes  $\|G^n\| < M$  for all  $n$ , as  $n \rightarrow \infty$ .

The norm  $\|G\|$  is a measure of the **size** of a matrix  $G$ .

# Aside: Spectral Radius

For the multi-dimensional case, the modulus of  $\rho$  is replaced by the **norm of the matrix**  $G$  and the stability condition becomes  $\|G^n\| < M$  for all  $n$ , as  $n \rightarrow \infty$ .

The norm  $\|G\|$  is a measure of the **size** of a matrix  $G$ .

Let  $\sigma(G)$  be the **spectral radius** of  $G$ , i.e.,  $\sigma(G) = \max_i |\lambda_i|$ , where  $\lambda_i$  are the eigenvalues of  $G$ ,

# Aside: Spectral Radius

For the multi-dimensional case, the modulus of  $\rho$  is replaced by the **norm of the matrix**  $\mathbf{G}$  and the stability condition becomes  $\|\mathbf{G}^n\| < M$  for all  $n$ , as  $n \rightarrow \infty$ .

The norm  $\|\mathbf{G}\|$  is a measure of the **size** of a matrix  $\mathbf{G}$ .

Let  $\sigma(\mathbf{G})$  be the **spectral radius** of  $\mathbf{G}$ , i.e.,  $\sigma(\mathbf{G}) = \max_i |\lambda_i|$ , where  $\lambda_i$  are the eigenvalues of  $\mathbf{G}$ ,

Then “it can be shown that”, for any norm,

$$[\sigma(\mathbf{G})]^n \leq \|\mathbf{G}^n\| \leq \|\mathbf{G}\|^n$$

# Aside: Spectral Radius

For the multi-dimensional case, the modulus of  $\rho$  is replaced by the **norm of the matrix**  $\mathbf{G}$  and the stability condition becomes  $\|\mathbf{G}^n\| < M$  for all  $p$ , as  $n \rightarrow \infty$ .

The norm  $\|\mathbf{G}\|$  is a measure of the **size** of a matrix  $\mathbf{G}$ .

Let  $\sigma(\mathbf{G})$  be the **spectral radius** of  $\mathbf{G}$ , i.e.,  $\sigma(\mathbf{G}) = \max_i |\lambda_i|$ , where  $\lambda_i$  are the eigenvalues of  $\mathbf{G}$ ,

Then “it can be shown that”, for any norm,

$$[\sigma(\mathbf{G})]^n \leq \|\mathbf{G}^n\| \leq \|\mathbf{G}\|^n$$

The equal sign is valid if  $\mathbf{G}$  is **normal**, i.e., if  $\mathbf{G}^* = \mathbf{G}^*\mathbf{G}$ , where  $\mathbf{G}^*$  is the transpose-conjugate of  $\mathbf{G}$ , but in general the amplification matrices arising from FDEs are not normal.

# Aside: Spectral Radius

For the multi-dimensional case, the modulus of  $\rho$  is replaced by the **norm of the matrix**  $\mathbf{G}$  and the stability condition becomes  $\|\mathbf{G}^n\| < M$  for all  $p$ , as  $n \rightarrow \infty$ .

The norm  $\|\mathbf{G}\|$  is a measure of the **size** of a matrix  $\mathbf{G}$ .

Let  $\sigma(\mathbf{G})$  be the **spectral radius** of  $\mathbf{G}$ , i.e.,  $\sigma(\mathbf{G}) = \max_i |\lambda_i|$ , where  $\lambda_i$  are the eigenvalues of  $\mathbf{G}$ ,

Then “it can be shown that”, for any norm,

$$[\sigma(\mathbf{G})]^n \leq \|\mathbf{G}^n\| \leq \|\mathbf{G}\|^n$$

The equal sign is valid if  $\mathbf{G}$  is **normal**, i.e., if  $\mathbf{G}^* = \mathbf{G}^*\mathbf{G}$ , where  $\mathbf{G}^*$  is the transpose-conjugate of  $\mathbf{G}$ , but in general the amplification matrices arising from FDEs are not normal.

End of digression



We found that, for stability, we must have  $|\rho|^n < M$  for all  $p$  as  $n \rightarrow \infty$ . Clearly, this requires

$$|\rho|^n \leq \exp \alpha \quad \text{for some constant } \alpha$$

We found that, for stability, we must have  $|\rho|^n < M$  for all  $p$  as  $n \rightarrow \infty$ . Clearly, this requires

$$|\rho|^n \leq \exp \alpha \quad \text{for some constant } \alpha$$

Thus, **a necessary condition for stability**, and therefore **a necessary condition for convergence**, is that

$$\lim_{\Delta t \rightarrow 0, n\Delta t \rightarrow t} |\rho|^n = \text{finite} = e^\alpha$$

We found that, for stability, we must have  $|\rho|^n < M$  for all  $p$  as  $n \rightarrow \infty$ . Clearly, this requires

$$|\rho|^n \leq \exp \alpha \quad \text{for some constant } \alpha$$

Thus, **a necessary condition for stability**, and therefore **a necessary condition for convergence**, is that

$$\lim_{\Delta t \rightarrow 0, n\Delta t \rightarrow t} |\rho|^n = \text{finite} = e^\alpha$$

Then

$$|\rho| \leq [|\rho|^n]^{1/n} \leq e^{\alpha/n} = e^{\alpha\Delta t/t} \approx 1 + \frac{\alpha\Delta t}{t}$$

or simply

$$|\rho| \leq 1 + O(\Delta t)$$

We found that, for stability, we must have  $|\rho|^n < M$  for all  $p$  as  $n \rightarrow \infty$ . Clearly, this requires

$$|\rho|^n \leq \exp \alpha \quad \text{for some constant } \alpha$$

Thus, **a necessary condition for stability**, and therefore **a necessary condition for convergence**, is that

$$\lim_{\Delta t \rightarrow 0, n\Delta t \rightarrow t} |\rho|^n = \text{finite} = e^\alpha$$

Then

$$|\rho| \leq [|\rho|^n]^{1/n} \leq e^{\alpha/n} = e^{\alpha\Delta t/t} \approx 1 + \frac{\alpha\Delta t}{t}$$

or simply

$$|\rho| \leq 1 + O(\Delta t)$$

This is the **von Neumann necessary condition** for computational stability.

**Comment:** The term  $O(\Delta t)$  allows for bounded growth which may arise from a physical instability present in the PDE.

If the exact solution grows with time, then the FDE cannot both satisfy  $|\rho| \leq 1$  and be consistent with the PDE.

**Comment:** The term  $O(\Delta t)$  allows for bounded growth which may arise from a physical instability present in the PDE.

If the exact solution grows with time, then the FDE cannot both satisfy  $|\rho| \leq 1$  and be consistent with the PDE.

**Sufficient conditions** are very complicated, and are known only for special cases.

**Comment:** The term  $O(\Delta t)$  allows for bounded growth which may arise from a physical instability present in the PDE.

If the exact solution grows with time, then the FDE cannot both satisfy  $|\rho| \leq 1$  and be consistent with the PDE.

**Sufficient conditions** are very complicated, and are known only for special cases.

In practice, we usually require  $|\rho| \leq 1$  to guarantee computational stability.

★ ★ ★

**Comment:** The term  $O(\Delta t)$  allows for bounded growth which may arise from a physical instability present in the PDE.

If the exact solution grows with time, then the FDE cannot both satisfy  $|\rho| \leq 1$  and be consistent with the PDE.

**Sufficient conditions** are very complicated, and are known only for special cases.

In practice, we usually require  $|\rho| \leq 1$  to guarantee computational stability.

★   ★   ★

For more complicated equations, the von Neumann criterion involves a matrix  $\mathbf{G}$  rather than the amplification factor  $\rho$ .

The stability criterion then involves the eigenvalues of the amplification matrix, and the von Neumann stability criterion is  $\|\mathbf{G}\| \leq 1 + O(\Delta t)$ .



# Application to Advection Equation

**PDE:**  $\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$

**FDE:**  $\frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0$  (upstream scheme)

# Application to Advection Equation

$$\text{PDE: } \frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$$

$$\text{FDE: } \frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0 \quad (\text{upstream scheme})$$

We have already studied consistency, and used the criterion of the maximum to get a sufficient condition for stability.

# Application to Advection Equation

$$\text{PDE: } \frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$$

$$\text{FDE: } \frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0 \quad (\text{upstream scheme})$$

We have already studied consistency, and used the criterion of the maximum to get a sufficient condition for stability.

Let us now apply the **von Neumann criterion**. Assume that

$$U_j^n = \sum_p Z_p^n e^{ipj} = \sum_p A_p \rho_p^n e^{ipj}$$

# Application to Advection Equation

$$\text{PDE: } \frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$$

$$\text{FDE: } \frac{U_j^{n+1} - U_j^n}{\Delta t} + c \frac{U_j^n - U_{j-1}^n}{\Delta x} = 0 \quad (\text{upstream scheme})$$

We have already studied consistency, and used the criterion of the maximum to get a sufficient condition for stability.

Let us now apply the **von Neumann criterion**. Assume that

$$U_j^n = \sum_p Z_p^n e^{ipj} = \sum_p A_p \rho_p^n e^{ipj}$$

Since the equation is linear we can consider a single term

$$U_j^n = A_p \rho_p^n e^{ipj} = A_p \rho^n e^{ipj}$$

We substitute  $U_j^n = A\rho^n e^{ipj}$  in the equation and divide by  $U_j^n$  to obtain

$$\frac{\rho - 1}{\Delta t} + c \frac{(1 - e^{-ip})}{\Delta x} = 0 \quad \text{for all } p$$

We substitute  $U_j^n = A\rho^n e^{ipj}$  in the equation and divide by  $U_j^n$  to obtain

$$\frac{\rho - 1}{\Delta t} + c \frac{(1 - e^{-ip})}{\Delta x} = 0 \quad \text{for all } p$$

The amplification factor  $\rho$  is the same as a  $1 \times 1$  amplification matrix  $G$ , and the stability condition is  $|\rho| \leq 1$  for all wavenumbers  $p$ .

We substitute  $U_j^n = A\rho^n e^{ipj}$  in the equation and divide by  $U_j^n$  to obtain

$$\frac{\rho - 1}{\Delta t} + c \frac{(1 - e^{-ip})}{\Delta x} = 0 \quad \text{for all } p$$

The amplification factor  $\rho$  is the same as a  $1 \times 1$  amplification matrix  $G$ , and the stability condition is  $|\rho| \leq 1$  for all wavenumbers  $p$ .

We need to estimate the **maximum value of  $\rho$** .

$$\rho = 1 - \mu(1 - e^{-ip}) = 1 - \mu(1 - \cos p + i \sin p)$$

Then the modulus squared is just

$$|\rho|^2 = [1 - \mu(1 - \cos p)]^2 + \mu^2 \sin^2 p$$

To repeat,

$$|\rho|^2 = [1 - \mu(1 - \cos p)]^2 + \mu^2 \sin^2 p$$



To repeat,

$$|\rho|^2 = [1 - \mu(1 - \cos p)]^2 + \mu^2 \sin^2 p$$

We make use of the trigonometrical relationships

$$\cos p = \cos^2 \frac{p}{2} - \sin^2 \frac{p}{2} = c^2 - s^2 \qquad \sin p = 2 \sin \frac{p}{2} \cos \frac{p}{2} = 2sc$$

To repeat,

$$|\rho|^2 = [1 - \mu(1 - \cos p)]^2 + \mu^2 \sin^2 p$$

We make use of the trigonometrical relationships

$$\cos p = \cos^2 \frac{p}{2} - \sin^2 \frac{p}{2} = c^2 - s^2 \quad \sin p = 2 \sin \frac{p}{2} \cos \frac{p}{2} = 2sc$$

Substituting these into  $|\rho|^2$  we have

$$\begin{aligned} |\rho|^2 &= [1 - \mu(1 - c^2 + s^2)]^2 + 4\mu^2 s^2 c^2 \\ &= [1 - 2\mu s^2]^2 + 4\mu^2 s^2 (1 - s^2) \\ &= [1 - 4\mu s^2 + 4\mu^2 s^4] + 4\mu^2 s^2 - 4\mu^2 s^4 \\ &= 1 - 4\mu s^2 + 4\mu^2 s^2 \\ &= 1 - 4\mu(1 - \mu)s^2 \end{aligned}$$

To repeat,

$$|\rho|^2 = [1 - \mu(1 - \cos p)]^2 + \mu^2 \sin^2 p$$

We make use of the trigonometrical relationships

$$\cos p = \cos^2 \frac{p}{2} - \sin^2 \frac{p}{2} = c^2 - s^2 \quad \sin p = 2 \sin \frac{p}{2} \cos \frac{p}{2} = 2sc$$

Substituting these into  $|\rho|^2$  we have

$$\begin{aligned} |\rho|^2 &= [1 - \mu(1 - c^2 + s^2)]^2 + 4\mu^2 s^2 c^2 \\ &= [1 - 2\mu s^2]^2 + 4\mu^2 s^2 (1 - s^2) \\ &= [1 - 4\mu s^2 + 4\mu^2 s^4] + 4\mu^2 s^2 - 4\mu^2 s^4 \\ &= 1 - 4\mu s^2 + 4\mu^2 s^2 \\ &= 1 - 4\mu(1 - \mu)s^2 \end{aligned}$$

Thus we obtain

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

**First**, consider the  $\sin^2 p/2$  term: **The shortest wave** that can be present in the finite difference solution is  $L = 2\Delta x$ .

To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

**First**, consider the  $\sin^2 p/2$  term: **The shortest wave** that can be present in the finite difference solution is  $L = 2\Delta x$ .

Therefore the maximum value that  $p = k\Delta x = 2\pi\Delta x/L$  can take is  $p = \pi$ , and the maximum value of  $\sin^2 p/2$  is 1.

To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

**First**, consider the  $\sin^2 p/2$  term: **The shortest wave** that can be present in the finite difference solution is  $L = 2\Delta x$ .

Therefore the maximum value that  $p = k\Delta x = 2\pi\Delta x/L$  can take is  $p = \pi$ , and the maximum value of  $\sin^2 p/2$  is 1.

**Second**, consider the factor  $4\mu(1 - \mu)$ . This is a parabola, whose maximum value is 1 when  $\mu = 0.5$ .

To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

**First**, consider the  $\sin^2 p/2$  term: **The shortest wave** that can be present in the finite difference solution is  $L = 2\Delta x$ .

Therefore the maximum value that  $p = k\Delta x = 2\pi\Delta x/L$  can take is  $p = \pi$ , and the maximum value of  $\sin^2 p/2$  is 1.

**Second**, consider the factor  $4\mu(1 - \mu)$ . This is a parabola, whose maximum value is 1 when  $\mu = 0.5$ .

So the **von Neumann condition** is satisfied provided

$$0 \leq \mu \leq 1.$$



To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

**First**, consider the  $\sin^2 p/2$  term: **The shortest wave** that can be present in the finite difference solution is  $L = 2\Delta x$ .

Therefore the maximum value that  $p = k\Delta x = 2\pi\Delta x/L$  can take is  $p = \pi$ , and the maximum value of  $\sin^2 p/2$  is 1.

**Second**, consider the factor  $4\mu(1 - \mu)$ . This is a parabola, whose maximum value is 1 when  $\mu = 0.5$ .

So the **von Neumann condition** is satisfied provided

$$0 \leq \mu \leq 1.$$

This coincides with the criterion of the maximum result.

To repeat,

$$|\rho|^2 = 1 - 4\mu(1 - \mu) \sin^2 \frac{p}{2}$$

**First**, consider the  $\sin^2 p/2$  term: **The shortest wave** that can be present in the finite difference solution is  $L = 2\Delta x$ .

Therefore the maximum value that  $p = k\Delta x = 2\pi\Delta x/L$  can take is  $p = \pi$ , and the maximum value of  $\sin^2 p/2$  is 1.

**Second**, consider the factor  $4\mu(1 - \mu)$ . This is a parabola, whose maximum value is 1 when  $\mu = 0.5$ .

So the **von Neumann condition** is satisfied provided

$$0 \leq \mu \leq 1.$$

This coincides with the criterion of the maximum result.

It is also consistent with the idea that **we should not extrapolate** but always interpolate to get the new values.

# Damping Effects of Scheme

The amplification factor  $\rho$  indicates how much the amplitude of each wavenumber will decrease or increase each time step.

# Damping Effects of Scheme

The amplification factor  $\rho$  indicates how much the amplitude of each wavenumber will decrease or increase each time step.

The upstream scheme decreases the amplitude of all Fourier wave components of the solution, since  $0 < \mu < 1 \implies |\rho| < 1$ .

It is therefore a very dissipative FDE: it has **strong numerical diffusion**.

# Damping Effects of Scheme

The amplification factor  $\rho$  indicates how much the amplitude of each wavenumber will decrease or increase each time step.

The upstream scheme decreases the amplitude of all Fourier wave components of the solution, since  $0 < \mu < 1 \implies |\rho| < 1$ .

It is therefore a very dissipative FDE: it has **strong numerical diffusion**.

The figure below shows the decrease in amplitude when using the upstream scheme after one time step and after 100 time steps (the Courant number is  $\mu = 0.5$ ).

# Damping Effects of Scheme

The amplification factor  $\rho$  indicates how much the amplitude of each wavenumber will decrease or increase each time step.

The upstream scheme decreases the amplitude of all Fourier wave components of the solution, since  $0 < \mu < 1 \implies |\rho| < 1$ .

It is therefore a very dissipative FDE: it has **strong numerical diffusion**.

The figure below shows the decrease in amplitude when using the upstream scheme after one time step and after 100 time steps (the Courant number is  $\mu = 0.5$ ).

Since the truncation error is large, the upstream scheme is in general not recommended except for special situations (e.g., for outflow boundary conditions).

# Damping Effects of Scheme

The amplification factor  $\rho$  indicates how much the amplitude of each wavenumber will decrease or increase each time step.

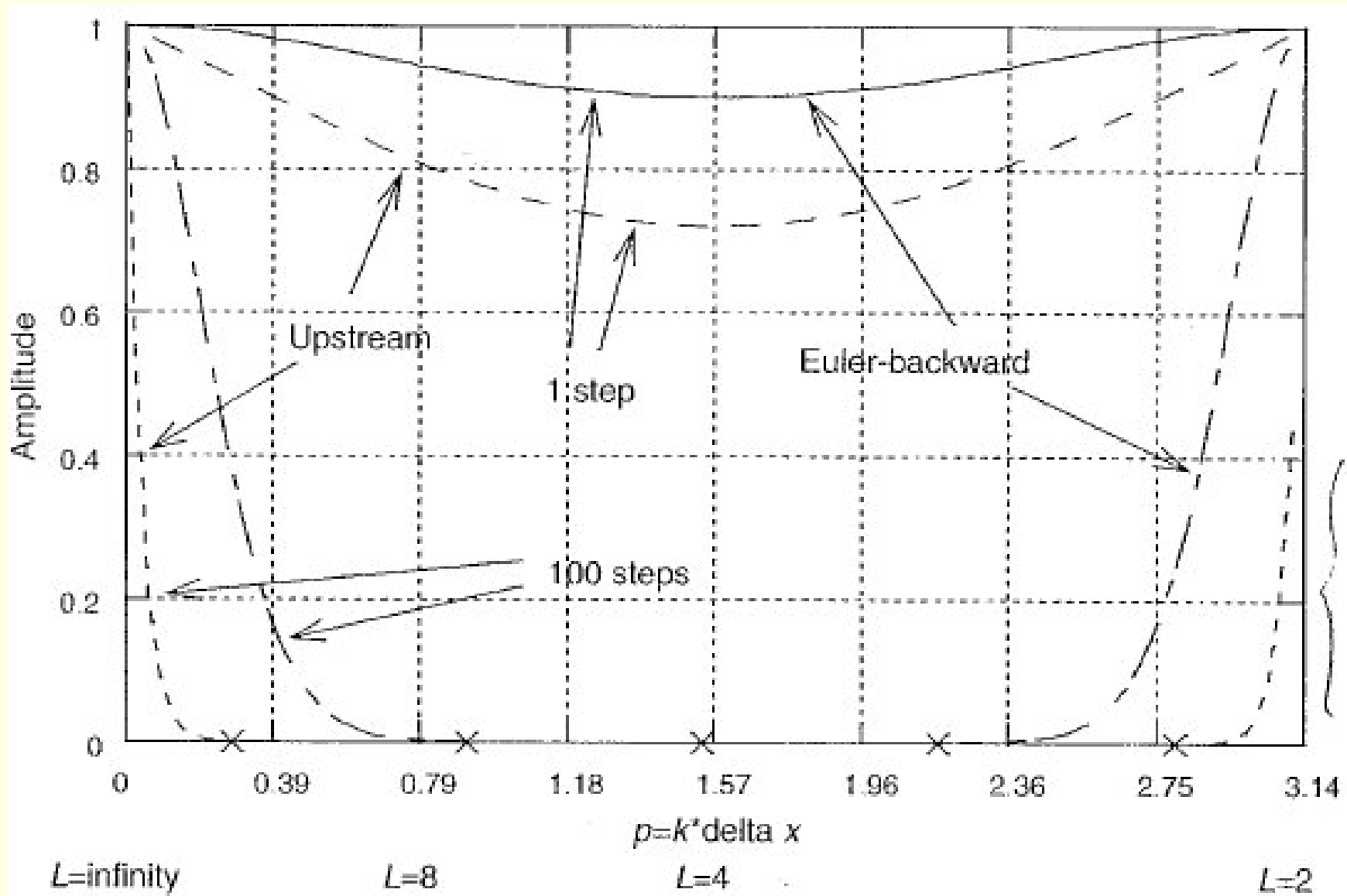
The upstream scheme decreases the amplitude of all Fourier wave components of the solution, since  $0 < \mu < 1 \implies |\rho| < 1$ .

It is therefore a very dissipative FDE: it has **strong numerical diffusion**.

The figure below shows the decrease in amplitude when using the upstream scheme after one time step and after 100 time steps (the Courant number is  $\mu = 0.5$ ).

Since the truncation error is large, the upstream scheme is in general not recommended except for special situations (e.g., for outflow boundary conditions).

An alternative, less damping scheme known as the **Matsuno** or **Euler-backward scheme** is also shown.



Amplification factor for the **upstream scheme** and the **Matsuno scheme**, with Courant Number  $\mu = 0.5$ . Response for 1 step and 100 steps shown.  $L$  is the wavelength in units of  $\Delta x$ .