

# ADVANCED NUMERICAL ANALYSIS

Part 1

FDM for the solution of boundary value problems for the second order ODE.

$$y'' = p(x)y + q(x) \text{ on } (a, b) \quad (1)$$

$$y(a) = \alpha, \quad y(b) = \beta \quad (2)$$

or

$$y'(a) - \alpha_0 y(a) = \alpha_1, \quad y'(b) + \beta_0 y(b) = \beta_1 \quad (3)$$

Assume that  $p(x), q(x) \in C^2[a, b]$ , and  $p(x) > 0$ ;  $\alpha_0, \beta_0 > 0$ ,  $\alpha_0^2 + \beta_0^2 > 0$ .

The solution (1), (2) and (1), (3)  $\in C^4[a, b]$ .

Let us divide  $[a, b]$  into  $n$  subintervals by the equally spaced points

$$x_i = a + ih \quad (i=0, 1, 2, \dots, n; h = \frac{b-a}{n})$$

~~The points  $x_i$  are the grid.~~

We call the points  $x_i$  the grid.

If we replace in (1) the second order derivative  $y''$

at  $x_i$  by

$$y''(x_i) \approx \frac{y_{i+1} - 2y_i + y_{i-1}}{h^2},$$

and  $y'(a)$  and  $y'(b)$  in (3) by

$$\frac{-y_2 + 4y_1 - 3y_0}{2h}$$

and

$$\frac{3y_n - 4y_{n-1} + y_{n-2}}{2h}$$

From these two formulas, we have

For the numerical solution (1), (2), we have

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} - p_i y_i = q_i \quad (i=1, 2, \dots, n-1) \quad (1')$$

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} - p_i y_i = q_i \quad (i=1, 2, \dots, n-1) \quad (2')$$

$$y_0 = \alpha, \quad y_n = \beta,$$

and for (1), (3)

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} - p_i y_i = q_i \quad (i=1, 2, \dots, n-1), \quad (1')$$

$$\frac{-y_2 + 4y_1 - 3y_0}{2h} - \alpha_0 y_0 = \alpha_1, \quad \frac{3y_n - 4y_{n-1} + y_{n-2}}{2h} + \beta_0 y_n = \beta_1. \quad (3')$$

- 2 -

We prove solvability of these systems. To do this we consider the following lemma:

Lemma 1. If for the system of numbers  $y_0, y_1, \dots, y_n$  ( $y_i \neq c$ )

$$L(y_i) = \frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} - p_i y_i \geq 0 \quad (i=1, 2, \dots, n-1) \quad (4)$$

then the maximum positive number among  $y_i$  can be only  $y_0$  or  $y_n$ .

Proof. If  $y_c = M$  is the maximal positive number from  $y_0, y_1, \dots, y_n$ , such that at least one of the numbers  $y_{k-1}$  or  $y_{k+1}$  less than  $M$ . By assumption

$$\frac{M - 2y_k + y_{k-1}}{h^2} - p_k y_k \geq 0.$$

If we replace  $y_{k-1}$  and  $y_{k+1}$  by  $M$ , then we increase the left hand side. Therefore,

$$\frac{M - 2M + M}{h^2} - p_k M = -p_k M > 0$$

But it's impossible, since  $p_k > 0$  and  $M > 0$ .

Similarly we prove the next

Lemma 2. If the system of numbers  $y_0, y_1, \dots, y_n$  ( $y_i \neq c$ ),  $L(y_i) \leq 0$  for all  $i=1, 2, \dots, n-1$ , then the minimum negative number among them can be only  $y_0$  or  $y_n$ .

Let us prove the solvability of (1'), (2'). For this it is enough to show that the corresponding homogeneous system

$$\frac{z_{i+1} - 2z_i + z_{i-1}}{h^2} - p_i z_i = 0 \quad (i=1, 2, \dots, n-1), \quad \} \quad (5)$$

$z_0 = 0, z_n = 0$   
has only trivial solution.

→ 3-

If the system (5) has nontrivial solution, then it should be among  $z_1, z_2, \dots, z_n$  must exist maximal positive or minimal negative number.

But this contradicts to lemmas 1 or 2, i.e.,

$\ell(z_i) = 0$  for all  $i = 1, 2, \dots, n-1$ .

Now we prove the solvability of the system (1')(3').

We consider the inhomogeneous system

$$\frac{z_{i+1} - 2z_i + z_{i-1}}{h^2} - p_i z_i = 0 \quad (i=1, 2, \dots, n-1), \quad (6)$$

$$-\frac{z_2 + 4z_1 - 3z_0}{2h} - d_0 z_0 = 0, \quad \frac{3z_n - 4z_{n-1} + z_{n-2}}{2h} + \beta_0 z_n = 0.$$

We eliminate  $z_2$  from eqs

$$\frac{z_2 - 2z_1 + z_0}{h^2} - p_1 z_1 = 0 \quad \text{and} \quad \frac{-z_2 + 4z_1 - 3z_0}{2h} - d_0 z_0 = 0$$

we have

$$\begin{aligned} z_2 &= h^2 p_1 z_1 + 2z_1 - z_0 \\ z_1 &= -2h d_0 z_0 + 4z_1 - 3z_0 \end{aligned} \quad \left\{ \begin{array}{l} h^2 p_1 z_1 + 2z_1 - z_0 = -2h d_0 z_0 + 4z_1 - 3z_0 \\ z_1 = \frac{-2z_0 - 2h d_0}{h^2 p_1 - 2} z_0 \\ = \frac{2 + 2h d_0}{2 - h^2 p_1} z_0 \end{array} \right.$$

$$z_1 = \frac{1 + d_0 h}{1 - \frac{1}{2} p_1 h^2} z_0 \quad (7)$$

$$z_{n-1} = \frac{1 + \beta_0 h}{1 - \frac{1}{2} p_{n-1} h^2} z_n \quad (8)$$

If  $z_1 \neq 0$  of (6). Since

$\ell(z_i) = 0$  ( $i = 1, 2, \dots, n-1$ ), then

max. <sup>positive</sup> or min. <sup>negative</sup> can be only  $z_0$  or  $z_n$ . Let  $z_0$  be max. positive value. Assume that  $h$  is so small that  $\frac{1}{2} p_1 h^2 < 1$ . From (7) follows that  $z_1 > z_0$ . Since  $z_1$  can not be greater than  $z_0$ , then

$z_1 = z_0 \Rightarrow d_0 = 0$  and  $p_1 = 0 \Rightarrow$  all  $z_i$  are the same, and  $z_{n-1} = z_n$ . From (8) follows that  $\beta_0 = 0$ ,  $p_{n-1} = 0$ , but  $d_0^2 + \beta_0^2 > 0$ .

By contradiction there is no nontrivial solution which takes the maximal positive value on the boundary. Similarly there is no  $\neq$  solution which takes the negative value on the boundary.

There mean that the system (6) has only trivial solution  $\Rightarrow$  the system (1'), (2') for any  $q_i, \alpha_i, \beta_i$  has a unique solution.

We estimate an error of the approximate solution of (1), (2) obtained by this method.

Lemma 3. If we have two systems  $y_0, y_1, \dots, y_n$  and  $\bar{y}_0, \bar{y}_1, \dots, \bar{y}_n$  such that

$$l(y_i) \leq -|l(\bar{y}_i)| \quad (i=1, 2, \dots, n-1), \quad (7)$$

$$|y_0| \leq \bar{y}_0, |y_n| \leq \bar{y}_n, \quad (8)$$

then for all  $i$

$$|y_i| \leq \bar{y}_i. \quad (9)$$

Proof: From (7) follows that

$$l(y_i) + |l(\bar{y}_i)| \leq 0 \Rightarrow l(y_i) \pm l(\bar{y}_i) \leq 0$$

$$\Rightarrow l(y_i \pm \bar{y}_i) \leq 0, \quad y_0 \pm \bar{y}_0 \geq 0, \quad y_n \pm \bar{y}_n \geq 0.$$

On the basis of Lemma 2 follows that  $y_i \pm \bar{y}_i$  cannot have negative values, i.e.,  $y_i \pm \bar{y}_i \geq 0$

$$\Rightarrow |y_i| \leq \bar{y}_i.$$

-5-

We denote by  $y_i$  the exact solution of the BVP (1), (2) at the grid  $x_i$ , i.e.,  $y_i = y(x_i)$ , by  $\tilde{y}_i$  the approximate solution obtained as the solution of the system (1'), (2'), and by  $\varepsilon_i$  an error of this solution, i.e.,  $\varepsilon_i = y_i - \tilde{y}_i$ .

It is easy to show that

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} - p_i y_i = q_i + R_i,$$

where  $R_i$  is the error of approximation, when the second order derivative is approximated by the central difference formula, and it is known that

$$|R_i| \leq \frac{h^2}{12} M_4, \quad M_4 = \max_{x \in [a, b]} |y^{(4)}(x)|.$$

Then  $\varepsilon_i$  will satisfy the system of equations

$$\frac{\varepsilon_{i+1} - 2\varepsilon_i + \varepsilon_{i-1}}{h^2} - p_i \varepsilon_i = R_i \quad (i=1, 2, \dots, n-1),$$

$$\varepsilon_0 = \varepsilon_n = 0.$$

We consider a system of numbers  $\gamma_i$ , which is a solution of the system

$$\frac{\gamma_{i+1} - 2\gamma_i + \gamma_{i-1}}{h^2} - p_i \gamma_i = -\frac{h^2}{12} M_4 \quad (i=1, 2, \dots, n-1),$$

$$\gamma_0 = \gamma_n = 0.$$

(On the basis of Lemma 3, we obtain that (Show!))

$$|\varepsilon_i| \leq \gamma_i.$$

If we consider the solution  $p_i$  of the system

$$\frac{p_{i+1} - 2p_i + p_{i-1}}{h^2} = -\frac{h^2}{12} M_4 \quad (i=1, 2, \dots, n-1),$$

$$p_0 = p_n = 0,$$

Then again by Lemma 3, we have  $\gamma_i \leq \xi_i$ , therefore  
 $|\xi_i| \leq \xi_i$  show! (Home work)

The solution  $\tilde{P}_2$  can be found easily. Indeed, since the second order finite difference of  $\tilde{P}_2$  is constant, then  $\tilde{P}_2$  can be consider as a value of some polynomial of second order at  $x_i$ . This polynomial is

$$\tilde{P}_2(x) = \frac{h^2}{24} M_4 (x-a)(b-x), \quad \tilde{P}_2(a) = \tilde{P}_2(b) = 0$$

and  $\tilde{P}_2 = P_2(x_i)$

$$\begin{aligned} & \tilde{P}_2(x) = x^2 ab - x^2 ab + ax = -x^2 + (b+a)x - ab \\ & \tilde{P}_2'(x) = -2x + (b+a) \\ & \tilde{P}_2''(x) = -2 \end{aligned}$$

$$\frac{\tilde{P}_{2,i+1} - 2\tilde{P}_{2,i} + \tilde{P}_{2,i-1}}{h^2} = \tilde{P}_2'' + \frac{h^2}{2} P_2''' = 0$$

$$M_4 \frac{h^2}{24} (-2) = -\frac{h^2}{12} M_4$$

$$\max_{x \in [a,b]} P_2(x) = ?$$

$$P_2'(x) = \frac{h^2}{24} M_4 (-2x + (b+a)) = 0$$

$$P_2''(x) = -\frac{h^2}{12} M_4 < 0$$

So

$$\max_{x \in [a,b]} P_2(x) = P_2\left(\frac{a+b}{2}\right) =$$

$$= \frac{h^2}{24} M_4 \left(\frac{a+b}{2} - a\right) \left(b - \frac{a+b}{2}\right)$$

$$= \frac{h^2}{24} M_4 \cdot \frac{b-a}{2} \cdot \frac{b-a}{2} = \frac{h^2 (b-a)^2}{96} M_4$$

$$\text{So } |\xi_i| \leq \frac{h^2 (b-a)^2}{96} M_4$$

$$\Rightarrow \xi_i \rightarrow 0 \text{ as } h \rightarrow 0.$$

$$\ell(\varepsilon_i) = R_i, \quad \varepsilon_0 = \varepsilon_n = 0$$

$$\ell(y_i) = -\frac{h^2}{12} M_4, \quad y_0 = y_n = 0$$

$$\text{By } \cancel{\text{[Lagrange]}} \quad \ell(y_i) = -\frac{h^2}{12} M \leq -|\ell(\varepsilon_i)| = -|R_i|$$

$$|\ell(\varepsilon_i)| = |R_i| \\ |R_i| \leq |\ell(\varepsilon_i)| \leq \frac{h^2}{12} M_4$$

$$-\frac{h^2}{12} M_4 \leq \ell(\varepsilon_i) \leq \frac{h^2}{12} M_4$$

$$|R_i| \leq \frac{h^2}{12} M_4 \quad -|R_i| \geq -\frac{h^2}{12} M_4 \\ -\frac{h^2}{12} M_4 \leq |R_i| = -|\ell(\varepsilon_i)|$$

$$\ell(y_i) = -\frac{h^2}{12} M \leq -|\ell(\varepsilon_i)|$$

$$\ell(y_i) \leq -|\ell(\varepsilon_i)|$$

$$\ell(y_i \pm \varepsilon_i) \leq 0, \quad y_0 \pm \varepsilon_0 = 0$$

$$y_i \pm \varepsilon_i \geq 0$$

$$|\varepsilon_i| \leq y_i$$

$$\frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} - p_i f_i = -\frac{h^2}{12} M_4 - p_i f_i = -\left(\frac{h^2}{12} M_4 + p_i f_i\right)$$

$$\ell(f_i + g_i) = \ell(f_i) + \ell(g_i) = -\left(\frac{h^2}{12} M_4 + p_i f_i\right) + \left(-\frac{h^2}{12} M_4\right)$$

$$-\frac{h^2}{6} - p_i f_i \leq 0$$

$$-\left(\frac{h^2}{12} M_4 + p_i f_i\right) - \left(-\frac{h^2}{12} M_4\right)$$

$$\geq -p_i f_i \leq 0$$

-8-

$$y'' - p(x)y = q(x) \quad \text{on } (a, b) \quad (1)$$

$$\begin{aligned} y'(a) - \alpha_0 y(a) &= \alpha, \\ y'(b) + \beta_0 y(b) &= \beta, \\ p(x) > 0, \quad p(x), q(x) &\in C^2[a, b], \quad \alpha_0^2 + \beta_0^2 > 0 \end{aligned} \quad (2)$$

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} - p_i y_i = q_i, \quad i = 1, 2, \dots, n-1 \quad (1')$$

$$-\frac{y_2 + 4y_1 - 3y_0}{2h} - \alpha_0 y_0 = \alpha, \quad \frac{3y_n - 4y_{n-1} + y_{n-2}}{2h} + \beta_0 y_n = \beta, \quad (2')$$

$$\frac{z_{i+1} - 2z_i + z_{i-1}}{h^2} - p_i z_i = 0, \quad i = 1, 2, \dots, n-1. \quad (*)$$

$$-\frac{z_2 + 4z_1 - 3z_0}{2h} - \alpha_0 z_0 = 0, \quad \frac{3z_n - 4z_{n-1} + z_{n-2}}{2h} + \beta_0 z_n = 0$$

$$i=1: \quad \frac{z_2 - 2z_1 + z_0}{h^2} - p_1 z_1 = 0, \quad -\frac{z_2 + 4z_1 - 3z_0}{2h} - \alpha_0 z_0 = 0$$

$$\Rightarrow \begin{cases} z_2 = h^2 p_1 z_1 + 2z_1 - z_0 \\ z_2 = -2h\alpha_0 z_0 + 4z_1 - 3z_0 \end{cases} \Rightarrow \begin{cases} h^2 p_1 z_1 + 2z_1 - z_0 = -2h\alpha_0 z_0 + 4z_1 - 3z_0 \\ (h^2 p_1 - 2)z_1 = (-2 - 2h\alpha_0)z_0 \end{cases}$$

$$z_1 = \frac{2h\alpha_0 h}{2 - P_1 h^2} z_0$$

Let  $z_i \neq 0$  of (\*). Since,  
 $\ell(z_i) = 0 \quad i = 1, 2, \dots, n-1$ , then  
max. positive or min. negative  
can be only  $z_0$  or  $z_n$ .

Let  $z_0$  be max. positive value.

Assume that  $h$  is so small that  $\frac{1}{2} P_1 h^2 < 1$ . From (3)  
follows that  $z_1 \geq z_0$ . Since  $z_1$  can not be greater than  $z_0$ ,  
then  $z_1 = z_0 \Rightarrow \alpha_0 = 0$ , and  $P_1 z_0 = 0 \Rightarrow$  all  $z_i$  are the same,  
and  $z_{n-1} = z_0$ . From (2) follows that  $\beta_0 = 0$ ,  $P_{n-1} = 0$ ,  
but  $\alpha_0^2 + \beta_0^2 > 0$ .

$$z_1 = \frac{1 + \alpha_0 h}{1 - \frac{1}{2} P_1 h^2} z_0 \quad (3)$$

$$z_{n-1} = \frac{1 + \beta_0 h}{1 - \frac{1}{2} P_{n-1} h^2} z_0 \quad (4)$$

-9-

By contradiction there is no nontrivial solution which takes the maximal positive value on the boundary. Similarly, there is no nontrivial solution which takes the negative value on the boundary. These mean that the system (\*) has only trivial solution. So, the system (1'), (2') for any  $q_i, d_i, \beta_i$  has a unique solution.

## Elliptic equations

### Finite difference method

The best known elliptic equations are Poisson's equation

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = f(x, y), \quad (\nabla^2 \phi = f(x, y)) \Rightarrow \Delta \phi = f(x, y) \quad (1)$$

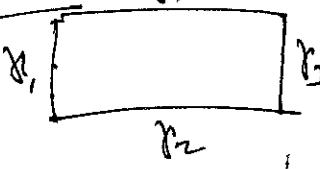
and Laplace's equation

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = 0 \Rightarrow \nabla^2 \phi = 0. \quad (\Delta \phi = 0) \quad (2)$$

#### 1) Dirichlet problem on rectangles

Let  $R$  be rectangle  $\{(x, y) : 0 < x < a, 0 < y < b\}$  and  $\gamma$  be the boundary of  $R$ ,  $R = R \cup \gamma$ .

$\gamma$  be sides of  $R$  without the end points



Dirichlet problem:

$$\Delta \phi = f(x, y) \text{ on } R \quad (3)$$

$$\phi = \varphi_i \text{ on } \gamma_k, \quad k=1, 2, 3, 4 \quad (4)$$

If is known that if  $\varphi \in C^2(\bar{R})$ ,  $u \in C^4(\bar{R})$ , then

$$\left( \frac{\partial^2 \phi}{\partial x^2} \right)_{(i,j)} \approx \frac{\phi_{i+1,j} - 2\phi_{i,j} + \phi_{i-1,j}}{h^2} \quad \begin{matrix} x = x_i = ih, & y = y_j = jh, \\ i=1, 2, \dots, M \\ j=1, 2, \dots, N \end{matrix}$$

$$\left( \frac{\partial^2 \phi}{\partial y^2} \right)_{(i,j)} \approx \frac{\phi_{i,j+1} - 2\phi_{i,j} + \phi_{i,j-1}}{h^2}$$

$$R_h = \{(x_i, y_j) : (x_{i+1}, y_{j+1}) \in \bar{R}\}$$

$$\gamma_h = \bigcup_{k=1}^4 \gamma_{k,h}$$

Then

$$\Delta_h \phi = [\phi_{i+1,j} + \phi_{i-1,j} + \phi_{i,j+1} + \phi_{i,j-1} - 4\phi_{i,j}] = f_{ij} \text{ on } R_h \quad (5)$$

is a 5-point approximation of (1)

$$\phi|_{\gamma_{k,h}} = \varphi_k, \quad k=1, 2, 3, 4 \quad (6)$$

$$\left[ \begin{array}{l} a \frac{\partial^2 \phi}{\partial x^2} + b \frac{\partial^2 \phi}{\partial xy} + c \frac{\partial^2 \phi}{\partial y^2} + d \frac{\partial \phi}{\partial x} + e \frac{\partial \phi}{\partial y} \\ b^2 - 4ac < 0 - \text{elliptic} \end{array} \right. \\ \left. \begin{array}{l} + f\phi + g = v \\ b^2 - 4ac > 0 - \text{hyperbolic} \end{array} \right. \\ \left. \begin{array}{l} b^2 - 4ac = 0 - \text{parabolic} \end{array} \right]$$

A Solution of (5), (6) exists and unique.

Since

Lemma 8. If  $\nabla \Phi_{ij} \neq \text{const.}$ , and

$$\Delta_h W_{ij} \geq 0 \quad (\leq 0) \text{ on } R_h, \text{ then}$$

positive max. (negative min) all values  $W_{ij}$  can (not) be

from Lemma 8 follows:

$\Rightarrow$  The homogeneous system

$$\Delta_h \Phi_{ij} = 0 \quad \text{on } R_h,$$

$$\Phi_{ij} = 0 \quad \text{on } \gamma_h$$

has only trivial solution

Analysis of the discretization error.

Differential problem:

$$\frac{\partial^2 U}{\partial x^2} + \frac{\partial^2 U}{\partial y^2} = f(x,y), \quad (x,y) \in R$$

$$U(x,y) = g(x,y), \quad (x,y) \in \gamma.$$

5-point approximation (FD) - problem

$$\Delta_h U_{ij} = f_{ij}, \quad (i,j) \in R_h,$$

$$U_{ij} = \varphi_{ij}, \quad (i,j) \in \gamma_h.$$

Let  $e_{ij} = U_{ij} - u_{ij}$  be discretization error at the  $(i,j)$ th mesh point  $R_h$ .

$$\text{Then } u_{ij} = U_{ij} - e_{ij},$$

and

$$\left. \begin{aligned} \Delta_h (U_{ij} - e_{ij}) &= f_{ij} \\ U_{ij} - e_{ij} &= \varphi_{ij} \end{aligned} \right\} \Rightarrow$$

$$\Delta_h e_{ij} = \Delta_h U_{ij} - f_{ij}$$

$$e_{ij} = U_{ij} - \varphi_{ij} = 0$$

$$\Rightarrow \Delta u_{ij} = T_{ij}, \quad (7)$$

where  $T_{ij}$  is the local truncation error of the difference approximation  $\Delta u - f = 0$  at the point  $(i,j)$ .

It can be shown (show!) that

$$T_{ij} = \frac{1}{12} h^2 \left\{ \left( \frac{\partial^4 U}{\partial x^4} \right)_{(i,j)} + \left( \frac{\partial^4 U}{\partial y^4} \right)_{(x_i,y_j)} \right\},$$

where  $x_i - h < x < x_i + h, y_j - h < y < y_j + h$ .

If  $M_4 = \max \left\{ \max_R \left| \frac{\partial^4 U}{\partial x^4} \right|, \max_R \left| \frac{\partial^4 U}{\partial y^4} \right| \right\}$ ,

then

$$\max_{R_h} |T_{ij}| \leq \frac{1}{6} h^2 M_4$$

From (7), we have

$$\max_{R_h} |\Delta u_{ij}| = \max_{R_h} |T_{ij}| \leq \frac{1}{6} h^2 M_4. \quad (8)$$

Lemma If  $v$  is any function defined on the set of mesh points  $R_h = R_n \cup Y_h$  in the rectangular region  $0 \leq x \leq a, 0 \leq y \leq b$ , then

$$\max_{R_h} |v| \leq \max_{Y_h} |v| + \frac{1}{4} (a^2 + b^2) \max_{R_h} |\Delta_h v|,$$

where  $\Delta_h v_{ij} = \frac{1}{h^2} (v_{0i+j} + v_{i-1,j} + v_{ij+1} + v_{ij-1} - 4v_{ij})$ .

Proof:

Define the function  $\Phi_{ij}$  by the equation

$$\Phi_{ij} = \frac{1}{4} (x_i^2 + y_j^2) = \frac{1}{4} (i^2 + j^2) h^2, \quad (i, j) \in R_h.$$

Clearly

$$0 \leq \Phi_{ij} \leq \frac{1}{4} (a^2 + b^2) \text{ for all } (i, j) \in R_h. \quad (9)$$

It also follows that  $\forall (i, j) \in R_h$ ,

$$\begin{aligned} \Delta_h \Phi_{ij} &= \frac{1}{4} \left\{ (i+1)^2 + j^2 + (i-1)^2 + j^2 + i^2 + (j+1)^2 + i^2 + (j-1)^2 - 4i^2 - 4j^2 \right\} \\ &= \frac{1}{4} \left( \underbrace{i^2 + 2(i+1)^2 + j^2}_{-4i^2} + \underbrace{i^2 - 2(i-1)^2 + j^2}_{-4j^2} + 1 + j^2 + i^2 + j^2 + 2(j+1)^2 + i^2 + j^2 - 2(j-1)^2 \right) = 1 \end{aligned} \quad (10)$$

Now define the functions  $w^+$  and  $w^-$  by

$$w^+ = v + N\Phi \text{ and } w^- = -v + N\Phi, \quad (11)$$

where

$$N = \max_{R_h} |\Delta_h v_{ij}|.$$

~~Then~~

$$\Delta_h w^\pm = \pm \Delta_h v + N \quad (i, j) \in R_h$$

$$\Rightarrow \Delta_h w_{ij}^\pm \geq 0$$

(If for any function  $\Delta_h w \geq 0$  on  $R_h$ , then  $\max_{R_h} w_{ij} \leq \max_{R_h} W_{ij}$ )

Then

$$\max_{R_h} w_{ij}^\pm \leq \max_{R_h} w_{ij}^\mp = \max_{R_h} (\pm v_{ij} + N\Phi_{ij}) \quad (\text{by (11)})$$

$$\leq \max_{R_h} (\pm v_{ij}) + \frac{1}{4} (a^2 + b^2) N \quad (\text{by (10)})$$

Since  $w_{ij}^\pm = \pm v_{ij} + N\Phi_{ij}$  and  $N\Phi_{ij} \geq 0$

$$\Rightarrow \pm v_{ij} \leq w_{ij}^\pm \text{ for all } (i, j) \in R_h.$$

Hence

$$\max_{R_h} (\pm v_{ij}) \leq \max_{R_h} (\pm v_{ij}) + \frac{1}{4} (a^2 + b^2) N,$$

i.e.,

$$\max_{R_h} |v_{ij}| \leq \max_{R_h} |v_{ij}| + \frac{1}{4} (a^2 + b^2) \max_{R_h} |\Delta_h v_{ij}|.$$

-10-

Applying this Lemma to the discretization error  $\epsilon_{ij}$  gives that

$$\max_{R_h} |\epsilon_{ij}| \leq \max_{\mathcal{V}_h} |\epsilon_{ij}| + \frac{1}{4} (a^2 + b^2) \max_{R_h} |\Delta_h \epsilon_{ij}| \quad (12)$$

But  $\epsilon_{ij}=0$  on  $\mathcal{V}_h$ , because  $V_{ij} = U_{ij} = \varphi_{ij}$ ,  $(i,j) \in \mathcal{V}_h$ .

Then from (12) and (8) we obtain

$$\max_{R_h} |\epsilon_{ij}| \leq \frac{1}{24} (a^2 + b^2) h^2 M_4.$$

$\mathcal{O}(h^2)$

The Variational difference methods (the Ritz method and Bubnov-Galerkin methods) and the finite element method (FEM)

Let  $A$  be a self-adjoint positive definite linear operator in Hilbert space  $H$  equipped with an inner product  $(\cdot, \cdot)$  and let  $f$  be a given element of the space  $H$ .

The problem of minimizing the functional

$$I[u] = (Au, u) - 2(u, f) \quad (1)$$

is equivalent to the problem of solving the equation

$$Au = f. \quad (2)$$

The element  $u_0 \in H$  satisfying the equation  $Au_0 = f$

and realizing

$$\min I[u] = I[u_0]$$

is unique.

The main idea behind the Ritz method is to take into consideration a sequence of finite-dimensional spaces  $V_n$  with basis functions  $\varphi_i^{(n)}$ ,  $i=1, 2, \dots, n$ , and look for an element  $u_n \in V_n$ , minimizing the functional  $I[u]$  in the space  $V_n$ .

We take an approximate solution  $u_n$  in the form

$$u_n = \sum_{j=1}^n y_j \varphi_j \quad (3)$$

with unknown coefficients  $y_1, y_2, \dots, y_n$ . By inserting this expression in the formula for  $I[u]$  we find that

$$I[u_n] = \sum_{i,j=1}^n \alpha_{ij} y_i y_j - 2 \sum_{j=1}^n \beta_j y_j. \quad (4)$$

$$(I[u_n] = (A \sum_{j=1}^n y_j \varphi_j, \sum_{j=1}^n y_j \varphi_j) - 2 \left( \sum_{j=1}^n y_j \varphi_j, f \right))$$

where

$$\alpha_{ij} = (A\varphi_i, \varphi_j), \beta_i = (f, \varphi_i). \quad (5)$$

Since  $A = A^*$  is a self-adjoint operator, we have  $\alpha_{ij} = \alpha_{ji}$ . The functional  $I[u_n]$  is a function of  $n$  coefficients  $y_1, y_2, \dots, y_n$ . By equating the derivatives  $\frac{\partial I[u_n]}{\partial y_i}$  to zero and using the symmetry of coefficients  $\alpha_{ij} = \alpha_{ji}$ , we obtain  $n$  eq.s for determination of  $y_i$ :

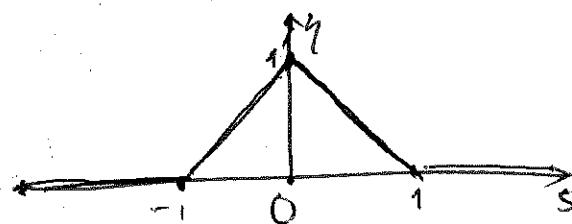
$$\sum_{j=1}^n \alpha_{ij} y_j - \beta_i = 0, \quad i = 1, 2, \dots, n. \quad (6)$$

Now let us apply the Ritz method to the problem

$$(Ku')' - qu = -f(x), \quad 0 < x < 1, \quad u(0) = 0, \quad u(1) = 0 \quad (7)$$

~~Not define~~ Let

$$\eta(s) = \begin{cases} 0, & s < -1, \quad s > 1 \\ 1+s, & -1 < s < 0 \\ 1-s, & 0 < s < 1 \end{cases} \quad (7)$$



We define

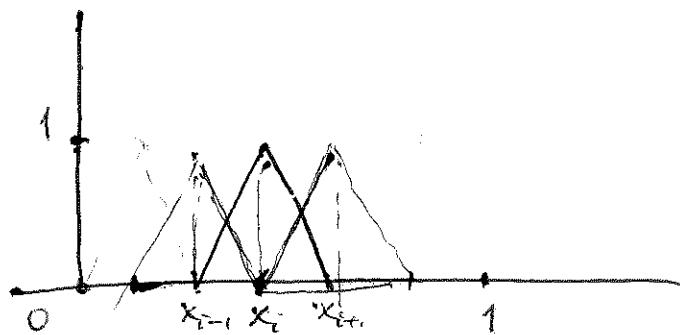
$$\varphi_i(x) = \eta\left(\frac{x-x_i}{h}\right) = \eta_i(x), \quad (8)$$

-47-

where  $x_i = ih$ ,  $i=1, 2, \dots, N-1$ , is a node of the grid  
 $\bar{\omega}_h = \{x_i = ih, i=0, 1, \dots, N, hN=1\}$ .

From (7) and (8) follows

$$\eta_i(x) = \begin{cases} 0 & \text{for } x < x_{i-1} \text{ and } x > x_{i+1}, \\ \frac{x - x_{i-1}}{h} & \text{for } x_{i-1} < x < x_i, \\ \frac{x_{i+1} - x}{h} & \text{for } x_i < x < x_{i+1}, \end{cases} \quad (7')$$



and, hence

$$\frac{d\eta_i}{dx} = \begin{cases} 0 & \text{for } x < x_{i-1} \text{ and } x > x_{i+1}, \\ \frac{1}{h} & \text{for } x_{i-1} < x < x_i, \\ -\frac{1}{h} & \text{for } x_i < x < x_{i+1}. \end{cases} \quad (8)$$

Upon substituting

$$Au = -\frac{d}{dx} \left( \kappa \frac{du}{dx} \right) + qu = f(x)$$

into (5) we get  $u(0) = u(1) = 0$

$$d_{ij} = \int_0^1 \left( \kappa \frac{d\eta_i}{dx} \frac{d\eta_j}{dx} + q\eta_i \eta_j \right) dx, \quad \beta_i = \int_0^1 f(x)\eta_i(x) dx$$

$$I(u) = \int_0^1 u \nu dx$$

$$I(u) = (Au, u) - 2(u, f)$$

$$= \int_0^1 \left( -\frac{d}{dx} \left( K \frac{du}{dx} \right) + q(u) \right) u dx - 2 \int_0^1 u f dx$$

$$\int_0^1 \left[ -\frac{d}{dx} \left( K \frac{du}{dx} \right) u dx + \int_0^1 q u^2 dx - 2 \int_0^1 u f dx \right]$$

$$= -K \frac{du}{dx} \cdot u \Big|_0^1 + \int_0^1 \left[ K \frac{du}{dx} \frac{du}{dx} + q u^2 \right] dx$$

$$- 2 \int_0^1 u f dx$$

$$\begin{aligned} dv &= -\frac{d}{dx} \left( K \frac{du}{dx} \right) dx \\ v &= -K \frac{du}{dx} \\ u &= u \\ du &= u' dx \end{aligned}$$

$$D(u) = \int_0^1 \left[ K \left( \frac{du}{dx} \right)^2 + q u^2 \right] dx - 2 \int_0^1 u f dx$$

$$U_{ij} = \sum_{j=1}^n y_j \varphi_j, \quad U_n = \sum y_j \varphi_j(x)$$

$$I(y_1, y_2, \dots, y_n) =$$

$$= \int_0^1 \left[ K \left( \sum_{j=1}^n y_j \varphi_j(x) \right)^2 + q \left( \sum_{j=1}^n y_j \varphi_j(x) \right)^2 \right] dx - 2 \int_0^1 \left( \sum_{j=1}^n y_j \varphi_j(x) \right) f dx$$

$$\frac{\partial I}{\partial y_i} = 0 \Rightarrow \int_0^1 \left[ K \cdot 2 \left( \sum_{j=1}^n y_j \frac{d \varphi_j}{dx} \right) \frac{d \varphi_i}{dx} + q_2 \left( \sum_{j=1}^n y_j \varphi_j(x) \right) \varphi_i(x) \right] dx$$

$$- 2 \int_0^1 \varphi_i(x) f(x) dx = 0$$

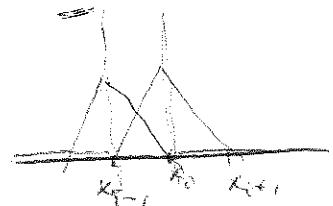
$$\begin{aligned}
 & \text{19-} \\
 & \sum_{j=1}^n \left[ \int_0^1 \left( K \frac{d\eta_i}{dx} \frac{d\eta_j}{dx} + q \eta_i(x) \eta_j(x) \right) dx \right] y_j - \underbrace{\int_0^1 f(x) \eta_i(x) dx}_{\beta_i} = 0 \\
 & \Rightarrow \boxed{\sum_{j=1}^n L_{ij} y_j - \beta_i = 0, \quad i=1,2,\dots,n} \quad (10)
 \end{aligned}$$

In light of the properties of the function  $\eta_i(x)$  and its derivatives the matrix  $\{L_{ij}\}$  is tridiagonal, because only the elements with  $j=i-1$ ,  $j=i$  and  $j=i+1$  are nonzero. Therefore from (10), we have

$$L_{i,i-1} y_{i-1} + L_{i,i} y_i + L_{i,i+1} y_{i+1} = \beta_i, \quad i=1,2,\dots,n \quad (11)$$

where

$$L_{i,i-1} = \int_0^1 \left( K \frac{d\eta_i}{dx} \frac{d\eta_{i-1}}{dx} + q \eta_i(x) \eta_{i-1}(x) \right) dx$$



$$= \int_{x_{i-1}}^{x_i} \left[ K \frac{1}{h} \left( -\frac{1}{h} \right) + q \frac{x-x_{i-1}}{h} \cdot \frac{x_i-x}{h} \right] dx$$

$$L_{i,i} = -\frac{1}{h^2} \int_{x_{i-1}}^{x_i} K(x) dx + \frac{1}{h^2} \int_{x_{i-1}}^{x_i} q(x)(x-x_{i-1})(x_i-x) dx$$

$$L_{i,i+1} = \int_{x_{i-1}}^{x_{i+1}} \left( K(x) \frac{1}{h^2} \right) dx + \frac{1}{h^2} \int_{x_{i-1}}^{x_i} q(x)(x-x_{i-1})^2 dx + \frac{1}{h^2} \int_{x_i}^{x_{i+1}} q(x)(x-x_{i+1})^2 dx$$

$$L_{i,i} = \frac{1}{h^2} \left[ \int_{x_{i-1}}^{x_{i+1}} K(x) dx + \int_{x_{i-1}}^{x_i} q(x)(x-x_{i-1})^2 dx + \int_{x_i}^{x_{i+1}} q(x)(x-x_{i+1})^2 dx \right]$$

$$\begin{aligned}
 \alpha_{i,i+1} &= \int_0^1 \left( K(x) \frac{d\psi_i}{dx} \frac{d\psi_{i+1}}{dx} + q(x) \psi_i(x) \psi_{i+1}(x) \right) dx \\
 &= \int_{x_i}^{x_{i+1}} \left[ \left( K(x) \left( -\frac{1}{h} \right) \left( \frac{1}{h} \right) + q(x) \frac{(x_{i+1}-x)(x-x_i)}{h} \right) \right] dx \\
 \alpha_{i,i+1} = \alpha_{i+1,i} &= \frac{1}{h^2} \left[ - \int_{x_i}^{x_{i+1}} K(x) dx + \int_{x_i}^{x_{i+1}} q(x) (x_{i+1}-x)(x-x_i) dx \right] \\
 \beta_i &= \int_0^1 f(x) \psi_i(x) dx = \int_{x_{i-1}}^{x_i} f(x) \frac{x-x_{i-1}}{h} dx + \int_{x_i}^{x_{i+1}} f(x) \frac{x_{i+1}-x}{h} dx \\
 \beta_i &= \frac{1}{h} \left[ \int_{x_{i-1}}^{x_i} f(x) (x-x_{i-1}) dx + \int_{x_i}^{x_{i+1}} f(x) (x_{i+1}-x) dx \right]
 \end{aligned}$$

We denote by

$$a_i = -h\alpha_{i,i+1}, \quad h^2 d_i = h\alpha_{i,i} + h(\alpha_{i,i-1} + \alpha_{i,i+1})$$

we obtain

$$\begin{aligned}
 a_i &= \frac{1}{h} \int_{x_{i-1}}^{x_i} K(x) dx - \frac{1}{h} \int_{x_{i-1}}^{x_i} q(x) (x-x_{i-1})(x_i-x) dx \\
 h^2 d_i &= h\alpha_{i,i} + h(\alpha_{i,i-1} + \alpha_{i,i+1}) = \frac{1}{h} \int_{x_{i-1}}^{x_i} q(x) (x-x_{i-1})^2 dx + \frac{1}{h} \int_{x_i}^{x_{i+1}} q(x) (x-x_{i+1})^2 dx \\
 &\quad + \frac{1}{h} \int_{x_{i-1}}^{x_i} q(x) (x-x_{i-1})(x_i-x) dx + \frac{1}{h} \int_{x_i}^{x_{i+1}} q(x) (x_{i+1}-x)(x-x_i) dx \quad (12) \\
 &= \int_{x_{i-1}}^{x_i} q(x) (x-x_{i-1}) dx + \int_{x_i}^{x_{i+1}} q(x) (x_{i+1}-x) dx \\
 d_i &= \frac{1}{h^2} \left( \int_{x_{i-1}}^{x_i} q(x) (x-x_{i-1}) dx + \int_{x_i}^{x_{i+1}} q(x) (x_{i+1}-x) dx \right).
 \end{aligned}$$

-24'

Then the system of eqs

$$d_{i,i-1}y_{i-1} + d_{i,i}y_i + d_{i,i+1}y_{i+1} - \beta_i = 0$$

can be rewritten as

$$a_{ii}y_{i-1} - (a_i + a_{i+1} + h^2 d_i)y_i + a_{i+1}y_{i+1} + h^2 \varphi_i = 0$$

or

$$(ay_x)_x - dy = -\varphi, \quad (13')$$

where

$$\varphi_i = \frac{1}{h^2} \left( \int_{x_{i-1}}^{x_i} f(x)(x-x_{i-1})dx + \int_{x_i}^{x_{i+1}} f(x)(x_{i+1}-x)dx \right). \quad (14')$$

So, the three-point scheme (12)-(14) constructed by the Ritz method is identical with scheme obtained by FIM.

## Bubnov-Galerkin Method

In contrast to the Ritz method the Bubnov-Galerkin method applies equally well to problems, where there are no fixed sign and non-self-adjoint. The coefficients  $y_i$  of the approximate solution

$$U_n = \sum_{j=1}^n y_j \varphi_j$$

are to be determined from the orthogonality conditions for the residual  $AU_n - f$  with respect to all of the basis functions  $\varphi_i(x)$ :

$$(AU_n - f, \varphi_i) = 0, \quad i=1, 2, \dots, n \quad (12)$$

We apply this method for a non-self-adjoint boundary-value problem in the form:

$$\frac{d}{dx} \left( K(x) \frac{du}{dx} \right) + r(x) \frac{du}{dx} - q(x)u = -f(x), \quad 0 < x < 1$$

$$u(0) = u(1) = 0 \quad (13)$$

$$K(x) > 0, \quad q(x) \geq 0$$

We introduce the grid  $\bar{\omega}_n = \{x_i = ih, i=0, 1, \dots, n, hN=1\}$ .

Then the dimension  $n$  of the space  $V_n$  equals  $N-1$ . The functions

$$\varphi_i(x) = \varphi_i \left( \frac{x-x_i}{h} \right), \quad i=1, 2, \dots, N-1,$$

where the function  $\varphi(s)$  was specified by (7).

$$\varphi(s) = \begin{cases} 0, & s < -1, s > 0 \\ 1+s, & -1 \leq s < 0, \\ 1-s, & 0 < s \leq 1 \end{cases}$$

$$23' - \int_0^1 \left( K(x) \frac{dy}{dx} \right) \eta_i(x) + r(x) \frac{dy}{dx} \cdot \eta_i(x) - q(x) u \eta_i(x) + f(x) \eta_i(x) dx$$

$$\int_0^1 K(x) \frac{dy}{dx} \eta_i(x) dx =$$

$$dV = \frac{d}{dx} \left( K \frac{dy}{dx} \right) dx$$

$$= K \frac{dy}{dx} \eta_i(x) \Big|_0^1 - \int_0^1 K(x) \frac{dy}{dx} \frac{d\eta_i}{dx} dx$$

$$\sum_{j=1}^{N-1} \int_0^1 \left( K(x) \frac{dy_j}{dx} \frac{d\eta_i}{dx} - r(x) \frac{d\eta_j}{dx} \eta_i(x) + q(x) y_j \eta_i(x) + f(x) \eta_i(x) \right) dx$$

$\alpha_{ij} \cdot y_j$

$$\sum_{j=1}^{N-1} \alpha_{ij} y_j - \beta_i = 0 \quad , \quad i=1, 2, \dots, N-1$$

?

where

$$\alpha_{ij} = \int_0^1 \left( K(x) \frac{dy_j}{dx} \frac{d\eta_i}{dx} - r(x) \frac{d\eta_j}{dx} \eta_i(x) + q(x) y_j \eta_i(x) \right) dx,$$

$$\beta_i = \int_0^1 f(x) \eta_i(x) dx, \quad i, j = 1, 2, \dots, N-1.$$

-24-

In this context, condition (12) becomes

$$\sum_{j=1}^{N-1} \alpha_{ij} y_j - \beta_i = 0, \quad i=1, 2, \dots, N-1, \quad (14)$$

where

$$\alpha_{ij} = \int_0^1 \left( K(x) \frac{dy_i}{dx} \frac{dy_j}{dx} - r(x) \frac{dy_i}{dx} y_j(x) + q(x) y_i(x) y_j(x) \right) dx, \quad (15)$$

$$\beta_i = \int_0^1 f(x) y_i(x) dx, \quad i, j = 1, 2, \dots, N-1$$

By definition (14)  $y(s)$  and  $y_t(x)$  of the function  $y_i(x)$ ,  
 the coefficients  $\alpha_{ij}$  are nonzero only for  $j=i-1, i, i+1$ .

At Home!

- 1) Show that, when  $r(x) \equiv 0$   $(B-G) \equiv Ritz$ .  
 2) Derive (obtain) the finite-difference approximation  
 by the  $(B-G)$  method, where  $K(x)$ ,  $r(x)$  and  $q(x)$   
 are constants.

Using the same notations (12) for  $a_i$  and  $d_i$  and (14') for  $\varphi_i$  and denoting by

$$\bar{b}_i = \frac{1}{h^2} \int_{x_{i-1}}^{x_i} r(x)(x-x_{i-1})dx =$$

$$x = x_i + sh \\ dx = hds$$

$$= \frac{1}{h^2} \int_0^1 r(x_i + sh)(1+s)ds$$

$$x = x_{i-1} + sh \\ x - x_{i-1} = h(1+s)$$

$$\bar{b}_i^+ = \int_{-1}^0 r(x_i + sh)(1+s)ds$$

$$x = x_{i+1} - h + sh$$

$$\bar{b}_i^+ = \frac{1}{h^2} \int_{x_{i-1}}^{x_{i+1}} r(x)(x_{i+1}-x)dx = \int_0^1 r(x_i + sh)(1-s)ds$$

$$x_{i+1} - x = h(1-s)$$

we reduce the system of eq.s (14) with the numbers  $b_i^+$  and  $\bar{b}_i^-$  to

$$\begin{aligned} & \frac{1}{h^2} [a_{i+1}(y_{i+1} - y_i) - a_i(y_i - y_{i-1})] + \frac{\bar{b}_i^-(y_i - y_{i-1})}{h} \\ & + \frac{\bar{b}_i^+(y_{i+1} - y_i)}{h} - dy_i = -\varphi_i, \quad i=1, 2, \dots, N-1 \end{aligned}$$

Thus, we arrive at the difference scheme

$$(ay_x)_x + \bar{b}_i^+ y_{i+1} + \bar{b}_i^- y_{i-1} - dy = -\varphi(x), \quad 0 < x = ih < 1,$$

$$y_0 = 0, y_N = 0.$$

whose coefficients can be find from (12), (13') and (16)

For  $r(x) \equiv 0$  this scheme is identical with scheme obtained by means of the Ritz method. In the case of constant coefficients  $K(x)$ ,  $r(x)$  and  $q(x)$

$$a_i = K - \frac{h^2}{2} q, \quad d_i = d = q, \quad \bar{b}_i^+ = \bar{b}_i^- = \frac{1}{2}, \quad b_i^- y_{i+1} + b_i^+ y_{i-1} = r y_i$$

where the coordinate functions  $\varphi_i(x) = \eta\left(\frac{x-x_i}{h}\right)$  are chosen by an approved rule as suggested before, the Ritz and Babuška-Galerkin methods coincide with the FEM.

## The elimination method

The problems we must solve take now the form

$$A_i y_{i+1} - C_i y_i + B_i y_{i+1} = -F_i, \quad i=1, 2, \dots, N-1 \quad (1)$$

$$y_0 = \alpha_1 y_1 + \mu_1, \quad y_N = \alpha_2 y_{N-1} + \mu_2,$$

where  $A_i \neq 0$  and  $B_i \neq 0$  for all  $i=1, 2, \dots, N-1$ .

Let

$$y_i = \alpha_i y_{i+1} + \beta_i y_{i+1}, \quad (2)$$

where  $\alpha_i$  and  $\beta_i$  unknown parameters.

From (2), we have  $y_{i+1} = \alpha_i y_i + \beta_i y_{i+1}$  and from (1), we obtain

$$A_i(\alpha_i y_i + \beta_i y_{i+1}) - C_i y_i + B_i y_{i+1} = -F_i$$

$$\Rightarrow (A_i \alpha_i - C_i) y_i + A_i \beta_i + B_i y_{i+1} = -F_i$$

which leads, because of (2), to

$$(A_i \alpha_i - C_i)(\alpha_i y_{i+1} + \beta_i y_{i+1}) + A_i \beta_i + B_i y_{i+1} = -F_i$$

$$[(A_i \alpha_i - C_i) \alpha_i + B_i] y_{i+1} + A_i \beta_i + (A_i \alpha_i - C_i) \beta_i y_{i+1} = -F_i$$

If the conditions

$$(A_i \alpha_i - C_i) \alpha_i + B_i = 0, \quad A_i \beta_i + (A_i \alpha_i - C_i) \beta_i y_{i+1} = 0$$

are fulfilled simultaneously, then the equation in view holds true for any  $y_i$ . Thus, assuming  $C_i - \alpha_i A_i \neq 0$ , we establish the recurrence formulae for determination of both  $\alpha_{i+1}$  and  $\beta_{i+1}$ :

$$\alpha_{i+1} = \frac{B_i}{C_i - \alpha_i A_i}, \quad i=1, 2, \dots, N-1, \quad (3)$$

$$\beta_{i+1} = \frac{A_i \beta_i + F_i}{C_i - \alpha_i A_i}, \quad i=1, 2, \dots, N-1, \quad (4)$$

Having substituted  $i=0$  into (2), we get

$$y_0 = \alpha_1 y_1 + \beta_1.$$

On the other hand,

$$y_0 = \alpha_1 y_1 + \mu_1,$$

giving

$$\alpha_1 = \alpha_i, \quad (5)$$

$$\beta_1 = \mu_i. \quad (6)$$

Therefore, by (3), (4), (5), (6), we calculate

$$\begin{matrix} \alpha_1, \alpha_2, \dots, \alpha_N; \\ \beta_1, \beta_2, \dots, \beta_N. \end{matrix}$$

At the second stage, knowing  $\alpha_i$  and  $\beta_i$ , the boundary value  $y_N$  is recovered from the system of the equations

$$\left. \begin{array}{l} y_N = \alpha_2 y_{N-1} + \mu_2, \\ y_{N-1} = \alpha_N y_N + \beta_N \end{array} \right\} \Rightarrow \begin{array}{l} y_N = \alpha_2 (\alpha_N y_N + \beta_N) + \mu_2 \\ (1 - \alpha_N \alpha_2) y_N = \mu_2 + \alpha_2 \beta_N \end{array}$$

When  $1 - \alpha_N \alpha_2 \neq 0$ , we have

$$\Rightarrow y_N = \frac{\mu_2 + \alpha_2 \beta_N}{1 - \alpha_N \alpha_2} \quad (7)$$

The computational formulae (2) and (7) constitute what is called the backward elimination path.

Therefore, we have

$$\overset{(1)}{\alpha_{i+1}} = \frac{\beta_i}{C_i - \alpha_i A_i}, \quad i=1, 2, \dots, N-1, \quad \alpha_1 = \alpha_i$$

$$\overset{(2)}{\beta_{i+1}} = \frac{A_i \beta_i + \mu_i}{C_i - \alpha_i A_i}, \quad i=1, 2, \dots, N-1, \quad \beta_1 = \mu_i$$

$$y_N = \frac{\mu_2 + \alpha_2 \beta_N}{1 - \alpha_N \alpha_2},$$

$$\overset{(3)}{y_i} = \alpha_i y_{i+1} + \beta_{i+1}, \quad i=N-1, N-2, \dots, 1, 0.$$

Now, when will satisfy the assumptions

$$C_i - \alpha_i A_i \neq 0 \text{ and } 1 - \alpha_n \alpha_2 \neq 0.$$

Sufficient conditions ~~are~~

$$|C_i| > |A_i| + |B_i|, \quad i=1, 2, \dots, n-1$$

$$\alpha_2 \leq 1, \quad \alpha = 1, 2, \quad |\alpha_1| + |\alpha_2| < 2,$$

yielding  $|\alpha_i| \leq 1$  for all  $i=1, 2, \dots, N$ .

The proof is carried out by induction.

Assuming  $|\alpha_i| \leq 1$  we will show that  $|\alpha_{i+1}| \leq 1$ .

Since  $|\alpha_i| = |\alpha_1| \leq 1$ , for all  $i=2, 3, \dots, N$ , we have

$$|C_i - \alpha_i A_i| - |B_i| \geq |C_i| - |\alpha_i||A_i| - |B_i| \geq$$

$$\geq |A_i| + |B_i| - |\alpha_i||A_i| - |B_i| = |A_i|(1 - |\alpha_i|) \geq 0$$

from which it follows that  $|C_i - \alpha_i A_i| > 0$ , hence  $B_i \neq 0$ .

$\Rightarrow$  relation  $|\alpha_i| \leq 1$ ,

$$|\alpha_{i+1}| = \frac{|B_i|}{|C_i - \alpha_i A_i|} \leq 1$$

Also  $|\alpha_{i+1}| < 1$  when  $|\alpha_i| < 1$ .

The assumption  $|\alpha_1| = |\alpha_2| < 1$  implies  $|\alpha_i| < 1$  for all  $i=1, 2, \dots, n-1$ .

For the denominator in formula (7):

$$|1 - \alpha_n \alpha_2| \geq 1 - |\alpha_n||\alpha_2| \geq 1 - |\alpha_2| > 0$$

because either  $|\alpha_2| < 1$  or  $|\alpha_2| > 1$ .

The Dirichlet difference problem in a domain of rather complicated configuration.

Let of

$\circ$ -nodes  
are called strictly inner or regular nodes  
we denote set of such points

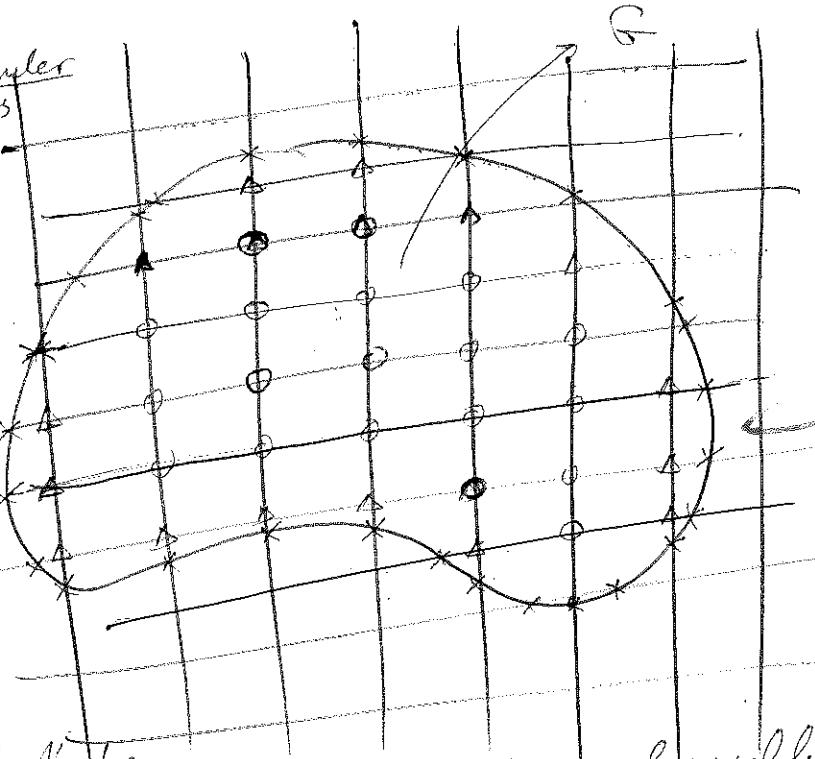
as  $C_n^0$

$\star$ -irregular nodes  
if at least one of 4 neighbor  
nodes is outside of  $G$

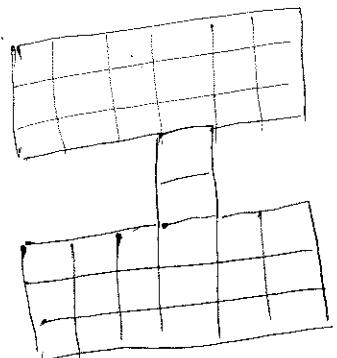
The set of such nodes we  
denote by  $W^*$   
 $\times$  intersection of the grid lines  
with the boundary of  $G$   
we denote these points as  $\Gamma_h$

Then

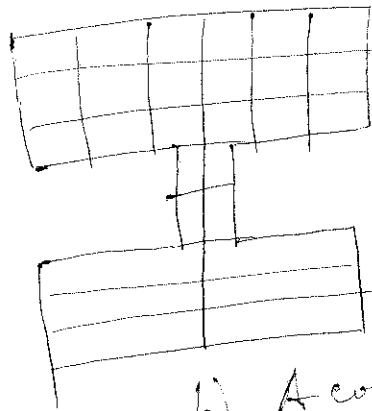
$$\bar{w}_n = C_n^0 + W^* + \Gamma_h$$



The grid  $\bar{w}_n$  is said to be  
connected if any two inner nodes can be joined by a polygonal line,  
the parts of which are parallel to the coordinate axes and  
vertices coincide with inner nodes of the grid. Then at least  
one of the four nodes  $X^{(\pm d)}$ ,  $d=1, 2$ , of the five-point  
pattern  $(X^{(\pm 1)}, X, X^{(\pm 2)})$  (regular or irregular) falls within  
the collection of inner nodes



a) A disconnected grid



b) A connected grid

The procedure of constructing a grid in the plane domain we have described above can easily be generalized to the case of an arbitrary  $P$ -dimensional domain. A grid is constructed as a result of the intersection of hyperplanes (planes for  $P=3$  or straight lines for  $P=2$ )

$$x_\alpha^{i_\alpha} = i_\alpha h_\alpha, \quad i_\alpha = 0, \pm 1, \dots, \alpha = 1, 2, \dots, P,$$

where  $h_\alpha > 0$ . The preceding classification of nodes remains unchanged here.

Our purpose here is to construct a difference scheme for solving the Dirichlet problem in the domain  $\bar{G} = G + \Gamma$ , the complete posing of which is to find an unknown solution to the equation

$$\Delta u = \frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2} = -f(x), \quad x = (x_1, x_2) \in G,$$

which is continuous in the closed domain  $\bar{G} = G + \Gamma$  and satisfies the b. c.  $u|_{\Gamma} = \mu(x)$ .

At each of the inner nodes  $x \in W_h$  we approximate the differential operator

$$L_h u = \frac{\partial^2 u}{\partial x_\alpha^2}$$

by the three-point difference operator  $\Lambda_\alpha$

If a node  $x \in W_h$  is regular with respect to  $x_\alpha$ , then the difference operator  $\Lambda_\alpha$  on the regular pattern  $(x^{(-h_\alpha)}, x, x^{(h_\alpha)})$  is

$$\Lambda_\alpha y = y_{x_\alpha x_\alpha} = \frac{y^{(h_\alpha)} - 2y + y^{(-h_\alpha)}}{h_\alpha^2} \quad (1)$$

But if a node  $x \in W_{h,\alpha}^*$ , that is, a node is irregular with respect to  $x_\alpha$  on the irregular pattern, ~~the~~ the difference operator  $\Lambda_\alpha$  can be rewritten as

$$\Lambda_\alpha^* y = \frac{1}{h_\alpha} \left( \frac{y^{(h_\alpha)} - y}{h_\alpha} - \frac{y - y^{(-h_\alpha)}}{h_\alpha^*} \right) \text{ for } x^{(-h_\alpha)} \in W_{h,\alpha}, \quad (2)$$

where  $h_\alpha^*$  is the distance between the nodes  $x$  and  $x^{(-\alpha)}$ , or

$$\Delta_\alpha^* y = \frac{1}{h_\alpha} \left( \frac{y^{(+\alpha)} - y}{h_\alpha^*} - \frac{y - y^{(-\alpha)}}{h_\alpha} \right) \text{ for } x^{(\pm\alpha)} \in \mathcal{Y}_{h,\alpha}, \quad (3)$$

where  $h_\alpha^*$  is the distance between the nodes  $x$  and  $x^{(+\alpha)}$ .

If  $x^{(-\alpha)} \in \mathcal{Y}_{h,\alpha}$  and  $x^{(+\alpha)} \in \mathcal{Y}_{h,\alpha}$ , then

$$\Delta_\alpha^* y = \frac{1}{h_\alpha} \left( \frac{y^{(+\alpha)} - y}{h_{\alpha+}^*} - \frac{y - y^{(-\alpha)}}{h_{\alpha-}} \right) \text{ for } x^{(\pm\alpha)} \in \mathcal{Y}_{h,\alpha} \quad (4)$$

where  $h_{\alpha\pm}^* \neq h_\alpha$  is the distance between  $x$  and  $x^{(\pm\alpha)}$ .

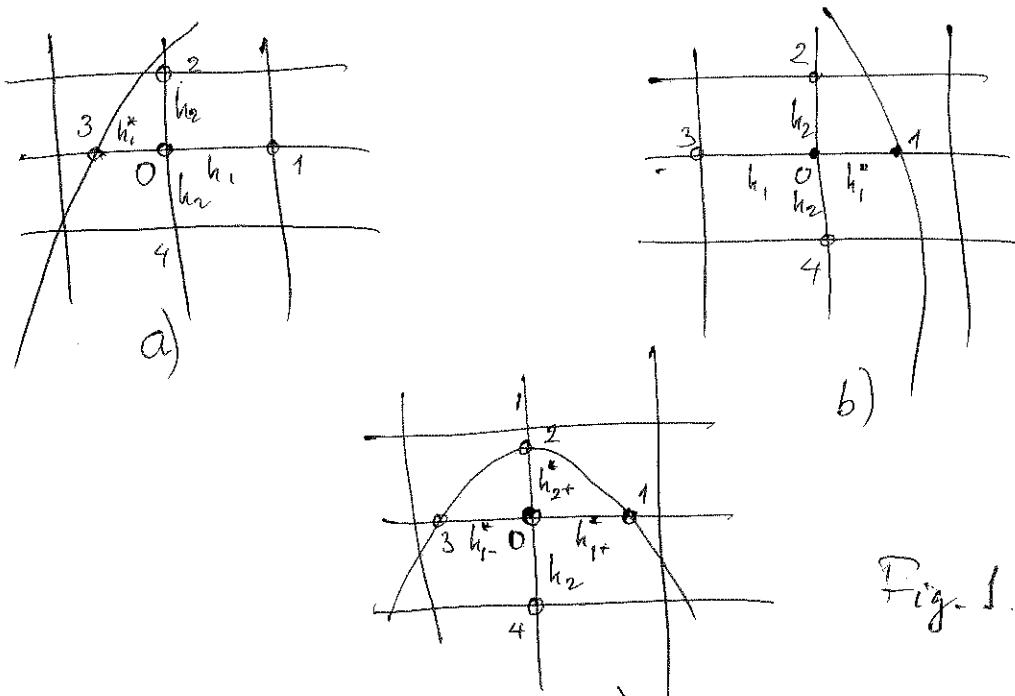


Fig. 1.

Finally, we arrive at the Dirichlet difference problem of determining a grid function  $\bar{y}(x)$  defined for  $x \in \bar{\Omega}_h = \Omega_h \cup \mathcal{Y}_h$ , satisfying at the inner nodes the equation

$$\begin{aligned} \Delta y + \varphi(x) &= 0 \quad \text{at the regular nodes,} \\ \Delta_\alpha^* y + \varphi(x) &= 0 \quad \text{at the irregular nodes} \end{aligned} \quad (5)$$

$$y = f(x), \quad x \in \mathcal{Y}_h.$$

where  $\varphi$  is the approximation error equal for  $\varphi(x) = f(x)$  to

a)  $\Lambda_1^* y = \frac{1}{h_1} \left( \frac{y_1 - y_0}{h_1^*} - \frac{y_0 - y_3}{h_1^*} \right), \Lambda_2 y = y_{\bar{x}_2 \bar{x}_2}, \Lambda = \Lambda_1^* + \Lambda_2,$

b)  $\Lambda_2^* y = \frac{1}{h_1} \left( \frac{y_1 - y_0}{h_1^*} - \frac{y_0 - y_3}{h_1^*} \right), \Lambda_2 y = y_{\bar{x}_2 \bar{x}_2}, \Lambda = \Lambda_2^* + \Lambda_2$

c)  $\Lambda_1^* y = \frac{1}{h_1} \left( \frac{y_1 - y_0}{h_{1+}^*} - \frac{y_0 - y_3}{h_{1-}^*} \right), \Lambda_2 y = \frac{1}{h_2} \left( \frac{y_2 - y_0}{h_2^*} - \frac{y_0 - y_3}{h_2^*} \right),$

$\Lambda^* y = \Lambda_1^* + \Lambda_2^*,$

$\Psi = \Lambda u + \varphi = \Lambda u - Lu$  at the regular nodes,

$\Psi^* = \Lambda^* u - Lu$  at the irregular nodes. (6)

Let  $u \in C^{(4)}(\bar{\Omega})$ , ~~where  $C^{(4)}$  is the~~ As stated before we have

$$|\Psi| \leq M_4 \frac{|h|^2}{12}, |h|^2 = h_1^2 + h_2^2 + \dots + h_p^2 \quad (7)$$

at the irregular nodes

$$\Psi^* = \sum_{\alpha=1}^p \Psi_\alpha^*, \Psi_\alpha^* = \Lambda_\alpha u - L_\alpha u,$$

$$\Psi_\alpha^* = \frac{h_\alpha + h_{\alpha-1}}{2h_\alpha} \frac{\partial^2 u}{\partial x^2} + O(h_\alpha) = O(1), \Psi^* = O(1),$$

meaning that at the irregular nodes the scheme does not approximate the equation  $\Delta u + f(x) = 0$ .

Thus, in the p-dimensional case a difference scheme such as

$$\Lambda y = \sum_{\alpha=1}^p \Lambda_\alpha y = -f(x) \text{ at the regular nodes.}$$

$$\Lambda^* y = \sum_{\alpha=1}^p \Lambda_\alpha^* y = -f(x) \text{ at the irregular nodes,}$$

where  $\Lambda_\alpha y = y_{\bar{x}_\alpha \bar{x}_\alpha}$  and  $\Lambda_\alpha^*$  is specified by the formula

$$-\Delta^* y = \frac{1}{h_x^2} \left( \frac{y^{(+x)} - y}{h_x^+} - \frac{y - y^{(-x)}}{h_x^-} \right),$$

is associated with problem (1).

$$\Delta u = \sum_{\alpha=1}^p \frac{\partial^2 u}{\partial x_\alpha^2} = -f(x), \quad x \in G,$$

$$u|_{\Gamma} = \mu(x)$$

where  $x = (x_1, x_2, \dots, x_p)$ ,  $G$  is a  $p$ -dimensional finite domain with the boundary  $\Gamma$ .

Remark. Quite often, the Dirichlet problem is approximated by the method based on the difference approximation at the near-boundary nodes of the Laplace operator on an irregular pattern, with the use of formulae with  $\frac{1}{h_x^2} ( )$  instead of  $\frac{1}{h_x} ( )$ . However, in this case <sup>the difference operator</sup> several important properties does not possess, such as the self-adjointness and the property of having fixed sign.

## 5. The canonical form of a difference equation

We now consider the  $(2p+1)$ -point scheme  $\Delta y = -f$  at a

regular node

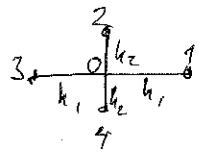
$$\sum_{\alpha=1}^p \frac{1}{h_\alpha^2} \left( y^{(+\alpha)} - 2y + y^{(-\alpha)} \right) = -f,$$

which admits an alternative form of writing

$$\sum_{\alpha=1}^p \frac{2}{h_\alpha^2} y^{(\alpha)} = \sum_{\alpha=1}^p \frac{1}{h_\alpha^2} \left( y^{(+\alpha)} + y^{(-\alpha)} \right) + f(x).$$

In two-dimensional case

$$2 \left( \frac{1}{h_1^2} + \frac{1}{h_2^2} \right) y_0 = \frac{1}{h_1^2} (y_1 + y_3) + \frac{1}{h_2^2} (y_2 + y_4) + f_0$$



- 34 -

Let  $x \in \omega_{h_1}^*$  be an irregular node. (in Fig 1 a))

$$\Lambda_1^* y_0 = \frac{1}{h_1} \left( \frac{y_1 - y_0}{h_1} - \frac{y_0 - y_3}{h_1^*} \right) = \frac{1}{h_1} \left( \frac{y_1}{h_1} + \frac{y_3}{h_1^*} - \frac{2t_1}{h_1 h_1^*} y_0 \right),$$

where  $t_1 = \frac{1}{2}(h_1 + h_1^*)$ , and

$$\Lambda_2 y_0 = \frac{1}{h_2^2} (y_2 - 2y_0 + y_4).$$

From the eq.  $\Lambda^* y = \Lambda_1^* y + \Lambda_2 y = -f$  we find that

$$\left( \frac{2t_1}{h_1^2 h_1^*} + \frac{2}{h_2^2} \right) y_0 = \frac{1}{h_1^2} y_1 + \frac{1}{h_1 h_1^*} y_3 + \frac{1}{h_2^2} (y_2 + y_4) + f_0.$$

In the case Fig 1, c), we deduce that

$$\left( \frac{2t_1}{h_1 h_1^* h_{1+}} + \frac{2t_2}{h_2^2 h_2^*} \right) y_0 = \frac{1}{h_1 h_{1+}} y_1 + \frac{1}{h_1 h_{1+}} y_3 + \frac{1}{h_2 h_2^*} y_2 + \frac{1}{h_2^2} y_4 + f_0,$$

where  $t_1 = \frac{1}{2}(h_{1+}^* + h_{1+})$  and  $t_2 = \frac{1}{2}(h_2 + h_2^*)$ .

$\Rightarrow$  all these eq.s. can be represented in the canonical form

$$A(x)y(x) = \sum_{\{s\} \in \text{Patt}'(x)} B(x, s)y(s) + F(x), \quad x \in \omega_h \quad (8)$$

where  $\text{Patt}'(x)$  is the set consisting of  $2p$  nodes of the  $(2p+1)$ -point "cross" pattern with center at the point  $x$ .

except for the node  $x$  itself, that is  $s \neq x$ .

We call the set  $\text{Patt}'(x)$  the neighborhood of the node  $x$ .

It is easily seen that

$$A(x) > 0, \quad B(x, s) > 0, \quad \sum_{s \in \text{Patt}'(x)} B(x, s) = A(x) \quad \text{for all } x \in \omega_h. \quad (9)$$

Eq. (8) is put together with the boundary conditions

$$y|_{\gamma_h} = \mu(x). \quad (10)$$

The Dirichlet difference problem is a special case of a more general problem in which it is required to find a grid function  $y(x)$  defined on the grid  $\bar{w}_h = w_h + \gamma_h$  and satisfying on  $w_h$  the equation

$$A(x)y(x) = \sum_{\{s\} \in \text{Patt}(x)} B(x,s)y(s) + F(x), \quad x \in w_h, \quad (11)$$

$$y(x) = \mu(x), \quad x \in \gamma_h,$$

where

$$A(x) > 0, \quad B(x,s) > 0, \quad D(x) = A(x) - \sum_{\{s\} \in \text{Patt}(x)} B(x,s) \geq 0 \quad (12)$$

for all  $x \in w_h$

To prove the existence and uniqueness of a solution of problem (11), (12), it suffices to check that the homogeneous eq.

$$\mathcal{L}[y] = A(x)y(x) - \sum_{\{s\} \in \text{Patt}(x)} B(x,s)y(s) = 0, \quad x \in w_h, \quad (13)$$

$$y(x) = 0, \quad x \in \gamma_h$$

has only the trivial solution  $y(x) = 0$  for  $x \in \bar{w}_h$ .

## The Maximum Principle

The maximum principle is suitable of difference elliptic and parabolic eq.s in the space  $C$  and is certainly true for grid eq.s of common structure, which will be

Let  $w$  be a finite set of nodes (a grid) in some bounded domain of the  $n$ -dimensional Euclidean space and let  $P \in w$  a point of the grid  $w$ .

Consider the eq.

$$A(P)y(P) = \sum_{Q \in \text{Patt}'(P)} B(P, Q)y(Q) + F(P), \quad P \in w \quad (1)$$

related to a function  $y(P)$  defined on the grid  $w$ . Here the coefficients of the equation  $A(P)$  and  $B(P, Q)$  and the right-hand side of the eq.  $F(P)$  are given grid functions;  $\text{Patt}'(P) \subset w$  being the set of all the nodes of the grid  $w$  except for the node  $P$ , is the neighborhood of all the nodes of the grid  $w$  except for the node  $P$ , is the neighborhood of the node  $P$ . The pattern of the grid eq. (1) at the node  $P$  consists, evidently, of the node  $P$  itself and its neighborhood  $\text{Patt}'(P)$ .

In what follows we will suppose that coefficients  $A(P)$  and

$B(P, Q)$  are subject to the conditions

$$A(P) > 0, \quad B(P, Q) > 0 \quad \text{for all } P \in w, \quad Q \in \text{Patt}'(P), \quad (2)$$

$$D(P) = A(P) - \sum_{Q \in \text{Patt}'(P)} B(P, Q) \geq 0.$$

A point  $P$  is called a boundary node of the grid  $w$  if at this point the value of the function  $y(P)$  is known in advance:

$$y(P) = \mu(P) \quad \text{for } P \in \gamma, \quad (3)$$

where  $\gamma$  is the set of all boundary nodes.

Comparing (3) and (1) we see that on the boundary  $\gamma$  we must set formally  $A(P) \equiv 1$ ,  $B(P, Q) \equiv 0$  and  $F(P) = \mu(P)$ .

The grid  $\bar{w}$  is taken to be connected, that is, for fixed points  $\bar{P} \in \bar{w}$  and  $\bar{P} \in \bar{w}$  there always exists a sequence of neighborhoods  $\{\text{Patt}'(P)\}$  such that the passage from  $\bar{P}$  to  $\bar{P}$  can be done using only the nodes of those neighborhoods or, in other words, one can select nodes  $P_1, P_2, \dots, P_m$  of the grid  $w$  such that

$$P_1 \in \text{Patt}'(\bar{P}), \quad P_2 \in \text{Patt}'(P_1), \dots, \quad P_m \in \text{Patt}'(P_{m-1}), \quad \bar{P} \in \text{Patt}'(P_m)$$

with  
 $B(P_i, P_{i+1}) \neq 0, \quad i=1, 2, \dots, m-1,$   
 $B(P_1, P_2) \neq 0, \quad B(P_m, \bar{P}) \neq 0.$

By using the notation

$$\mathcal{L}y(P) = A(P)y(P) - \sum_{Q \in \text{Patt}'(P)} B(P, Q)y(Q), \quad (5)$$

the eq. (1) can be written in the form

$$\mathcal{L}y(P) = F(P), \quad (6)$$

An alternative form of  $\mathcal{L}y(P)$  may be useful in the further development:

$$\mathcal{L}y(P) = D(P)y(P) + \sum_{Q \in \text{Patt}'(P)} B(P, Q)(y(P) - y(Q)). \quad (7)$$

Consider as one possible example the so-called scheme with weights for the heat conduction equation

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + f(x, t), \quad 0 < x < 1, \quad t > 0,$$

$$u(x, 0) = u_0(x), \quad u(0, t) = \mu_1(t), \quad u(1, t) = \mu_2(t)$$

On the grid  $\bar{\omega}_{h\tau} = \{(x_i = ih, t_j = j\tau), \quad i=0, 1, \dots, N, hN=1, j=0, 1, \dots\}$

This scheme takes the form

$$\frac{y_i^{j+1} - y_i^j}{\tau} = \lambda (5y_i^{j+1} + (1-\sigma)y_i^j) + \varphi_i^j, \quad (8)$$

$$\lambda y = y_{xx}, \quad y_i^0 = u_0(x_i), \quad y_0^j = \mu_1(t_j), \quad y_N^j = \mu_2(t_j).$$

$$\Rightarrow \left(\frac{1}{\tau} + \frac{2\sigma}{h^2}\right)y_i^{j+1} = \frac{\sigma}{h^2}(y_{i-1}^{j+1} + y_{i+1}^{j+1})$$

$$+ \left(\frac{1}{\tau} + \frac{2(\sigma-1)}{h^2}\right)y_i^j$$

$$+ \frac{1-\sigma}{h^2}(y_{i-1}^j + y_{i+1}^j) + \varphi_i^j.$$

From here it is easily seen that  $B(P, Q) \geq 0$  only if  $\tau \leq \frac{h^2}{2(1-\sigma)}$  and  $0 \leq \sigma \leq 1$ .  $D(P) = 0$ .

## The maximum principle

Theorem 1 Let  $y(P) \neq \text{const}$  be a grid function defined on a connected grid  $\bar{\omega}$  and let both conditions (2) <sup>and (4)</sup> and (4) hold.

Then the condition  $\mathcal{L}y(P) \leq 0$  ( $\mathcal{L}y(P) > 0$ ) on the grid  $\bar{\omega}_h$  implies that  $y(P)$  cannot attain the maximal positive (minimal negative) value at the inner nodes  $P \in \bar{\omega}_h$ .

Proof. Let  $\mathcal{L}y(P) \leq 0$  at all of the inner nodes  $P \in \bar{\omega}_h$ . On the contrary, let  $y(P)$  attain its maximal positive value at an inner node  $\bar{P} \in \bar{\omega}$ , so that

$$y(\bar{P}) = \max_{\bar{\omega}} y(P) = M_0 > 0.$$

Since  $y(\bar{P}) > y(Q)$  for all  $Q \in \text{Patt}'(\bar{P})$ , we find by virtue of the relation  $\mathcal{D}(\bar{P}) \geq 0$  and  $y(\bar{P}) > 0$  that

$$\mathcal{L}y(\bar{P}) = \mathcal{D}(\bar{P}) y(\bar{P}) + \sum_{Q \in \text{Patt}'(\bar{P})} \beta(\bar{P}, Q) (y(\bar{P}) - y(Q)) \geq \mathcal{D}(\bar{P}) y(\bar{P}) \geq 0.$$

Then  $\Rightarrow \mathcal{L}y(\bar{P}) = 0$  needs investigation.

We now face the node  $P \in \text{Patt}'(\bar{P})$  at which  $y(P) = y(\bar{P}) = M_0$ , and adopt these ideas. Since  $y(P) \neq \text{const}$  on the grid  $\bar{\omega}$  and the grid is connected, there exists a sequence of nodes  $P_1, P_2, \dots, P_m, \bar{P}$ , satisfying conditions (4), such that

$$y(P_m) = y(\bar{P}) = M_0, \quad y(\bar{P}) < M_0,$$

but  $\bar{P} \in \text{Patt}'(P_m)$ .

Then

$$\begin{aligned} \mathcal{L}y(P_m) &\geq \mathcal{D}(P_m) y(P_m) + \beta(P_m, \bar{P}) (y(P_m) - y(\bar{P})) \geq \\ &\geq \beta(P_m, \bar{P}) (y(\bar{P}) - y(\bar{P})) > 0 \end{aligned}$$

meaning  $\bar{P} = P_m$  and justifying the first assertion of the theorem.

The second assertion will be reduced to the first one once we replace  $y(P)$  by  $-y(P)$ .

Corollary 1. Let conditions (2) and (4) hold and let a grid function  $y(P)$  defined on  $\omega + \gamma$  be nonnegative on the boundary, that is,  $y(P) \geq 0$  for  $P \in \gamma$  and  $\mathcal{L}y(P) \geq 0$  on  $\omega$ . Then  $y(P)$  is nonnegative on  $\omega + \gamma$ , that is  $y(P) \geq 0$  for all  $P \in \omega + \gamma$ . But if  $y(P) \leq 0$  on  $\gamma$  and  $\mathcal{L}y(P) \leq 0$  on  $\omega$ , then  $y(P) \leq 0$  on  $\omega + \gamma$ .

Proof. -

Corollary 2. The homogeneous eq. (1) subject to the boundary condition

$\mathcal{L}y(P) = 0$  on  $\omega$ ,  $y(P) = 0$  on  $\gamma$ , (9)  
has only the trivial solution  $y(P) \equiv 0$ .

Proof.

Corollary 3. Problem (1)-(4) possesses a unique solution.

Comparison theorem. The majorant.

Theorem. Let  $y(P)$  be a solution of problem (1)-(4) and let  $\bar{Y}(P)$  be a solution of the problem

$$\mathcal{L}\bar{Y}(P) = \bar{F}(P), \quad P \in \omega, \quad \bar{Y}(P) = \bar{\mu}(P), \quad P \in \gamma. \quad (10)$$

Then the conditions

$$|F(P)| \leq \bar{F}(P), \quad P \in \omega, \quad |\mu(P)| \leq \bar{\mu}(P), \quad P \in \gamma, \quad (11)$$

provide the validity of the inequality

$$|y(P)| \leq \bar{Y}(P), \quad \text{for } P \in \omega + \gamma. \quad (12)$$

Proof. By Corollary 1 the inequality  $\mathcal{Y}(P) \geq 0$  is valid on  $W+\mathcal{Y}$ . The functions  $U(P) = \mathcal{Y}(P) + y(P)$  and  $V(P) = \mathcal{Y}(P) - y(P)$  solve the equations  $\mathcal{L}U = F_u = \bar{F} + F \geq 0$  and  $\mathcal{L}V = F_v = \bar{F} - F \geq 0$  subject to the boundary conditions  $U|_{\mathcal{Y}} = (y+y)|_{\mathcal{Y}} = \bar{\mu} + \mu \geq 0$  and  $V|_{\mathcal{Y}} = (y-y)|_{\mathcal{Y}} = \bar{\mu} - \mu \geq 0$ .

Since the conditions of Corollary 1 are satisfied, we have  $U > 0$  or  $y \geq -\mathcal{Y}$ ,  $V > 0$  or  $y \leq \mathcal{Y}$ . It follows from the foregoing that  $-\mathcal{Y} \leq y \leq \mathcal{Y}$  or  $|y(P)| \leq \mathcal{Y}(P)$  on  $W+\mathcal{Y}$ . The function  $\mathcal{Y}(P)$  is called the majorant for a solution of problem (1)-(3).

Corollary For a solution of the problem

$$\mathcal{L}(y) = 0 \text{ on } W, \quad y(P) = \mu(P) \text{ on } \mathcal{Y}, \quad (13)$$

The estimate

$$\max_{P \in W+\mathcal{Y}} |y(P)| = \|y\|_C \leq \|\mu\|_{C_{\mathcal{Y}}} \quad (14)$$

is valid with  $\|\mu\|_{C_{\mathcal{Y}}} = \max_{P \in \mathcal{Y}} |\mu(P)|$ .

Indeed, let us specify the majorant  $\mathcal{Y}(P)$  by the conditions  $\mathcal{L}\mathcal{Y}=0$  on  $W$  and  $\mathcal{Y} = \|\mu\|_{C_{\mathcal{Y}}}$  on  $\mathcal{Y}$ . The function  $\mathcal{Y}(P)$  is nonnegative on  $W+\mathcal{Y}$  and attains its maximum at some node of the grid. This node is none the inner nodes if  $\mathcal{Y}(P) \neq \text{const}$  and, hence,

$$\|\mathcal{Y}\|_C = \max_{P \in W+\mathcal{Y}} \mathcal{Y}(P) = \max_{P \in \mathcal{Y}} \mathcal{Y}(P) = \|\mu\|_{C_{\mathcal{Y}}}.$$

If  $\mathcal{Y}(P) = \text{const}$ , then  $\mathcal{Y}(P) = \|\mu\|_{C_{\mathcal{Y}}}$ . In both cases  $\|\mathcal{Y}\|_C = \|\mu\|_{C_{\mathcal{Y}}}$ .

Combination of this relation and the inequality  $\|y\|_C \leq \|\mathcal{Y}\|_C$  gives estimate (14).

The estimate of a solution to the nonhomogeneous equation.

In the further development a solution of problem (1)-(3) is viewed as a sum

$$y = \bar{y} + v,$$

where  $\bar{y} = \bar{y}(P)$  is a solution to the homogeneous equation

$$\mathcal{L}\bar{y} = 0 \text{ on } \omega, \quad \bar{y} = \mu(P) \text{ on } \gamma, \quad (15)$$

and  $v = v(P)$  is a solution to the nonhomogeneous equation

$$\mathcal{L}v(P) = F(P) \text{ on } \omega, \quad v(P) = 0 \text{ on } \gamma. \quad (16)$$

We have already obtained estimate (14) for  $\bar{y}(P)$  and so it remains to evaluate the function  $v(P)$ .

Theorem 3. If  $D(P) > 0$  everywhere on the grid  $\omega$ , then a solution of problem (16) admits the estimate

$$\|v\|_C \leq \left\| \frac{F}{D} \right\|_C. \quad (17)$$

Proof. Let a majorant  $\Upsilon(P)$  be taken such that

$$\mathcal{L}\Upsilon = |F(P)|, \quad \Upsilon|_{\gamma} = 0, \quad \Upsilon(P) \geq 0 \text{ for } P \in \omega \cup \gamma$$

and  $\Upsilon(P)$  attain its maximum at a node  $P_0 \in \omega$ . As far as

$\Upsilon(P_0) = \|Y\|_C$  is concerned, the equation

$$D(P_0)\Upsilon(P_0) + \sum_{Q \in \text{Patt}'(P_0)} B(P_0, Q)(\Upsilon(P_0) - \Upsilon(Q)) = |F(P_0)|$$

$$\left\{ \begin{array}{l} (D(P_0)) \\ (A(P_0) - \sum_{Q \in \text{Patt}'(P_0)} B(P_0, Q))\Upsilon(P_0) + \sum_{Q \in \text{Patt}'(P_0)} B(P_0, Q)(\Upsilon(P_0) - \Upsilon(Q)) \end{array} \right\}$$

$$\text{yields } D(P_0)\Upsilon(P_0) \leq |F(P_0)|, \quad \Upsilon(P_0) \leq \frac{|F(P_0)|}{D(P_0)} \leq \left\| \frac{F}{D} \right\|_C$$

thereby completing the proof.

Remark Estimate (17) is still valid for the solution of problem (16) provided that instead of (2) other conditions

$$|A(P)| \neq 0, |B(P, Q)| \neq 0,$$

$$D(P) = |A(P)| - \sum_{Q \in \text{Path}'(P)} |B(P, Q)| > 0$$

hold.

Indeed, let  $|V(P)| \geq 0$  take the maximal value at a node  $P_0$ . Because of this fact,

$$\begin{aligned} |A(P_0)| \cdot |V(P_0)| &= \left| \sum_{Q \in \text{Path}'(P_0)} B(P_0, Q) V(Q) + F(P_0) \right| \\ &\leq \sum_{Q \in \text{Path}'(P_0)} |B(P_0, Q)| \cdot |V(Q)| + |F(P_0)|, \end{aligned}$$

whence it follows that

$$D(P_0) |V(P_0)| \leq |F(P_0)|, \quad \|V\|_C = |V(P_0)| \leq \frac{|F(P_0)|}{|D(P_0)|} \leq \|\frac{F}{D}\|_C.$$

It may happen that  $D(P) = 0$  on a subset  $\tilde{\omega}$  of the grid  $\omega$  and  $D(P) > 0$  on the complement of  $\tilde{\omega}$  to  $\omega$ :  $\tilde{\omega} + \omega^* = \omega$ . This type of situation is covered by the following assertion.

Theorem 4 Let the conditions

$$D(P) = 0 \text{ for } P \in \tilde{\omega}, \quad D(P) > 0 \text{ for } P \in \omega^*$$

hold, where  $\tilde{\omega}$  is a connected grid. Then for a solution of problem (16) with the right-hand side

$$F(P) = 0 \text{ for } P \in \tilde{\omega}, \quad F(P) \neq 0 \text{ for } P \in \omega^*,$$

the estimate

$$\|V\|_C \leq \left\| \frac{F}{D} \right\|_{C^*} \quad (18)$$

is valid in the norm  $\|f\|_{C^*} = \max_{P \in \omega^*} |f(P)|$ .

-43-

Proof. Let  $\mathcal{Y}(P)$  be a majorant and  $\mathcal{L}Y = |F(P)|$  on the grid  $w$ ,  $Y|_Y = 0$ ,  $Y \geq 0$ . The function  $\mathcal{Y}(P)$  should attain its maximum on a finite set  $w + Y$  at some node, not belonging to the boundary, because  $Y|_Y = 0$ . Also, it does not enter the grid  $w$  due to the connectedness of  $w$  and the maximum principle. Hence,

$$\max_{P \in w} Y(P) = \max_{P \in w^*} Y(P) = Y(P_0),$$

where  $P_0$  is a node on the set  $w^*$ .

By the initial hypothesis,  $D(P_0) > 0$ . Arguing as in the proof of T. 3 we arrive at (18). An analog of the remark to T. 3 is still valid for that case.

### Stability and convergence of the Dirichlet difference problem

#### 1. Estimation of a solution of the Dirichlet difference problem.

We construct a uniform estimate of a solution of the Dirichlet difference problem

$$\Delta Y = -\varphi \text{ at the center nodes}$$

$$\Delta^* Y = -\varphi \text{ at the irregular nodes.}$$

$$Y = \mu(x) \text{ on the boundary,}$$

where

$$\Delta Y = \sum_{\alpha=1}^D \Delta_\alpha Y, \quad \Delta_\alpha Y = \frac{Y_{x_\alpha+1} - 2Y_{x_\alpha} + Y_{x_\alpha-1}}{h_\alpha^2},$$

$$\Delta^* Y = \sum_{\alpha=1}^D \Delta_\alpha^* Y, \quad \Delta_\alpha^* Y = \frac{1}{h_\alpha} \left( \frac{Y^{(\alpha)} - Y}{h_\alpha^+} - \frac{Y - Y^{(-\alpha)}}{h_\alpha^-} \right).$$

The problem (1) can be written

$$\bar{\Delta} Y = -\varphi, \quad x \in w, \quad Y|_Y = \mu(x), \quad (2)$$

where  $\bar{\Delta}$  coincides with  $\Delta^*$  at the near-boundary node, and with  $\Delta$  at the remaining inner nodes.

The problem (1) also can be written as

$$A(x)y(x) = \sum_{\{j\} \in \text{Patt}(x)} B(x,j)y(j) + f(x), \quad x \in \omega, \quad y|_{\partial\Omega} = \mu(x), \quad (3)$$

where

$$A(x) > 0, \quad B(x,j) > 0, \quad D(x) = A(x) - \sum_{\{j\} \in \text{Patt}(x)} B(x,j) > 0.$$

We now represent a solution of problem (1) as a sum

$$y = \bar{y} + \tilde{y},$$

where  $\bar{y}$  and  $\tilde{y}$  are, respectively, solutions of the appropriate problems,

$$\bar{\Lambda}\bar{y} = 0, \quad x \in \omega_n, \quad \bar{y} = \mu \quad \text{for } x \in \partial\Omega, \quad (4)$$

$$\bar{\Lambda}\tilde{y} = -\varphi, \quad x \in \omega_n, \quad \tilde{y} = 0 \quad \text{for } x \in \partial\Omega. \quad (5)$$

An estimate for a solution of problem (4) such as

$$\|\bar{y}\|_C \leq \|\mu\|_{C_p}, \quad (6)$$

has been derived before.

Having decomposition the right-hand side  $\varphi$  as

$$\varphi = \ddot{\varphi} + \varphi^*,$$

where  $\ddot{\varphi} = \varphi$  and  $\varphi^* = 0$  at the strictly inner nodes  $x \in \omega_h^*$  and

$\ddot{\varphi} = 0$  and  $\varphi^* = \varphi$  at the near-boundary nodes  $x \in \omega_h^+$ , we

consider

$$\tilde{y} = v + w$$

with  $v$  and  $w$  being solutions of the problems

$$\bar{\Lambda}v = -\ddot{\varphi} \quad \text{for } x \in \omega_n, \quad v|_{\partial\Omega} = 0, \quad (7)$$

$$\bar{\Lambda}w = -\varphi^* \quad \text{for } x \in \omega_n, \quad w|_{\partial\Omega} = 0. \quad (8)$$

We are going to evaluate separately each of the terms  $v(x)$  and  $w(x)$ .

In order to estimate  $v(x)$ , it is necessary to construct a majorant  $\mathcal{D}(x)$ . Assume that the origin is inside the domain  $\Omega$ , we try to determine a majorant of the type

- 45 -

$$Y(x) = K(R^2 - r^2), \quad r^2 = \sum_{\alpha=1}^P x_\alpha^2,$$

where  $K > 0$  is a constant and  $R$  is the radius of a  $P$ -dimensional ball (a circle for  $P=2$ ) with center at the origin containing entirely the domain  $\Omega$ .  $K$  will be chosen.  
Since,

$$\sum_{\alpha \neq \beta} x_\alpha x_\beta = 0 \text{ for } \alpha \neq \beta \text{ and}$$

$$\Lambda_x x_\alpha^2 = \frac{(x_\alpha + h_\alpha)^2 - 2x_\alpha^2 + (x_\alpha - h_\alpha)^2}{h_\alpha^2} = 2,$$

$$\Lambda_x^* x_\alpha^2 = \frac{1}{h_\alpha} \left[ \frac{(x_\alpha + h_\alpha)^2 - x_\alpha^2}{h_{\alpha+}} - \frac{x_\alpha^2 - (x_\alpha - h_\alpha)^2}{h_{\alpha-}} \right] =$$

$$= \frac{1}{h_\alpha} \left[ \frac{x_\alpha^2 + 2x_\alpha h_{\alpha+} + h_{\alpha+}^2 - x_\alpha^2}{h_{\alpha+}} - \frac{x_\alpha^2 - x_\alpha^2 + 2x_\alpha h_{\alpha-} - h_{\alpha-}^2}{h_{\alpha-}} \right]$$

$$= \frac{1}{h_\alpha} \left[ 2x_\alpha + h_{\alpha+} - 2x_\alpha + h_{\alpha-} \right] = \frac{h_{\alpha+} + h_{\alpha-}}{h_\alpha} = 2 \cdot \frac{h_{\alpha+} + h_{\alpha-}}{2h_\alpha} = 2\Theta_\alpha,$$

$$\Theta_\alpha = \frac{h_{\alpha+} + h_{\alpha-}}{2h_\alpha}$$

We find that

$$\Lambda Y = \sum_{\alpha=1}^P \Lambda_x Y = -2PK \quad \text{for } x \in \tilde{\omega}_h,$$

$$\Lambda^* Y = -2P\Theta K \quad \text{for } x \in \tilde{\omega}_h^*,$$

where  $\Theta = P^{-1} \sum_{\alpha=1}^P \Theta_\alpha$ . Here  $\Theta_\alpha = 1$  if a node  $x \in \tilde{\omega}_h^*$  is regular along the direction  $x_\alpha$ . Thus,  $Y$  is a solution of the problem

$$\bar{\Lambda} Y = -\bar{F}(x), \quad Y|_{\tilde{\omega}_h} = K(R^2 - r^2) \Big|_{\tilde{\omega}_h} \geq 0,$$

where  $\bar{F}(x) = 2PK$  for  $x \in \tilde{\omega}_h$  and  $\bar{F}(x) = 2P\Theta K$  for  $x \in \tilde{\omega}_h^*$ .

Comparison with problem (7), where  $F = \varphi$ , that is,  $F = 0$  for  $x \in \omega_h^*$  and  $w|_{\partial\omega_h} = 0$ , shows that  $\bar{F}(x) \geq |F(x)| = |\varphi(x)|$  if we accept  $K = \frac{1}{2P} \|\varphi\|_C$ . Here the conditions of the comparison theorem are valid as long as  $\bar{F}(x) > |F(x)| = 0$  for  $x \in \omega^*$ , assuring the relation  $\|w\|_C \leq \|F\|_C$ .

Since  $\|F\|_C \leq KR^2$ , so

$$\|w\|_C \leq \frac{R^2}{2P} \|\varphi\|_C = \frac{R^2}{2P} \|F\|_C \quad (9)$$

is valid in the norm  $\|F\|_C = \max_{x \in \omega_h^*} |\varphi(x)|$ .

Our next step is the estimation of the function  $w(x)$ , which is a solution of problem (8):

$$(8) \quad \bar{A}w = -\varphi^*, \quad x \in \omega_h, \quad w|_{\partial\omega_h} = 0, \quad \varphi^* = \begin{cases} 0, & x \in \omega_h^* \\ \varphi, & x \in \omega_h^* \end{cases}$$

First, we are going to show that for problem (8)

$$D(x) \geq \frac{1}{h^2}, \quad \text{where } h = \max_x, \quad x \in \omega_h^*, \quad (10)$$

$$D(x) = 0 \quad \text{for } x \in \omega_h^*. \quad (11)$$

(11) is (clearly) simple to show.

After that, we look at eq. (8) at a near-boundary node

$$x \in \omega_h^*$$

$$A(x)w(x) = \sum_{q \in \text{latt}(x)} B(x, q)w(q) + F(x), \quad (12)$$

$$F(x) = \varphi^*(x), \quad w|_{\partial\omega_h} = 0.$$

If one of the nodes  $q = q_0$ , say  $q_0 = x^{(+1)}$ , happens to be a boundary node, then  $w(q_0) = 0$  and the neighborhood  $\text{latt}(x)$  contains no point  $q_0$ .

In this case the function  $D(x)$  takes on the value

$$D(x) = A(x) - \sum_{\substack{\beta \in \text{Patt}(x) \\ \beta \neq \beta_0}} B(x, \beta)$$

$$= A(x) - \left[ \sum_{\beta \in \text{Patt}(x)} B(x, \beta) - B(x, \beta_0) \right] = B(x, \beta_0),$$

since  $A(x) = \sum_{\beta \in \text{Patt}(x)} B(x, \beta)$  for the Laplace equation.

$$\Rightarrow D(x) = B(x, x^{(+\omega)}) > 0.$$

If a node  $x$  is near-boundary not only with respect to  $\mathbb{X}_\alpha$ , but also in other directions, then sum (12) contains no other terms for  $\beta = \beta_1, \beta_2, \dots, \beta_p$ , so that

$$D(x) = B(x, \beta_0) + B(x, \beta_1) + \dots + B(x, \beta_p) > 0.$$

Let  $x \in \mathbb{W}_h^*$  be a near-boundary and irregular node only in some direction  $\mathbb{X}_\alpha$  and  $\beta_0 = x^{(+\omega)} \in \mathbb{Y}_h, x^{(-\omega)} \in \mathbb{W}_h$ .

From the equation

$$\Lambda_\alpha^* w + \sum_{\substack{\beta=1 \\ \beta \neq \alpha}}^p \Lambda_\beta w = -\varphi^*(x),$$

where

$$\Lambda_\beta y = y_{\mathbb{X}_\beta} x_\beta, \\ \Lambda_\alpha^* w = \frac{1}{h_\alpha} \left( \frac{w^{(+\omega)} - w}{h_{\alpha+}} - \frac{w - w^{(-\omega)}}{h_\alpha} \right) = -\frac{1}{h_\alpha} \left( \frac{w}{h_{\alpha+}} + \frac{w - w^{(-\omega)}}{h_\alpha} \right),$$

we establish the relations

$$A(x) = \frac{1}{h_\alpha h_{\alpha+}} + \frac{1}{h_\alpha^2} + \sum_{\substack{\beta=1 \\ \beta \neq \alpha}}^p \frac{2}{h_\beta^2}, \quad \sum_{\beta \in \text{Patt}(x)} B(x, \beta) = \frac{1}{h_\alpha^2} + \sum_{\substack{\beta=1 \\ \beta \neq \alpha}}^p \frac{2}{h_\beta^2}$$

$$\Rightarrow D(x) = A(x) - \sum_{\beta \in \text{Patt}(x)} B(x, \beta) = \frac{1}{h_\alpha h_{\alpha+}} \geq \frac{1}{h^2}$$

48-

If  $x$  is a regular near-boundary node only with respect to  $\alpha$ , then

$$D(x) = \frac{1}{h^2} \geq \frac{1}{h^2}.$$

When, in addition, it turns out to be near-boundary in other directions, the function  $D(x)$  can only increase.

Then

$$\|w\|_C \leq \left\| \frac{\varphi^*}{D} \right\|_{C^*} \leq h^2 \|\varphi^*\|_C. \quad (13)$$

Collecting estimates (6), (9), (13) and then involving the well-known inequality

$$\|y\|_C \leq \|\bar{g}\|_C + \|v\|_C + \|w\|_C,$$

we establish the following theorem.

Theorem 1. For a solution of the Dirichlet difference problem the estimate

$$\|y\|_C \leq \|u\|_C + \frac{R^2}{2P} \|\varphi\|_C + h^2 \|\varphi\|_{C^*}. \quad (14)$$

holds with

$$\|u\|_C = \max_{x \in w_{n+1}} |f(x)|, \quad \|f\|_C = \max_{x \in w_n} |f(x)|,$$

$$\|\varphi\|_{C^*} = \max_{x \in w_n^*} |f(x)|, \quad \|\varphi\|_C = \max_{x \in w_n} |f(x)|.$$

This theorem express the stability of the Dirichlet difference problem (1) with respect to the boundary data and the right-hand side.

2. The uniform convergence and the order of accuracy of a difference scheme.

2. The uniform convergence and the order of accuracy of a difference scheme.

Let

$$Z = y - u,$$

where  $y$  is a solution of problem (1), and  $u = u(x)$  is a solution of problem (1). Substitution  $y = z + u$  into (1) or (2) yields

$$\bar{\Lambda}y = -\varphi \Rightarrow \bar{\Lambda}(z+u) = \bar{\Lambda}z + \Lambda u = -\varphi$$

$$\Rightarrow \bar{\Lambda}z = -\underbrace{(\bar{\Lambda}u + \varphi)}_{\Psi(x)}, \quad x \in \omega, \quad z|_{\bar{\gamma}} = 0.$$

$$\Rightarrow \bar{\Lambda}z = -\varphi(x), \quad x \in \omega, \quad z|_{\bar{\gamma}} = 0, \quad (15)$$

where  $\varphi(x) = \bar{\Lambda}u + \varphi(x)$  is the residual.

We stated before that

$$\varphi(x) = O(|h|^2) = O(h^2), \quad \text{at the regular nodes},$$

$$\varphi(x) = O(1) \quad \text{at the irregular nodes},$$

or, more specifically,

$$|\varphi| \leq \frac{M_4|h|^2}{12} \leq p \frac{M_4}{12} h^2 \quad \text{at the regular nodes},$$

$$|\varphi| \leq p M_2 \quad \text{at the irregular nodes}.$$

where

$$M_k = \max_{\substack{x \in \bar{\Gamma} \\ 1 \leq d \leq p}} \left| \frac{\partial^k u}{\partial x^k} \right|, \quad k=2,3,4, \dots, \quad |h|^2 = \sum_{d=1}^p h_d^2, \quad h = \max_{1 \leq d \leq p} h_d$$

By Theorem 1 (estimate (4)), we have

$$\Rightarrow \|Z\|_C = \|y - u\|_C \leq \underbrace{\|Z\|_C \leq \frac{R^2}{2p} \|\varphi\|_C + h^2 \|\varphi\|_C}_{\leq \left( \frac{R^2}{24} M_4 + p M_2 \right) h^2}$$

Therefore, we have proved

Theorem 2. If  $u(x) \in C^4(\bar{\Gamma})$ , that is, a solution possesses continuous derivatives in  $\bar{\Gamma} = \Gamma + \Gamma$  of the first four orders, then the difference scheme converges uniformly with the rate  $O(h^2)$ . That is, it is of second-order accuracy, so that estimate (16) is valid.

---

## Higher-accuracy schemes for Poisson's equation.

On the basis of the "cross" scheme it is possible to construct a scheme with the error of approximation  $O(h^4)$  or  $O(h^6)$  on a solution in the core of a square (cube) grid. In order to raise the order of approximation, we use the fact that  $U = U(x)$  is a solution of Poisson's eq.

$$\Delta U = -f(x) \quad (1)$$

Without loss of generality we may restrict ourselves to the careful analysis of the two-dimensional case ( $D=2$ ) where

$$\Delta U = L_1 U + L_2 U, \quad L_2 U = \frac{\partial^2 U}{\partial x_2^2},$$

Let

$$\Lambda U = (\lambda_1 + \lambda_2)U, \quad \lambda_2 U = U_{x_2 x_2}.$$

Then

$$\Lambda U - L_2 U = \frac{h_1^2}{12} L_1^2 U + \frac{h_2^2}{12} L_2^2 U + O(h^4). \quad (2)$$

From the equation

$$L_1 U + L_2 U = -f(x),$$

we find that

$$L_1^2 U = -L_1 f - L_1 L_2 U, \quad L_2^2 U = -L_2 f - L_1 L_2 U.$$

So that

$$\Lambda U = L_2 U + \frac{h_1^2}{12} (-L_1 f - L_1 L_2 U) + \frac{h_2^2}{12} (-L_2 f - L_1 L_2 U) + O(h^4)$$

$$= L_2 U - \frac{h_1^2}{12} L_1 f - \frac{h_2^2}{12} L_2 f - \frac{h_1^2 + h_2^2}{12} L_1 L_2 U + O(h^4). \quad (3)$$

We substitute here  $L_2 U = -f$  and replace  $L_1 L_2 U$  by the difference operator. We use the approximation

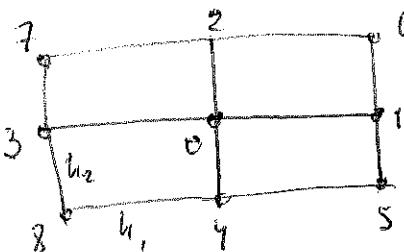
$$L_1 L_2 U = \frac{\partial^2 U}{\partial x_1^2 \partial x_2^2} - \Lambda_1 \Lambda_2 U = U_{\bar{x}_1 \bar{x}_2 \bar{x}_2 \bar{x}_1}$$

The expression for  $\Lambda_1 \Lambda_2 U$  such that  $\Lambda_1 \Lambda_2 U = \Lambda_1 \left[ \frac{U(x_1, x_2 - h_2) - 2U(x_1, x_2) + U(x_1, x_2 + h_2)}{h_2^2} \right]$

$$= \frac{1}{h_1^2 h_2^2} \left\{ U(x_1 - h_1, x_2 - h_2) - 2U(x_1 - h_1, x_2) + U(x_1 - h_1, x_2 + h_2) \right. \\ \left. - 2U(x_1, x_2 - h_2) + 4U(x_1, x_2) - 2U(x_1, x_2 + h_2) \right. \\ \left. + U(x_1 + h_1, x_2 - h_2) - 2U(x_1 + h_1, x_2) + U(x_1 + h_1, x_2 + h_2) \right\}$$

is needed in the estimation  
of the error of approximation

to  $\Lambda_1 \Lambda_2 U - L_1 L_2 U$  by virtue  
of the well-established expansion



$$\Delta V = V_{\bar{x}\bar{y}} = \frac{V(x+h) - 2V(x) + V(x-h)}{h^2} = V''(\xi), \quad \xi = x + \theta h, \quad |\theta| \leq 1, \quad (4)$$

assuming that  $V(x)$  has a continuous second derivative  
on the segment  $[x-h, x+h]$ ;

$$\Delta V = V_{\bar{x}\bar{y}} = V''(x) + \frac{h^2}{12} V^{(4)}(\xi^*), \quad \xi^* = x + \theta^* h, \quad |\theta^*| \leq 1, \quad (5)$$

$V(x)$  has a continuous fourth derivative on the segment  
 $[x-h, x+h]$ . By relating  $x_1$  to be fixed we might have

$$\Lambda_2 U = L_2 U(x_1, x_2) + \frac{h_2^2}{12} \frac{\partial^4 U}{\partial x_2^4}(x_1, \xi_2), \quad \xi_2 = x_1 + \theta_2 h_2, \quad |\theta_2| \leq 1.$$

Then

$$\Lambda_1 \Lambda_2 U(x_1, x_2) = \Lambda_1 L_2 U(x_1, x_2) + \frac{h_2^2}{12} \Lambda_1 \frac{\partial^4 U}{\partial x_2^4}(x_1, \xi_2).$$

Applying formula (5) with  $V = L_2 U$  and  $x = x_1$  to the first  
summand yields

$$\Lambda_1 L_2 U(x_1, x_2) = L_1 L_2 U(x_1, x_2) + \frac{h_1^2}{12} \frac{\partial^4 U}{\partial x_1^4}(\xi_1^*, x_2), \quad \xi_1^* = x_1 + \theta_1^* h_1, \quad |\theta_1^*| \leq 1.$$

By the formula (4):

$$\frac{h_1^2}{12} \Lambda_1 \frac{\partial^4 U}{\partial x_1^4}(x_1, \xi_1^*) = \frac{h_1^2}{12} \frac{\partial^6 U}{\partial x_1^2 \partial x_2^4}(\xi_1, \xi_2), \quad \xi_1 = x_1 + \theta_1 h_1, \quad |\theta_1| \leq 1.$$

Therefore,

$$(A, A_2 - L, L_2)U = O(h_1^2) + O(h_2^2) = O(|h|^2).$$

Substituting into (3) the expression for  $A, A_2 U$  in place of  $L, L_2 U$

$$L, L_2 U = A, A_2 U + O(|h|^2),$$

and involving the eq.  $LU = -f(x)$ , we finally get

$$\begin{aligned} AU &= LU - \frac{h_1^2 + h_2^2}{12} A, A_2 U - \frac{h_1^2}{12} L_1 f - \frac{h_2^2}{12} L_2 f + O(|h|^4) \\ &= -\left(f + \frac{h_1^2}{12} L_1 f + \frac{h_2^2}{12} L_2 f\right) - \frac{h_1^2 + h_2^2}{12} A, A_2 U + O(|h|^4). \end{aligned} \quad (6)$$

Because of this, the equation

$$\begin{cases} A'y = -\varphi, & A'y = Ay + \frac{h_1^2 + h_2^2}{12} A, A_2 y, \\ \varphi = f + \frac{h_1^2}{12} L_1 f + \frac{h_2^2}{12} L_2 f, \end{cases} \quad (7)$$

provides an approximation of order 4 on a solution  $U = U(x)$  of Poisson's equation (1). Indeed, formula (8) gives

$$A'U + \varphi = (A'U + \varphi) - (LU + f) = O(|h|^4), \quad L = L_1 + L_2.$$

The operator  $A'$  is defined on the nine-point "box" pattern.

(7) is representable by

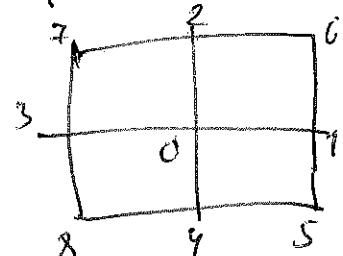
$$\begin{aligned} \frac{5}{3} \left( \frac{1}{h_1^2} + \frac{1}{h_2^2} \right) y &= \frac{1}{6} \left( \frac{5}{h_1^2} - \frac{1}{h_2^2} \right) (y^{(+1_1)} + y^{(-1_1)}) \\ &\quad + \frac{1}{6} \left( \frac{5}{h_2^2} - \frac{1}{h_1^2} \right) (y^{(+1_2)} + y^{(-1_2)}) \\ &\quad + \frac{1}{12} \left( \frac{1}{h_1^2} + \frac{1}{h_2^2} \right) (y^{(+1_1, +1_2)} + y^{(+1_1, -1_2)}) \\ &\quad + (y^{(-1_1, -1_2)} + y^{(-1_1, +1_2)}) + \varphi. \end{aligned}$$

where  $y^{(\pm 1_i)} = y(x_i \pm h_i, x_2)$ ,  $y^{(+1_i, -1_2)} = y(x_i + h_1, x_2 - h_2)$ , etc.

On the square grid ( $h_1 = h_2 = h$ ) the final result is:

$$y_0 = \frac{4(y_1 + y_2 + y_3 + y_4) + y_5 + y_6 + y_7 + y_8}{20} + \frac{3}{10} h^2 \varphi.$$

$$\varphi = f + \frac{h_1^2}{12} \Lambda_1 f + \frac{h_2^2}{12} \Lambda_2 f.$$



2. Estimation of a solution of the difference boundary-value problem:

Consider now the difference Dirichlet problem for the scheme of accuracy  $O(h^4)$  in the rectangle  $\bar{\Omega} = \{0 \leq x_i \leq l_2, i=1,2\}$ :

$$\begin{cases} \Lambda' y = -\varphi, \quad x \in \omega_h, \quad y|_{\gamma_h} = \mu(x) \\ \varphi = f + \frac{h_1^2}{12} \Lambda_1 f + \frac{h_2^2}{12} \Lambda_2 f, \end{cases} \quad (g)$$

where  $\Lambda' y$  is given by (7). Each of the grid nodes is regular, the boundary  $\gamma_h$  of the grid contains all the nodes on the boundary  $\Gamma$  excluding the vertices of the rectangle.

Let  $z = y - u$ :

$$\Lambda' z = -\varphi, \quad x \in \omega_h, \quad z = 0 \text{ on } \gamma_h, \quad (10)$$

where  $\varphi = \Lambda' u + \psi = O(h^4)$  for  $x \in \omega_h$ , if  $u \in C^6$ .

Also, we have

$$\beta(x, z) \geq 0 \quad \text{for} \quad \frac{1}{\sqrt{5}} \leq \frac{h_1}{h_2} \leq \sqrt{5} \quad (11)$$

$$\left( \frac{5}{h_1^2} - \frac{1}{h_2^2} \geq 0 \Rightarrow 5 > \left(\frac{h_1}{h_2}\right)^2 \Rightarrow \frac{h_1}{h_2} \leq \sqrt{5} \quad \left( \frac{5}{h_1^2} - \frac{1}{h_2^2} \geq 0 \quad \left(\frac{h_1}{h_2}\right)^2 \geq \frac{1}{5} \Rightarrow \frac{h_1}{h_2} \geq \frac{1}{\sqrt{5}} \right) \right) \Rightarrow \frac{1}{\sqrt{5}} \leq \frac{h_1}{h_2} \leq \sqrt{5}$$

To estimate the solution of problem (10), we take the majorant of the type

$$Y(x) = K(l_1^2 - x_1^2 + l_2^2 - x_2^2).$$

Taking into account that  $\Lambda Y = -4K$ ,  $\Lambda_1 \Lambda_2 Y = 0$ ,

$$\|Y\| \leq K(l_1^2 + l_2^2) \text{ and accepting } 4K = \|Y\|_C,$$

~~we obtain~~ for the solution of problem (10), we have

$$\|Z\|_C \leq \frac{l_1^2 + l_2^2}{4} \|Y\|_C \text{ provided } \frac{1}{\sqrt{5}} \leq \frac{l_1}{h_1} \leq \sqrt{5}$$

This implies that

Scheme (9) is of fourth-order accuracy when  $u \in C^{(6)}$ ,  
 $f \in C^{(4)}$  and condition (11) is satisfied.

Remark | The condition (11) is automatically fulfilled, when  $h_1 = h_2 = h$  (on the square grid). When  $u \in C^{(6)}$  and with a proper choice of  $\varphi$  guarantees the sixth order of accuracy of scheme (9) on the square grid.

# ADVANCED NUMERICAL ANALYSIS

Part 2

## Parabolic equations

1D case

Transformation to non-dimensional form

$$\frac{\partial U}{\partial T} = K \frac{\partial^2 U}{\partial X^2}, \quad K - \text{constant}, \quad (1)$$

The solution of which gives the temperature  $U$  at a distance  $X$  from one end of a thin uniform rod after a time  $T$ .

Let  $L$  represent the length of a rod and  $U_0$  some particular temperature such as the maximum or minimum temperature at zero time.

Put

$$x = \frac{X}{L} \quad \text{and} \quad u = \frac{U}{U_0}$$

Then

$$\frac{\partial U}{\partial X} = \frac{\partial U}{\partial x} \cdot \frac{dx}{dX} = \frac{\partial U}{\partial x} \cdot \frac{1}{L}$$

$$\frac{\partial^2 U}{\partial X^2} = \frac{\partial}{\partial X} \left( \frac{\partial U}{\partial x} \right) = \frac{\partial}{\partial x} \left( \frac{\partial U}{\partial X} \cdot \frac{1}{L} \right) = \frac{\partial}{\partial x} \left( \frac{\partial U}{\partial x} \cdot \frac{1}{L} \right) \frac{dx}{dX}$$

$$= \frac{1}{L^2} \frac{\partial^2 U}{\partial x^2}$$

$$\Rightarrow \frac{\partial(u \cdot U_0)}{\partial T} = \frac{K}{L^2} \frac{\partial^2(u \cdot U_0)}{\partial x^2} \Rightarrow \frac{L^2}{K} \frac{\partial u}{\partial T} = \frac{\partial^2 u}{\partial x^2}$$

$$\text{Writing } f = \frac{KT}{L^2}, \text{ we have } \frac{\partial u}{\partial T} = \frac{\partial u}{\partial t} \cdot \frac{dt}{dT} = \frac{\partial u}{\partial t} \cdot \frac{K}{L^2}$$

$$\Rightarrow \frac{L^2}{K} \cdot \frac{K}{L^2} \frac{\partial u}{\partial t} = \frac{\partial u}{\partial t}$$

(2)

$$\Rightarrow \left\{ \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} \right\} \text{ is the non-dimensional form of (1).}$$

## An explicit method of solution

One of the explicit methods of finite-difference approximations to

to

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} \quad (3)$$

i.e.

$$\frac{u_{i,j+1} - u_{i,j}}{\kappa} = \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{h^2},$$

where

$$x = ih, \quad (i=0,1,2,\dots),$$

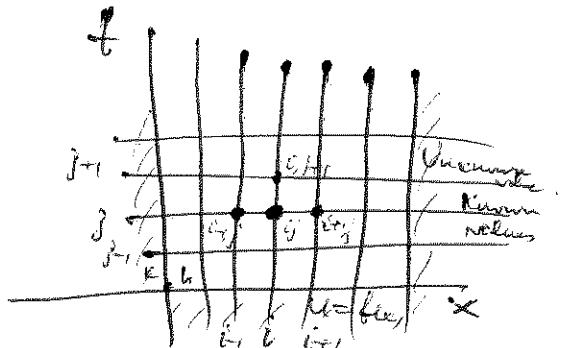
and

$$t = j\kappa, \quad (j=0,1,2,\dots).$$

This can be written as

$$u_{i,j+1} = u_{i,j} + r(u_{i+1,j} - 2u_{i,j} + u_{i-1,j}), \quad (4)$$

where  $r = \frac{\kappa t}{(h)^2} = \frac{r}{h^2}$ , and gives a formula for the unknown "temperature"  $u_{i,j+1}$  at the  $(i,j+1)$ -th mesh point in terms of known "temperature" along the  $j$ -th time row.



Example 1.

(see G. D. Smith)

P. 13

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

(i)  $u=0$  when  $x=0$  and 1 for all  $t$ . (The boundary conditions)

(ii)  $u=2x$  for  $0 \leq x \leq \frac{1}{2}$ , and  $u=2(1-x)$  for  $\frac{1}{2} \leq x \leq 1$  for  $t=0$ . (The initial condition)

It will be proved that this explicit method is valid only when  $0 < r \leq \frac{1}{2}$  (for stability and convergence)

$$0 < r = \frac{k}{h^2} \leq \frac{1}{2}$$

Case 1.  $h = \frac{1}{10}$ ,  $K = \frac{1}{1000}$

So  $r = \frac{\frac{1}{1000}}{\frac{1}{100}} = \frac{1}{10} \left[ < \frac{1}{2} \right]$

Case 2

$$h = \frac{1}{10}, K = \frac{5}{1000} \Rightarrow r = \frac{\frac{5}{1000}}{\frac{1}{100}} = \frac{5}{10} = 0.5$$

Case 3.

$$h = \frac{1}{10}, K = \frac{1}{100} \Rightarrow r = \frac{\frac{1}{100}}{\frac{1}{100}} = 1$$

Solve at FD eq.3 meaning less

Home work.

Take of  $h = \frac{1}{4}$ ,  $K = \frac{1}{20}$

b)  $h = \frac{1}{4}$ ,  $K = \frac{1}{10}$

The analytical soln.

$$u = \frac{8}{\pi^2} \sum_{n=1}^{\infty} \frac{1}{n^2} \left( \sin \frac{n\pi}{2} \right) \left( \sin n\pi x \right) e^{-4n^2 \pi^2 t}$$

### Crank-Nicolson implicit method.

Although the explicit method is computationally simple it has one serious drawback. The time step  $\delta t = \kappa$  is necessarily very small because the process is valid only for  $0 < \frac{\kappa}{h^2} \leq \frac{1}{2}$ , i.e.,  $\kappa \leq \frac{1}{2} h^2$ , and  $U_2 - \bar{U}_X$  must be kept small in order to attain reasonable accuracy.

Crank and Nicolson (1947) proposed, and used, a method that reduces the total volume of calculation and is valid (i.e., convergent and stable) for all finite values of  $\frac{\kappa}{h^2}$  (i.e., convergent and stable) for all finite values of  $\frac{\kappa}{h^2}$ . They replaced  $\frac{\partial U}{\partial X^2}$  by the mean of its FD representations on the  $(j+1)$ th and  $j$ th time rows and approximated the eq.

$$\frac{\partial U}{\partial t} = \frac{\partial^2 U}{\partial X^2} \quad (1)$$

$$\therefore \frac{U_{i,j+1} - U_{i,j}}{\kappa} = \frac{1}{2} \left\{ \frac{U_{i+1,j+1} - 2U_{i,j+1} + U_{i-1,j+1}}{h^2} + \frac{U_{i+1,j} - 2U_{i,j} + U_{i-1,j}}{h^2} \right\}$$

gives

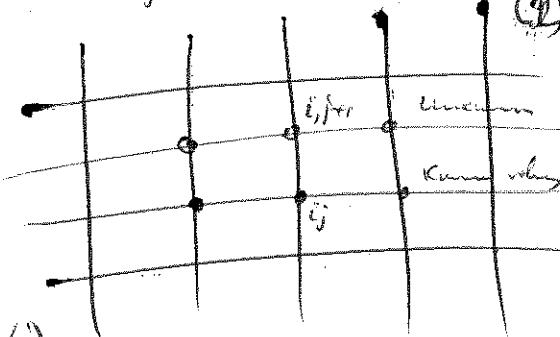
$$2U_{i,j+1} - 2U_{i,j} = r (U_{i+1,j+1} - 2U_{i,j+1} + U_{i-1,j+1} + U_{i+1,j} - 2U_{i,j} + U_{i-1,j})$$

$$\Rightarrow -rU_{i-1,j+1} + (2 + 2r)U_{i,j+1} - rU_{i+1,j+1} = rU_{i-1,j} + (2 - 2r)U_{i,j} + rU_{i+1,j} \quad (2)$$

where  $r = \frac{\kappa}{h^2}$ .

$\Rightarrow$  Implicit Method

See Example 2. (In the next page)



## Derivative boundary conditions

Example 1.

Solve the equation

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} \quad (0 < x < 1) \quad (1)$$

satisfying the initial condition

$$u=1 \text{ for } 0 \leq x \leq 1 \text{ when } t=0,$$

and the b.c.s

$$\frac{\partial u}{\partial x} = u \text{ at } x=0, \text{ for all } t,$$

$$\frac{\partial u}{\partial x} = -u \text{ at } x=1, \text{ for all } t.$$

Using an explicit method and employing central-differences for the boundary conditions.

One explicit finite-difference representation of eq. (1)

(i)

$$\frac{u_{i,j+1} - u_{i,j}}{\delta t} = \frac{u_{i-1,j} - 2u_{i,j} + u_{i+1,j}}{(\delta x)^2}, \quad i=1, 2, \dots, N-1,$$

i.e.,

$$u_{i,j+1} = u_{i,j} + r(u_{i-1,j} - 2u_{i,j} + u_{i+1,j}), \quad (2)$$

$$\text{where } r = \frac{\delta t}{(\delta x)^2}.$$

At  $x=0$ , (i.e.,  $i=0$ )

$$u_{0,j+1} = u_{0,j} + r(u_{-1,j} - 2u_{0,j} + u_{1,j}). \quad (3)$$

The b.c. at  $x=0$ , in terms of central differences, can be written as

$$\frac{u_{1,j} - u_{-1,j}}{2\delta x} = u_{0,j}. \quad (4)$$

Eliminating  $U_{i,j}$  between (3) and (4) gives

$$U_{i,j} = U_{j,j} - 2\delta x U_{o,j}$$

$$\Rightarrow U_{o,j+1} = U_{o,j} + r(U_{j,j} - 2\delta x U_{o,j} - 2U_{o,j} + U_{i,j})$$

$$\Rightarrow \boxed{U_{o,j+1} = U_{o,j} + 2r\{U_{j,j} - (1 + \delta x)U_{o,j}\}} \quad (5)$$

At  $x=1$  (i.e.,  $i=N$ ), from (2), we have

$$U_{N,j+1} = U_{N,j} + r(U_{N-1,j} - 2U_{N,j} + \underbrace{U_{N+1,j}}_?) \quad (6)$$

and the central difference approximation of the b.c. at  $x=1$ ,

we have

$$\frac{\overbrace{U_{N+1,j} - U_{N-1,j}}^?}{2\delta x} = -U_{N,j} \quad (7)$$

From (6) and (7) we obtain

$$\left\{ \begin{array}{l} U_{N+1,j} = U_{N-1,j} - 2\delta x U_{N,j} \\ U_{N,j+1} = U_{N,j} + r(U_{N-1,j} - 2U_{N,j} + U_{N-1,j} - 2\delta x U_{N,j}) \\ \Rightarrow \boxed{U_{N,j+1} = U_{N,j} + 2r\{U_{N-1,j} - (1 + \delta x)U_{N,j}\}} \end{array} \right. \quad (8)$$

The system of eq.s (2), (5), and (8) is an approximate  
of the boundary value problem.

Example 2, Re-solve the example 1 using an explicit method and employing a forward-difference for the boundary conditions at  $x=0$  ( $x=1$ )

From (2) :

$$U_{i,j+1} = U_{ij} + r(U_{i+1,j} - 2U_{ij} + U_{i-1,j})$$

for  $i=1$

$$U_{1,j+1} = U_{1,j} + r(U_{0,j} - 2U_{1,j} + U_{2,j}) \quad (9)$$

The boundary condition at  $x=0$ , namely  $\frac{\partial u}{\partial x} = u$ , in terms of a forward difference is

$$\frac{U_{1,j} - U_{0,j}}{\delta x} = U_{0,j}$$

$$\text{so } U_{0,j} = U_{1,j} - \delta x U_{0,j} \Rightarrow \\ \Rightarrow U_{0,j} = \frac{U_{1,j}}{1 + \delta x} \quad (10)$$

From (9) and (10), we obtain

$$U_{1,j+1} = U_{1,j} + r \left( \frac{U_{0,j}}{1 + \delta x} - 2U_{1,j} + U_{2,j} \right) \\ \Rightarrow U_{1,j+1} = \left( 1 - 2r + \frac{r}{1 + \delta x} \right) U_{1,j} + r U_{2,j} \quad (11)$$

Similarly the finite-difference eq. for  $\frac{\partial u}{\partial x} = -u$  at  $x=1$  is ~~the~~ can be obtained.

Example 3 Solve example 1 by the Crank-Nicolson method:

$$\frac{U_{i,j+1} - U_{i,j}}{\delta t} = \frac{1}{2} \left\{ \frac{U_{i-1,j+1} - 2U_{i,j+1} + U_{i+1,j+1}}{(\delta x)^2} + \frac{U_{i+1,j} - 2U_{i,j} + U_{i-1,j}}{(\delta x)^2} \right\}$$

$$\Rightarrow -rU_{i-1,j+1} + (2+2r)U_{i,j+1} - rU_{i+1,j+1} = rU_{i-1,j} + (2-2r)U_{i,j} + rU_{i+1,j} \quad (12)$$

The central difference representation of the b.c. at  $x=0$  is

$$\frac{U_0 - U_{-1,j}}{2\delta x} = U_{0,j},$$

from which it follows that

$$U_{-1,j} = U_{1,j} - 2\delta x U_{0,j}$$

$$\text{and } U_{-1,j+1} = U_{1,j+1} - 2\delta x U_{0,j+1}.$$

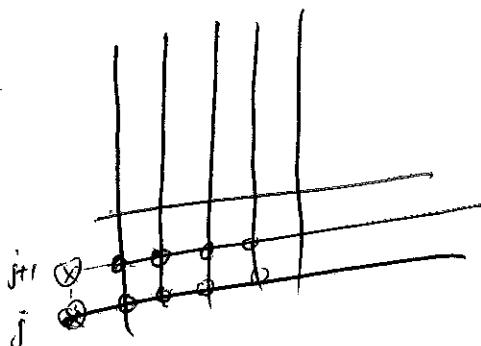
From (12) for  ~~$i=0$~~   $\underline{i=0}$ :

$$-rU_{-1,j+1} + (2+2r)U_{0,j+1} - rU_{1,j+1} = rU_{-1,j} + (2-2r)U_{0,j} + rU_{1,j}$$

$$-r(U_{-1,j+1} - 2\delta x U_{0,j+1}) + (2+2r)U_{0,j+1} - rU_{1,j+1} = r(U_{1,j} - 2\delta x U_{0,j}) + (2-2r)U_{0,j} + rU_{1,j}$$

$$-2rU_{-1,j+1} + (2r\delta x + 2+2r)U_{0,j+1} = 2rU_{1,j} + (2r\delta x + 2-2r)U_{0,j}$$

The obtained system can be solved  
by the direct elimination method.



The local truncation error

Consider the DE:

$$h(U_{ij}) = \left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{i,j} = 0$$

and let

$$F(U_{ij}) = 0 = \frac{U_{i,j+1} - U_{i,j}}{K} - \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} \quad (1)$$

be the classical explicit eq.

The local truncation error is defined as:

$$T_{ij} = F(U_{ij}) = \frac{U_{i,j+1} - U_{i,j}}{K} - \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} \quad (2)$$

This can be written

$$U_{i,j+1} = K T_{ij} + r U_{i-1,j} + (1-2r) U_{i,j} + r U_{i+1,j} \quad (3)$$

Assume now that the solution  $U$  of the PDE is known at all mesh points up to and including those at the  $j$ -th time-level. We could then calculate a local approximation  $U_{i,j+1}$  to  $U_{i,j+1}$  by means of (1) and get

$$U_{i,j+1} = r U_{i-1,j} + (1-2r) U_{i,j} + r U_{i+1,j} \quad (4)$$

By (3) and (4) it follows that

$$U_{i,j+1} - U_{i,j+1} = (\text{local}) \text{ discretization error at } (i,j+1) = \underline{F T_{ij}} \quad (5)$$

This shows that the local truncation error at the point  $(i,j)$  is a measure of the (local) discretization error at the point  $(i,j+1)$  when the finite-difference scheme is applied once only to the exact solution values of the PDE, all arithmetic being exact, i.e., without rounding errors.

By means of Taylor's expansion it is easy to express  $T_{ij}$  in terms of powers of  $h$  and  $K$  and part-derivatives of  $U$  at  $(i,j)$ .

~~QF~~

Example 3. Calculate the order of the local truncation error of the classical explicit difference approximation to

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0.$$

By Taylor's expansion

$$U_{i\pm 1,j} = U\{ (i\pm 1)h, k \} = U(x_i \pm h, t_j)$$
$$= U_{ij} \pm h \left( \frac{\partial U}{\partial x} \right)_{ij} + \frac{1}{2} h^2 \left( \frac{\partial^2 U}{\partial x^2} \right)_{ij} \pm \frac{1}{6} h^3 \left( \frac{\partial^3 U}{\partial x^3} \right)_{ij} + \frac{h^4}{24} \left( \frac{\partial^4 U}{\partial x^4} \right)_{ij} \pm \dots$$

$$U_{i,j\pm 1} = U(x_i, t_j \pm k) = U_{ij} \pm k \left( \frac{\partial U}{\partial t} \right)_{ij} + \frac{1}{2} k^2 \left( \frac{\partial^2 U}{\partial t^2} \right)_{ij} \pm \frac{k^3}{6} \left( \frac{\partial^3 U}{\partial t^3} \right)_{ij} \pm \dots$$

Hence

$$T_{ij} = \frac{U_{i+1,j} - U_{i-1,j}}{2h} - \frac{U_{i+1,j} - 2U_{ij} + U_{i-1,j}}{h^2} = \left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{ij} + \frac{k}{2} \left( \frac{\partial^2 U}{\partial t^2} \right)_{ij}$$
$$- \frac{h^2}{12} \left( \frac{\partial^4 U}{\partial x^4} \right)_{ij} + \frac{k^2}{6} \left( \frac{\partial^3 U}{\partial t^3} \right)_{ij} - \frac{1}{360} h^6 \frac{\partial^6 U}{\partial x^6} + \dots \quad (6)$$

But  $U$  is the solution of the DE, so

$$\left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{ij} = 0$$

Therefore the principal part of the local truncation error is

$$\left( \frac{1}{2} k \frac{\partial^2 U}{\partial t^2} - \frac{1}{12} h^2 \frac{\partial^4 U}{\partial x^4} \right)_{ij}.$$

Hence  $T_{ij} = O(k) + O(h^2)$ .

When  $k = rh^2$ ,  $0 < r \leq \frac{1}{2}$ ,  $T_{ij}$  is  $O(k)$  or  $O(h^2)$ .

This error may be further reduced by choosing a special value for  $\frac{k}{h^2}$  because eq. (6) can be written as

$$T_{ij} = \frac{1}{12} h^2 \left( 6 \frac{\kappa}{h^2} \frac{\partial^2 U}{\partial t^2} - \frac{\partial^4 U}{\partial x^4} \right)_{ij} + O(\kappa^2) + O(h^4).$$

By DB

$$\frac{\partial}{\partial t} = \frac{\partial^2}{\partial x^2},$$

$$\frac{\partial}{\partial t} \left( \frac{\partial U}{\partial t} \right) = \frac{\partial^2}{\partial x^2} \left( \frac{\partial^2 U}{\partial x^2} \right)$$

assuming that these derivatives exist. If we put  $\frac{6\kappa}{h^2} = 1$ ,  
then

$$T_{ij} = \underline{O(\kappa^2) + O(h^4)}.$$

$$\overbrace{\left( \kappa = \frac{1}{6} h^2 \right)}$$

## Consistency or compatibility

It is sometimes possible to approximate a parabolic or hyperbolic equation by a finite-difference scheme that is stable but which has a solution that converges to the solution of a different DE as the mesh lengths tend to zero. Such a difference scheme is said to be inconsistent or incompatible with the PDE.

Example. The equation

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0$$

is approximated at the point  $(ih, jk)$  by the difference eq.

$$\frac{U_{i+1,j} - U_{i,j-1}}{2k} - \frac{U_{i,j+1} - 2\{\theta U_{i,j} + (1-\theta)U_{i,j-1}\} + U_{i-1,j}}{h^2} = 0.$$

Show that the local truncation error at this point is

$$\frac{k^2}{6} \frac{\partial^3 U}{\partial t^3} - \frac{h^2}{12} \frac{\partial^4 U}{\partial x^4} + (2\theta-1) \frac{2k}{h^2} \frac{\partial U}{\partial t} + \frac{k^2}{h^2} \frac{\partial^2 U}{\partial t^2} + O\left(\frac{k^2}{h^2}, h^4, k^4\right).$$

Discuss the consistency of this scheme with the PDE when:

(i)  $k=rh$  and (ii)  $k=rh^2$ .

where  $r$  is a positive constant and  $\theta$  a variable parameter.

Expansion of the terms  $U_{i,j+1}$ ,  $U_{i,j-1}$ ,  $U_{i+1,j}$ , and  $U_{i-1,j}$  about  $(\bar{t}, \bar{x})$  by Taylor's series.

$$T_{ij} = \frac{U_{i,j+1} - U_{i,j-1}}{2k} - \frac{U_{i+1,j} - 2\{\theta U_{i,j+1} + (1-\theta)U_{i,j-1}\} + U_{i-1,j}}{h^2}$$

leads to  $\rightarrow V$

$$T_{ij} = \left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{ij} + \left\{ \frac{k^2}{6} \frac{\partial^3 U}{\partial t^3} - \frac{h^2}{12} \frac{\partial^4 U}{\partial x^4} \right. \\ \left. + (2\theta-1) \frac{2k}{h^2} \frac{\partial U}{\partial t} + \frac{k^2}{h^2} \frac{\partial^2 U}{\partial t^2} \right\} + O\left(\frac{k^3}{h^2}, h^4, k^4\right)$$

Case(i)  $k=rh$

As  $h \rightarrow 0$ ,

$$T_{ij} = F(U_{ij}) \rightarrow \left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + (2\theta-1) \frac{2r}{h} \frac{\partial U}{\partial t} + r^2 \frac{\partial^2 U}{\partial t^2} \right)_{ij}$$

When  $\theta \neq \frac{1}{2}$  the third term tends to infinity. When  $\theta = \frac{1}{2}$

the limiting value of  $T_{ij}$  is

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + r^2 \frac{\partial^2 U}{\partial t^2}.$$

In this case the finite-difference equation is consistent with the hyperbolic equation

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + r^2 \frac{\partial^2 U}{\partial t^2} = 0.$$

Hence the difference equation is always inconsistent with

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0 \text{ when } k=rh.$$

Case (ii)  $K = rh^2$ As  $h \rightarrow 0$ ,

$$T_{ij} \rightarrow \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + 2(2\theta - 1)r \frac{\partial U}{\partial t}.$$

When  $\theta \neq \frac{1}{2}$  the difference scheme is consistent with the parabolic eq.

$$\left\{ 1 + 2(2\theta - 1)r \right\} \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0.$$

It is only when  $\theta = \frac{1}{2}$  that the difference scheme is consistent with the given DE. This is then the well-known Du Fort and Frankel three-level explicit scheme which is also stable for all  $r > 0$  (show! At home)

Home work

Show that the explicit scheme

$$\frac{U_{i,j+1} - U_{i,j-1}}{2k} - \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} = 0$$

gives a local truncation error of  $O(k^2) + O(h^2)$ , but unconditional unstable.

## Lax's equivalence theorem

Given a properly posed linear initial-value problem and a linear finite-difference approximation to it that satisfies the consistency condition, stability is the necessary and sufficient condition for convergence.

Def A problem is properly posed if:

- (i) The solution is unique when it exists.
- (ii) The solution depends continuously on the initial data

## Convergence, stability, and consistency

Let  $U$  be the exact solution of a PDE with independent variables  $x$  and  $t$ , and  $u$  the exact solution of the difference equations used to approximate the PDE. Then the FDE is said to be convergent when  $u$  tends to  $U$  at a fixed point or along a fixed  $t$ -level as  $\Delta x$  and  $\Delta t$  both tend to zero.

The difference  $(U-u)$  is called the discretization error.

Let  $F_{ij}(U)=0$  represent the difference equation at the  $(i,j)$ th mesh point. The value  $F_{ij}(U)$  is called the local truncation error at the  $(i,j)$ th mesh point. If this tends to zero as the mesh lengths tend to zero the difference equation is said to be consistent with the PDE.

In practice the solution actually computed is not  $u$  but,  $N$ , (say).  $N$  will be called the numerical solution and  $(u-N)$  the global rounding error  $R$ . The total error at the  $(i,j)$ th mesh point is

$$U_{ij} - N_{ij} = (U_{ij} - u_{ij}) + (u_{ij} - N_{ij}) = \text{discretization error} + \text{global rounding error.}$$

If the growth of  $R_{ij}$  is bounded for all  $i$  as  $j$  tends to infinity, then we say that the difference eq.s are stable.

(The local truncation error and consistency)

Let  $F_{ij}(u) = 0$  represent the difference equation approximating the PDE at the  $(i,j)$ -th mesh point, with exact solution  $u$ . If  $u$  is replaced by  $V$  at the mesh points of the difference equation, where  $V$  is the exact solution of PDE, the value  $F_{ij}(V)$  is called the local truncation error  $T_{ij}$  at the  $(i,j)$  mesh point.

Example Calculate the order of the local truncation error of the classical explicit difference approximation to

$$\frac{\partial V}{\partial t} - \frac{\partial^2 V}{\partial x^2} = 0$$

at the point  $(ih, jk)$ .

$$F_{ij}(u) = \frac{u_{i,j+1} - u_{i,j}}{\kappa} - \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{h^2} = 0$$

Therefore,

$$T_{ij} = F_{ij}(V) = \frac{V_{i,j+1} - V_{i,j}}{\kappa} - \frac{V_{i+1,j} - 2V_{i,j} + V_{i-1,j}}{h^2}$$

By Taylor's expansion

$$V_{i+1,j} = V((i+1)h, jk) = V(x_i + h, t_j) = V_{i,j} + h \left( \frac{\partial V}{\partial x} \right)_{i,j} + \frac{h^2}{2} \left( \frac{\partial^2 V}{\partial x^2} \right)_{i,j} + \frac{h^3}{6} \left( \frac{\partial^3 V}{\partial x^3} \right)_{i,j} + \dots$$

$$V_{i-1,j} = V((i-1)h, jk) = V(x_i - h, t_j) = V_{i,j} - h \left( \frac{\partial V}{\partial x} \right)_{i,j} + \frac{h^2}{2} \left( \frac{\partial^2 V}{\partial x^2} \right)_{i,j} - \frac{h^3}{6} \left( \frac{\partial^3 V}{\partial x^3} \right)_{i,j} + \dots$$

$$V_{i,j+1} = V(x_i, t_j + \kappa) = V_{i,j} + \kappa \left( \frac{\partial V}{\partial t} \right)_{i,j} + \frac{\kappa^2}{2} \left( \frac{\partial^2 V}{\partial t^2} \right)_{i,j} + \frac{\kappa^3}{6} \left( \frac{\partial^3 V}{\partial t^3} \right)_{i,j} + \dots$$

Substitution into the expression for  $T_{ij}$  the gives

$$T_{ij} = \left( \frac{\partial V}{\partial t} - \frac{\partial^2 V}{\partial x^2} \right)_{i,j} + \frac{\kappa}{2} \left( \frac{\partial^2 V}{\partial t^2} \right)_{i,j} - \frac{h^2}{12} \left( \frac{\partial^4 V}{\partial x^4} \right)_{i,j} + \frac{\kappa^2}{6} \left( \frac{\partial^3 V}{\partial t^3} \right)_{i,j} - \frac{h^4}{360} \frac{\partial^6 V}{\partial x^6} + \dots$$

But  $V$  is the solution of the differential eq. so

$$\left( \frac{\partial V}{\partial t} - \frac{\partial^2 V}{\partial x^2} \right)_{i,j} = 0$$

Therefore the principal part of the local truncation error is

$$\left( \frac{1}{2} \kappa \frac{\partial^2 U}{\partial t^2} - \frac{1}{12} h^2 \frac{\partial^4 U}{\partial x^4} \right)_{ij}.$$

Hence

$$T_{ij} = O(\kappa) + O(h^2).$$

when  $\kappa = rh^2$ ,  $0 < r \leq \frac{1}{2}$ ,  $T_{ij}$  is  $O(\kappa)$  or  $O(h^2)$ , as one would expect by inspection.

$$\frac{\partial^2 U}{\partial x^2} \approx \frac{U_{i+1,j} - 2U_{i,j} + U_{i-1,j}}{h^2}, \quad \frac{\partial U}{\partial t} \approx \frac{U_{j+1,i} - U_{j,i}}{\kappa}.$$

This error may be further reduced by choosing a special value for  $\frac{\kappa}{h^2}$  because the eq. for  $T_{ij}$  can be written as

$$T_{ij} = \frac{1}{12} h^2 \left( 6 \frac{\kappa}{h^2} \frac{\partial^2 U}{\partial t^2} - \frac{\partial^4 U}{\partial x^4} \right)_{ij} + O(\kappa^2) + O(h^4)$$

By the differential eq.

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

so

$$\frac{\partial}{\partial t} \left( \frac{\partial U}{\partial t} \right) = \frac{\partial^2}{\partial x^2} \left( \frac{\partial^2 U}{\partial x^2} \right),$$

assuming that these derivatives exist. If we put  $\frac{6\kappa}{h^2} = 1$ , the expression in the brackets is then zero and  $T_{ij}$  is  $O(\kappa^2) + O(h^4)$ . This is of little use in practice because  $\kappa = \frac{1}{6} h^2$  is very small for small  $h$  so the volume of arithmetic needed to advance the solution to a large time-level is substantial.

## Consistency or compatibility

It is sometimes possible to approx. a parabolic or hyperbolic eq. by a FD scheme that is stable, but which has solution that converges to the solution of a different differential eq. as the mesh lengths tend to zero. Such a difference scheme is said to be inconsistent or inconsistency incompatible with the PDE.

Let  $L(U)=0$  represent the PDE in the independent variables  $x$  and  $t$ , with exact solution  $U$ .

Let  $F(u)=0$  represent the approximating PDE with exact solution  $u$ .

Let  $v$  be a continuous function of  $x$  and  $t$  with a sufficient number of continuous derivatives to enable  $L(v)$  to be evaluated at the point  $(ih,jk)$ .

Then the truncation error  $T_{ij}(v)$  at the point  $(ih,jk)$  is defined by

$$T_{ij}(v) = F(v_{ij}) - L(v_{ij}). \quad (1)$$

If  $T_{ij}(v) \rightarrow 0$  as  $h \rightarrow 0, k \rightarrow 0$ , the difference eq. is said to be consistent or compatible with the PDE.

Most authors put  $v=U$  because  $L(U)=0$ . It then follows by eq. (1) that

$$T_{ij}(U) = F(U_{ij}) \rightarrow \text{local truncation error}$$

The difference eq. is then consistent if the limiting value of the local truncation error is zero as  $h \rightarrow 0, k \rightarrow 0$ .

Example 1. The eq.

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0$$

is approximated at the point  $(ih, jc)$  by the difference eq.

$$\frac{U_{i+1,j} - U_{i-1,j}}{2K} - \frac{U_{i+1,j} - 2\{\theta U_{i,j+1} + (1-\theta)U_{i,j-1}\} + U_{i-1,j}}{h^2} = 0,$$

Show that the local truncation error at this point is

$$\frac{K^2}{6} \frac{\partial^3 U}{\partial t^3} - \frac{h^2}{12} \frac{\partial^4 U}{\partial x^4} + (2\theta-1) \frac{2K}{h^2} \frac{\partial U}{\partial t} + \frac{K^2}{h^2} \frac{\partial^2 U}{\partial t^2} + O\left(\frac{K^3}{h^2}, h^4, K^4\right)$$

Discuss the consistency of this scheme with the PDE when:

$$(i) \quad K=rh \quad \text{and} \quad (ii) \quad K=rh^2,$$

where  $r$  is a positive constant and  $\theta$ , a variable parameter.

Expansion of the terms  $U_{i,j+1}$ ,  $U_{i,j-1}$ ,  $U_{i+1,j}$ , and  $U_{i-1,j}$  about  $(ih, jc)$  by Taylor's series, and substitution into

$$T_{ij} = \frac{U_{i+1,j} - U_{i-1,j}}{2K} - \frac{U_{i+1,j} - 2\{\theta U_{i,j+1} + (1-\theta)U_{i,j-1}\} + U_{i-1,j}}{h^2}$$

$$\begin{aligned} U_{i,j\pm 1} &= U(x_i, t_j \pm K) = U_{ij} \pm K \frac{\partial U}{\partial t} + \frac{K^2}{2} \frac{\partial^2 U}{\partial t^2} \pm \frac{K^3}{6} \frac{\partial^3 U}{\partial t^3} + \frac{K^4}{24} \frac{\partial^4 U}{\partial t^4} \dots \\ U_{i\pm 1,j} &= U_{ij} \pm h \frac{\partial U}{\partial x} + \frac{h^2}{2} \frac{\partial^2 U}{\partial x^2} \pm \frac{h^3}{6} \frac{\partial^3 U}{\partial x^3} + \frac{h^4}{24} \frac{\partial^4 U}{\partial x^4} \pm \frac{h^5}{5!} \frac{\partial^5 U}{\partial x^5} + \dots \end{aligned}$$

$$\text{leads to } \frac{U_{i+1,j} - 2\{\theta(U_{i,j+1} - U_{i,j-1}) + U_{i-1,j}\} + U_{i-1,j}}{h^2} = 0$$

$$T_{ij} = \frac{K^2}{6} \frac{\partial^3 U}{\partial t^3} - \frac{2U_{ij} + h^2 \frac{\partial^2 U}{\partial x^2} + \frac{h^4}{12} \frac{\partial^4 U}{\partial x^4} + \dots - 2\theta \cdot 2K \frac{\partial U}{\partial t} - 2U_{ij}}{h^2}$$

$$\left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{(ij)} + \frac{+2K \frac{\partial U}{\partial t} - K^2 \frac{\partial^2 U}{\partial t^2}}{h^2} + O\left(\frac{K^3}{h^2}, h^4, K^4\right)$$

$$T_{ij} = \left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{(i,j)} + \left\{ \frac{k^2}{6} \frac{\partial^3 U}{\partial t^3} - \frac{h^2}{12} \frac{\partial^4 U}{\partial x^4} \right.$$

$$\left. + (2\theta-1) \frac{2k}{h^2} \frac{\partial U}{\partial t} + \frac{k^2}{h^2} \frac{\partial^2 U}{\partial t^2} \right\} + O\left(\frac{k^3}{h^2}, h^4, k^4\right).$$

Case (i)  $k=rh$

As  $h \rightarrow 0$ ,

$$T_{ij} = F(U_{ij}) \rightarrow \left\{ \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + (2\theta-1) \frac{2r}{h} \frac{\partial U}{\partial t} + r^2 \frac{\partial^2 U}{\partial t^2} \right\}_{(i,j)}$$

+ ~~O(h)~~

when  $\theta \neq \frac{1}{2}$  the third term tends to infinity.

when  $\theta = \frac{1}{2}$  the limiting value of  $T_{ij}$  is

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + r^2 \frac{\partial^2 U}{\partial t^2}.$$

In this case the FD eq. is consistent with the hyperbolic equation

$$\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + r^2 \frac{\partial^2 U}{\partial t^2} = 0$$

Hence the difference equation is always inconsistent with  $\frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0$  when  $k=rh$ .

Case (ii)  $k=rh^2$ .

As  $h \rightarrow 0$

$$T_{ij} \rightarrow \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} + 2(2\theta-1)r \frac{\partial U}{\partial t}$$

when  $\theta \neq \frac{1}{2}$  the difference scheme is consistent with the parabolic eq.

$$\left\{ 1 + 2(2\theta-1)r^2 \right\} \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} = 0.$$

It is only when  $\theta = \frac{1}{2}$  that the difference scheme is consistent with the given differential equation.

This is then well-known Du Fort and Frankel three-level explicit scheme which is also stable for all  $r > 0$ .

---

### The Richardson explicit scheme

$$\frac{U_{i,j+1} - U_{i,j-1}}{2k} - \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2} = 0.$$

has a local truncation error of  $O(k^2) + O(h^2)$   
but it is unconditionally ~~unstable~~

---

# Analytical Treatment of Convergence (Direct proof).

Consider the eq.

$$\frac{\partial U}{\partial t} = \frac{\partial^2 U}{\partial x^2}, \quad 0 < x < 1, \quad (1)$$

where  $U$  is known for  $0 < x < 1$  when  $t=0$ , and at  $x=0$  and  $1$  when  $t>0$ .

The simplest explicit finite-difference approximation to (1) is

$$\frac{U_{i,j+1} - U_{i,j}}{K} = \frac{U_{i-1,j} - 2U_{i,j} + U_{i+1,j}}{h^2}. \quad (2)$$

Let  $\epsilon = U - u$ , and at the mesh points.

$$u_{i,j} = U_{i,j} - \epsilon_{i,j}, \quad U_{i,j+1} = U_{i,j+1} - \epsilon_{i,j+1}, \text{ etc.}$$

Substitution into (2) leads to

$$\begin{aligned} (U_{i,j+1} &= U_{i,j} - r\epsilon_{i-1,j} + (1-2r)\epsilon_{i,j} + r\epsilon_{i+1,j}) \\ (U_{i,j+1} - \epsilon_{i,j+1} &= r(U_{i,j+1} - \epsilon_{i,j+1}) + (1-2r)(U_{i,j} - \epsilon_{i,j}) + r(U_{i+1,j} - \epsilon_{i+1,j}) \\ \Rightarrow \epsilon_{i,j+1} &= r\epsilon_{i-1,j} + (1-2r)\epsilon_{i,j} + r\epsilon_{i+1,j} + U_{i,j+1} - U_{i,j} \\ &\quad + r(2U_{i,j} - U_{i-1,j} - U_{i+1,j}) \end{aligned} \quad (3)$$

By Taylor's theorem

$$\begin{aligned} U_{i+1,j} &= U(x_{i+h}, t_j) = U_{i,j} + h \left( \frac{\partial U}{\partial x} \right)_{i,j} + \frac{h^2}{2} \frac{\partial^2 U(x_{i+\theta_1 h}, t_j)}{\partial x^2}, \\ U_{i-1,j} &= U(x_{i-h}, t_j) = U_{i,j} - h \left( \frac{\partial U}{\partial x} \right)_{i,j} + \frac{h^2}{2} \frac{\partial^2 U(x_{i-\theta_2 h}, t_j)}{\partial x^2}, \\ U_{i,j+1} &= U(x_i, t_j + K) = U_{i,j} + K \frac{\partial U(x_i, t_j + \theta_3 K)}{\partial t}, \end{aligned}$$

where  $0 < \theta_1 < 1$ ,  $0 < \theta_2 < 1$  and  $0 < \theta_3 < 1$ .

$$\Rightarrow e_{i,j+1} = r e_{i-1,j} + (1-2r) e_{ij} + r e_{i+1,j} + \\ + K \left\{ \frac{\partial U(x_i, t_j + \theta_3 h)}{\partial t} - \frac{\partial^2 U(x_i + \theta_4 h, t_j)}{\partial x^2} \right\}, \quad (4)$$

where  $-1 < \theta_4 < 1$ .

Let  $E_j = \max_i |e_{ij}|$ ,  $M = \max_{ij} \{ \dots \}$ . When

$$r \leq \frac{1}{2}, \Rightarrow 1-2r > 0, \text{ so}$$

$$|e_{i,j+1}| \leq r |e_{i-1,j}| + (1-2r) |e_{ij}| + r |e_{i+1,j}| + KM$$

$$\leq E_j + KM \quad \text{for all values of } i \Rightarrow \text{this is true}$$

for  $\max_i |e_{ij}|$ .

$$\text{Hence } E_{j+1} \leq E_j + KM \leq (E_{j-1} + KM) + KM \leq E_{j-2} + 3KM$$

$$\leq E_{j-3} + 4KM \leq \dots \leq$$

$$\leq E_0 + jKM = fM,$$

$E_j \leq E_{j-1} + KM \leq E_{j-2} + 2KM \leq E_0 + jKM = fM$ ,  
because the initial values for  $U$  and  $V$  are the same, i.e.,  $E_0 = 0$ ,

when  $h$  tends to zero,  $K = rh^2$  also tends to zero and  $M$

tends to

$$\left( \frac{\partial U}{\partial t} - \frac{\partial^2 U}{\partial x^2} \right)_{ij} \Rightarrow 0$$

$\Rightarrow E_j \rightarrow 0$ . As  $|U_{ij} - U_{ij}| \leq E_j$ , this proves that  
 $U$  converges to  $U$  as  $h$  tends to zero when  $r \leq \frac{1}{2}$  and  $t$  is finite.

~~when  $r > \frac{1}{2}$  it can be shown that the~~

# Analytic treatment of stability

There are two standard ways of investigating the convergence of the solution of the finite-difference equation.

① Matrix method

② The Fourier series method (von Neumann's method)

## Matrix method

Explicit finite-difference scheme.

Consider the equation

$$\frac{\partial U}{\partial t} = \frac{\partial^2 U}{\partial x^2}, \quad 0 < x < 1,$$

where  $U=0$  at  $x=0$  and  $t=0$ ,  $t>0$ , and  $U$  is known when  $t=0$ .

The explicit finite-difference approximation

$$\frac{U_{i,j+1} - U_{i,j}}{\Delta t} = \frac{U_{i+1,j} - 2U_{i,j} + U_{i-1,j}}{\Delta x^2} \quad (1)$$

$$\Rightarrow U_{i,j+1} = rU_{i-1,j} + (1-2r)U_{i,j} + rU_{i+1,j}$$

leads to the eqs

$$i=1, \quad U_{1,j+1} = 0 + (1-2r)U_{1,j} + rU_{2,j}$$

$$i=2, \quad U_{2,j+1} = rU_{1,j} + (1-2r)U_{2,j} + rU_{3,j}$$

$$i=N-1, \quad U_{N-1,j+1} = rU_{N-2,j} + (1-2r)U_{N-1,j} + 0,$$

where  $N\Delta x = 1$ , and  $U_{0,j} = U_{N,j} = 0$ . These can be written as

$$\begin{bmatrix} U_{1,j+1} \\ U_{2,j+1} \\ U_{3,j+1} \\ \vdots \\ U_{N-1,j+1} \end{bmatrix} = \begin{bmatrix} (1-2r) & r \\ r & (1-2r) & r \\ & r & (1-2r) & r \\ & & & \ddots \\ & & & & r & (1-2r) \end{bmatrix} \begin{bmatrix} U_{1,j} \\ U_{2,j} \\ U_{3,j} \\ \vdots \\ U_{N-1,j} \end{bmatrix}$$

$\uparrow$   
 $U_{j+1}$

A

$\downarrow$   
 $U_j$

or as

$$\vec{U}_{j+1} = A \vec{U}_j.$$

Hence

$$\vec{U}_j = A \vec{U}_{j-1} = A(A \vec{U}_{j-2}) = A^2 \vec{U}_{j-2} = A^3 \vec{U}_{j-3} = \dots = A^{j-1} \vec{U}_0,$$

where  $\vec{U}_0$  is the vector of initial values.

Now suppose we introduce errors at every pivotal point along  $t=0$  and start the computation with the vector of values  $\vec{U}_0^*$  instead of  $\vec{U}_0$ . We shall then calculate

$$\vec{U}_1^* = A \vec{U}_0^*, \quad \vec{U}_2^* = A \vec{U}_1^* = A^2 \vec{U}_0^*,$$

and finally

$$\vec{U}_j^* = A^j \vec{U}_0^*,$$

assuming we introduce no further errors.

Define the error vector  $\vec{\epsilon}$  by

$$\vec{\epsilon} = \vec{U} - \vec{U}^*.$$

Then

$$\vec{\epsilon}_j = \vec{U}_j - \vec{U}_j^* = A^j (\vec{U}_0 - \vec{U}_0^*) = A^j \vec{\epsilon}_0,$$

showing that the formula for the propagation of the errors is the same as that for the calculation of  $\vec{U}$ .

The finite difference scheme will be stable when  $\vec{\epsilon}_j$  remains bounded as  $j$  increases indefinitely. This can always be investigated by expressing the error vector in terms of the eigenvectors of  $A$ .

Since  $A$  is real and symmetric, the eigenvalues  $\lambda$  are all distinct. Then  $A$  has  $(N-1)$  LI eigenvectors  $\vec{v}_s$ . So these eigenvectors can be used as a basis for our  $(N-1)$ -dimensional vector space and the error  $\vec{\epsilon}_0$ , with its  $(N-1)$  components, can be expressed uniquely as a linear combination of them, namely,

$$\vec{e}_0 = \sum_{s=1}^{N-1} c_s \vec{v}_s$$

~~88~~ - 89 -

where the  $c_s$ ,  $s=1(1)(N-1)$ , are known scalars.

The errors along the time-level  $t=k$ , resulting from the initial errors  $\vec{e}_0$  will be given by

$$\vec{e}_1 = A\vec{e}_0 = A \sum c_s \vec{v}_s = \sum c_s A \vec{v}_s = \sum c_s \lambda_s \vec{v}_s$$

Similarly,  $(\vec{e}_j = A \sum c_s \lambda_s \vec{v}_s = \sum c_s \lambda_s A \vec{v}_s = \sum c_s \lambda_s^j \vec{v}_s)$

$$\vec{e}_j = \sum c_s \lambda_s^j \vec{v}_s \quad (2)$$

This shows that the errors will not increase exponentially with  $j$  provided the eigenvalue with the largest modulus has a modulus less than or equal to unity.

The matrix  $A$  can be written as

$$\left[ \begin{array}{ccc|cc} (1-2r) & r & & & \\ r & (1-2r) & r & & \\ & & \ddots & \ddots & \\ & & & r & (1-2r) & r \\ & & & r & (1-2r) & \\ & & & & r & (1-2r) \end{array} \right] = \left[ \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \ddots & \vdots \end{array} \right] + r \left[ \begin{array}{ccc} -2 & 1 & 0 \\ 1 & -2 & 1 \\ 0 & 1 & -2 \\ \vdots & \ddots & \vdots \\ 0 & 0 & 1 \end{array} \right]$$

i.e., as

$$A = I + r T_{N-1}$$

where  $T_{N-1}$  is an  $(N-1) \times (N-1)$  matrix whose eigenvalues  $\lambda_s$  and eigenvectors  $\vec{v}_s$  are given by

$$\lambda_s = -4 \sin^2 \left( \frac{s\pi}{2N} \right), \quad (s=1, 2, \dots, N-1),$$

$$\vec{v}_s = \left( \sin \frac{s\pi}{N}, \sin \frac{2s\pi}{N}, \dots, \sin \frac{(N-1)s\pi}{N} \right)$$

The eigenvalues and vectors of a common tridiagonal matrix.

Let

$$A = \begin{bmatrix} a & b & & & \\ c & a & b & & \\ & c & a & b & \\ & & \ddots & & \\ & & & c & a \end{bmatrix}$$

be a square matrix of order  $N$ , where  $a, b$ , and  $c$  may be real or complex numbers.

$$A\vec{v} = \lambda \vec{v}, \quad \vec{v} = (v_1, v_2, \dots, v_N)$$

$$\begin{bmatrix} ab & & & & \\ c & a & b & & \\ & c & a & b & \\ & & \ddots & & \\ & & & c & a \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_N \end{bmatrix} = \lambda \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_N \end{bmatrix} \Rightarrow \begin{array}{lcl} av_1 + bv_2 & = & \lambda v_1 \\ cv_1 + av_2 + bv_3 & = & \lambda v_2 \\ cv_2 + av_3 + bv_4 & = & \lambda v_3 \\ \vdots & & \vdots \\ cv_{N-1} + av_N & = & \lambda v_N \end{array}$$

$$(a-\lambda)v_1 + bv_2 = 0$$

$$cv_1 + (a-\lambda)v_2 + bv_3 = 0$$

.....

$$cv_{j-1} + (a-\lambda)v_j + bv_{j+1} = 0$$

$$cv_{N-1} + (a-\lambda)v_N = 0$$

If we define  $V_0 = V_{N+1} = 0$  then these  $N$  eq.s can be combined into the single difference equation.

$$cV_{j-1} + (a-\lambda)V_j + bV_{j+1} = 0, \quad j=1(1)N. \quad (3)$$

$$V_j = m^j \neq 0$$

$$\begin{aligned} c m^{j-1} + (a-\lambda)m^j + b m^{j+1} &= 0 \\ \Rightarrow c + (a-\lambda)m + b m^2 &= 0 \end{aligned} \quad (3')$$

If  $m_1 \neq m_2$  (are real), then

$$V_j = B m_1^j + C m_2^j$$

is the solution of (3'). Since  $V_0 = V_{N+1} = 0$ , then

$$\begin{cases} 0 = B + C \rightarrow B = -C \\ 0 = B m_1^{N+1} + C m_2^{N+1} \Rightarrow 0 = B m_1^{N+1} - B m_2^{N+1} \end{cases}$$

Hence

$$\left(\frac{m_1}{m_2}\right)^{N+1} = 1 = e^{i\frac{2s\pi}{N+1}}, \quad s=1(1)N, \quad i=\sqrt{-1}.$$

$$\Rightarrow \frac{m_1}{m_2} = e^{\frac{i2s\pi}{N+1}}. \quad (4)$$

By eq. (3')

$$m_1 m_2 = \frac{c}{b}, \quad (5)$$

and elimination of  $m_2$  between (4) and (5) leads to

$$m_2 = \frac{c}{b m_1}, \quad \frac{m_1}{c} = e^{\frac{i2s\pi}{N+1}} \Rightarrow m_1^2 = \frac{c}{b} e^{\frac{i2s\pi}{N+1}}$$

$$\Rightarrow m_1 = \left(\frac{c}{b}\right)^{\frac{1}{2}} e^{\frac{i s \pi}{N+1}}.$$

$$m_2 = \frac{c}{b \left(\frac{c}{b}\right)^{\frac{1}{2}} e^{\frac{i s \pi}{N+1}}} = \left(\frac{c}{b}\right)^{\frac{1}{2}} e^{-\frac{i s \pi}{N+1}}.$$

-99-

Again, by (3'),

$$m_1 + m_2 = \frac{\lambda - a}{b} \Rightarrow \lambda = a + b(m_1 + m_2)$$

$$\Rightarrow \lambda = a + b \left(\frac{c}{b}\right)^{\frac{1}{2}} \left(e^{\frac{i s \pi}{N+1}} + e^{-\frac{i s \pi}{N+1}}\right)$$

Hence the  $N$  eigenvalues are given by

$$\lambda_s = a + 2b \left(\frac{c}{b}\right)^{\frac{1}{2}} \cos \frac{s\pi}{N+1}, \quad s=1(1)N$$

$$\boxed{\lambda_s = a + 2\sqrt{bc} \cos \frac{s\pi}{N+1}}$$

$s=1(1)N$

The  $j$ -th component of the eigenvector is

$$v_j = B m_1^j + C m_2^j = B \left(\frac{c}{b}\right)^{\frac{1}{2}j} \left(e^{\frac{is\pi}{N+1}} - e^{-\frac{is\pi}{N+1}}\right)$$

$$= 2iB \left(\frac{c}{b}\right)^{\frac{1}{2}j} \sin \frac{s\pi}{N+1}$$

so the eigenvector  $\vec{v}_s$  corresponding to  $\lambda_s$  can be taken as

$$\vec{v}_s^T = \left\{ \left(\frac{c}{b}\right)^{\frac{1}{2}} \sin \frac{s\pi}{N+1}, \frac{c}{b} \sin \frac{2s\pi}{N+1}, \left(\frac{c}{b}\right)^{\frac{3}{2}} \sin \frac{3s\pi}{N+1}, \dots, \left(\frac{c}{b}\right)^{\frac{N}{2}} \sin \frac{Ns\pi}{N+1} \right\}$$

It is easily shown that the roots (3') cannot be equal because if we assume  $m_1 = m_2$  the solution (3') is then

$$v_j = (B + Cj)m_1^j$$

and  $v_0 = v_{N+1} = 0 \Rightarrow B = C = 0$ , giving  $\vec{v} = 0$ , which is not possible.

$$a = -2, b = 1, c = 1$$

$$\Rightarrow \lambda_s = -2 + 2i \left(\frac{1}{1}\right)^{\frac{1}{2}} \cos \frac{s\pi}{N} = -2 \left(1 - \cos \frac{s\pi}{N}\right) = -4 \sin^2 \frac{s\pi}{2N}$$

$$\Rightarrow \boxed{\lambda_s = -4 \sin^2 \frac{s\pi}{2N}} \quad (s=1, 2, \dots, N-1)$$

$$\vec{v}_s = \left\{ \sin \frac{s\pi}{N}, \sin \frac{2s\pi}{N}, \dots, \sin \frac{(N-1)s\pi}{N} \right\}$$

## A note on eigenvalues and eigenvectors

Let  $\vec{x}$  be an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda$ . Then  $A\vec{x} = \lambda\vec{x}$ . Hence

$$A(A\vec{x}) = A^2\vec{x} = \lambda A\vec{x} = \lambda^2\vec{x}.$$

$\Rightarrow A^2$  has an eigenvalue  $\lambda^2$  corresponding to the eigenvector  $\vec{x}$ . Similarly,  $A^p\vec{x} = \lambda^p\vec{x}$ ,  $p = 3, 4, \dots$ .

(i) If  $f(A) = a_p A^p + a_{p-1} A^{p-1} + \dots + a_0 I$  is a polynomial in  $A$  with scalar coefficients  $a_p, \dots, a_0$ , then

$$f(A)\vec{x} = (a_p\lambda^p + \dots + a_0) \vec{x} = f(\lambda)\vec{x}$$

$\Rightarrow f(A)$  has an eigenvalue  $f(\lambda)$  corresponding to the eigenvector  $\vec{x}$ .

(ii) The eigenvalue of  $[f_1(A)]^{-1} f_2(A)$  corresponding to the eigenvector  $\vec{x}$  is  $\frac{f_2(\lambda)}{f_1(\lambda)}$ , where  $f_1(A)$  and  $f_2(A)$  are polynomials in  $A$ .

Proof       $f_1(A)\vec{x} = f_1(\lambda)\vec{x}$   
 $f_2(A)\vec{x} = f_2(\lambda)\vec{x}$

Multiply both eq.s by  $[f_1(A)]^{-1}$  and write  $\Leftrightarrow$

$$[f_1(A)]^{-1}\vec{x} = \frac{\vec{x}}{f_1(\lambda)} \quad \text{and} \quad [f_1(A)]^{-1}f_2(A)\vec{x} = f_2(\lambda)[f_1(A)]^{-1}\vec{x}$$

$$\Rightarrow [f_1(A)]^{-1}f_2(A)\vec{x} = f_2(\lambda) \cdot \frac{1}{f_1(\lambda)} \vec{x} = \frac{f_2(\lambda)}{f_1(\lambda)} \vec{x}$$

Hence the eigenvalues of  $A$ , as shown later are

$$1 + r \left\{ -4 \sin^2 \frac{\pi}{2N} \right\}.$$

Therefore the condition for stability of the explicit scheme is

$$\left| 1 - 4r \sin^2 \frac{\pi}{2N} \right| \leq 1.$$

$$\Rightarrow -1 \leq 1 - 4r \sin^2 \frac{\pi}{2N} \leq 1,$$

$$\Rightarrow -1 \leq 1 - 4r \sin^2 \frac{\pi}{2N} \Rightarrow 4r \sin^2 \frac{\pi}{2N} \leq 2$$

$$\Rightarrow r \leq \frac{1}{2 \sin^2 \frac{\pi}{2N}} > \frac{1}{2}$$

proving that the scheme is stable for  $r \leq \frac{1}{2}$ .

Crank-Nicolson method for  $\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}$

$$\frac{u_{i,j+1} - u_{i,j}}{\delta t} = \frac{1}{2} \left\{ \frac{u_{i+1,j+1} - 2u_{i,j+1} + u_{i-1,j+1}}{(\delta x)^2} + \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{(\delta x)^2} \right\}_{i-w+1}^{i+w+1} \quad \begin{matrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix} \quad \begin{matrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix} \quad \begin{matrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix} \quad \begin{matrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix}$$

$$\Rightarrow 2u_{i,j+1} - 2u_{i,j} = r(u_{i+1,j+1} - 2u_{i,j+1} + u_{i-1,j+1}) + r(u_{i+1,j} - 2u_{i,j} + u_{i-1,j}) \\ - r u_{i-1,j+1} + (2+2r)u_{i,j+1} - ru_{i+1,j+1} = ru_{i,j} + (2-2r)u_{i,j} + r u_{i+1,j}$$

$$\Rightarrow \begin{bmatrix} (2+2r) & -r & & & & & \\ -r & (2+2r) & -r & & & & \\ & & & & & & \\ & & & & (2+2r) & -r & & \\ & & & & -r & (2+2r) & -r & \\ & & & & -r & -r & (2+2r) & \\ & & & & & & & \end{bmatrix} \begin{bmatrix} u_{1,j+1} \\ u_{2,j+1} \\ \vdots \\ u_{N-1,j+1} \end{bmatrix} = \begin{bmatrix} u_{i,j} \\ u_{i,j} \\ \vdots \\ u_{i,j} \end{bmatrix}$$

$$= \begin{bmatrix} (2-2r) & r & & & & & \\ r & (2-2r) & r & & & & \\ & & & & & & \\ & & & & r & (2-2r) & r & \\ & & & & r & -r & (2-2r) & \\ & & & & r & r & -r & \end{bmatrix} \begin{bmatrix} u_{i,j} \\ u_{i,j} \\ \vdots \\ u_{i,j} \end{bmatrix}$$

$$\Rightarrow (2I - rT_{N-1})\vec{U}_{j+1} = (2I + rT_{N-1})\vec{U}_j$$

$$\Rightarrow \vec{U}_{j+1} = (2I - rT_{N-1})^{-1}(2I + rT_{N-1})\vec{U}_j$$

By the previous argument these finite-difference eq.s will be stable when the moduli of the eigenvalues of

95-

$$A = (2I - rT_{N-1})^{-1} (2I + rT_{N-1})$$

are each less than one. As the eigenvalues of  $T_{N-1}$  are  $-4\sin^2\left(\frac{s\pi}{2N}\right)$ , the eigenvalues of A are

$$\frac{2 - 4r\sin^2\left(\frac{s\pi}{2N}\right)}{2 + 4r\sin^2\left(\frac{s\pi}{2N}\right)}, \quad s=1, 2, \dots, N-1$$

and these are clearly less than one for all positive values of r proving that the Crank-Nicholson eq.s have unrestricted stability.

Useful theorems on bounds for eigenvalues

### Gershgorin's first theorem

The largest of the moduli of the eigenvalues of the square matrix  $A$  cannot exceed the largest sum of the moduli of the elements along any row or any column.

Proof.

Let  $\lambda_i$  be an eigenvalue of the  $N \times N$  matrix  $A$ , and  $\vec{x}_i$  the corresponding eigenvector with components  $v_1, v_2, \dots, v_n$ . Then the eq.

$$A \vec{x}_i = \lambda_i \vec{x}_i$$

In detail, is

$$a_{11}v_1 + a_{12}v_2 + \dots + a_{1n}v_n = \lambda_i v_1,$$

$$a_{21}v_1 + a_{22}v_2 + \dots + a_{2n}v_n = \lambda_i v_2,$$

⋮

$$a_{s1}v_1 + a_{s2}v_2 + \dots + a_{sn}v_n = \lambda_i v_s,$$

Let  $v_s$  be the largest in modulus of  $v_1, v_2, \dots, v_n$ . Select the  $s$ -th equation and divide by  $v_s$ , giving

$$\lambda_i = a_{s1}\left(\frac{v_1}{v_s}\right) + a_{s2}\left(\frac{v_2}{v_s}\right) + \dots + a_{sn}\left(\frac{v_n}{v_s}\right).$$

Therefore

$$|\lambda_i| \leq |a_{s1}| + |a_{s2}| + \dots + |a_{sn}|.$$

because

$$\left| \frac{v_i}{v_s} \right| \leq 1, \quad i=1, 2, \dots, n.$$

In particular this holds for  $|\lambda_i| = \max_s |\lambda_{s,i}|$ ,  $s=1(1)N$ .

Since the eigenvalues of the transpose of  $A$  are the same as those of  $A$  the theorem is also true for columns.

### Gershgorin's circle theorem or Brauer's theorem

Let  $P_s$  be the sum of the moduli of the elements along the  $s$ -th row excluding the diagonal element  $a_{s,s}$ .

Then each eigenvalue of  $A$  lies inside or on the boundary of at least one of the circles  $|\lambda - a_{s,s}| = P_s$ .

Proof

Since

$$\lambda_i = a_{s,1} \left( \frac{v_1}{v_s} \right) + a_{s,2} \left( \frac{v_2}{v_s} \right) + \dots + a_{s,s} + \dots + a_{s,n} \left( \frac{v_n}{v_s} \right),$$

Hence

$$|\lambda_i - a_{s,s}| = \left| a_{s,1} \left( \frac{v_1}{v_s} \right) + \dots + 0 + \dots + a_{s,n} \left( \frac{v_n}{v_s} \right) \right| \\ \leq |a_{s,1}| + |a_{s,2}| + \dots + 0 + \dots + |a_{s,n}| = P_s$$

which completes the proof.

As an illustrative example consider the Crank-Nicolson eq.s.

$$(2I - rT_{N-1})\vec{U}_{j+1} = (2I + rT_{N-1})\vec{U}_j = \{4I - (2I - rT_{N-1})\vec{U}_j$$

which can be written as

$$B \vec{U}_{j+1} = (4I - B) \vec{U}_j ,$$

giving

$$\vec{U}_{j+1} = (4B^{-1} - I) \vec{U}_j ,$$

where

$$B = \begin{bmatrix} (2+2r) & -r & & & \\ -r & (2+2r) & -r & & \\ & & \ddots & \ddots & \\ & & -r & (2+2r) & -r \\ & & & -r & (2+2r) \end{bmatrix}$$

The finite-difference eqs will be stable when the modulus of every eigenvalue of  $(4B^{-1} - I)$  does not exceed one, that is, when

$$\left| \frac{4}{\lambda} - 1 \right| \leq 1 ,$$

where  $\lambda$  is an eigenvalue of  $B$ .

$$\Rightarrow -1 \leq \frac{4}{\lambda} - 1 \leq 1 \Rightarrow \underline{\lambda \geq 2} .$$

For the matrix  $B$ ,  $a_{ss} = 2+2r$ ,  $\max p_s = 2r$ , so Brauer's theorem leads to

$$|\lambda - 2 - 2r| \leq 2r \Rightarrow -2r \leq \lambda - 2 - 2r \leq 2r$$

$$\Rightarrow 2 \leq \lambda \leq 2+4r \Rightarrow \underline{\lambda \geq 2} \text{ for all value of } r.$$

-99-

## Stability criteria for derivative boundary conditions

Consider the eq.

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}, \quad 0 < x < 1,$$

and the cond.s

$$\frac{\partial u}{\partial x} = A_1(u - b_1) \text{ at } x=0, t>0,$$

$$\frac{\partial u}{\partial x} = -A_2(u - b_2) \text{ at } x=1, t>0.$$

where  $A_1, A_2, b_1, b_2$  are constants,  $A_1 > 0, A_2 > 0$ .

When the b.c.s are approximated by the central-difference

eq.s

$$\frac{u_{i+1,j} - u_{i-1,j}}{2\delta x} = A_1(u_{0,j} - b_1), \quad \Rightarrow u_{-1,j} = u_{0,j} - 2A_1\delta x u_{0,j} + 2A_1\delta x b_1$$

$$\frac{u_{N+1,j} - u_{N-1,j}}{2\delta x} = -A_2(u_{N,j} - b_2). \quad (N\delta x = 1), \quad \Rightarrow u_{N+1,j} = u_{N-1,j} - 2A_2\delta x u_{N,j} + 2A_2\delta x b_2$$

and the differential eq. by the explicit scheme

$$u_{i,j+1} = r u_{i,j} + (1-2r) u_{i,j} + r u_{i+1,j}$$

elimination of  $u_{-1,j}, u_{N+1,j}$ , leads to the eq.s

$$u_{0,j+1} = r u_{1,j} + (1-2r) u_{0,j} + r u_{1,j}, \quad \Rightarrow u_{-1,j} = u_{0,j} - 2\delta x h_1 u_{0,j} + 2\delta x h_1 v_1$$

$$\Rightarrow u_{0,j+1} = r u_{1,j} - 2h_1 \delta x u_{0,j} + 2h_1 \delta x v_1 + (1-2r) u_{0,j} + r u_{1,j}$$

$$\boxed{u_{0,j+1} = 2r u_{1,j} + (1-2r(1+h_1 \delta x)) u_{0,j} + 2h_1 \delta x v_1}$$

$$i=N:$$

$$u_{N,j+1} = r u_{N-1,j} + (1-2r) u_{N,j} + r u_{N+1,j}$$

$$u_{N,j+1} = r u_{N-1,j} - 2h_2 \delta x u_{N,j} + 2h_2 \delta x v_2 + r u_{N+1,j} - (1-2r) u_{N,j} + r u_{N+1,j}$$

$$\boxed{u_{N,j+1} = 2r u_{N-1,j} + (1-2r(1+h_2 \delta x)) u_{N,j} + 2r h_2 \delta x v_2}$$

~~100~~

$$\Rightarrow \begin{bmatrix} u_{0,j+1} \\ u_{1,j+1} \\ u_{2,j+1} \\ \vdots \\ u_{N-1,j+1} \\ u_{N,j+1} \end{bmatrix} = \begin{bmatrix} f_1 - 2r(1 + h_1 \frac{\partial}{\partial x}) \\ r \\ r \\ \vdots \\ r \\ 2r \end{bmatrix} \quad \begin{bmatrix} 2r \\ (1-2r) \\ r \\ \vdots \\ r \\ (1-2r) \\ r \\ 2r \end{bmatrix} \quad \begin{bmatrix} f_1 - 2r(1 + h_2 \frac{\partial}{\partial x}) \end{bmatrix}$$

$$x \begin{bmatrix} u_{0,j} \\ u_{1,j} \\ u_{2,j} \\ \vdots \\ u_{N-1,j} \\ u_{N,j} \end{bmatrix} + \begin{bmatrix} a_1 b_1 h \\ 2rh, v_1, \delta x \\ 0 \\ 0 \\ \vdots \\ 0 \\ 2rh_1 v_2 \frac{\partial}{\partial x} \end{bmatrix}$$

As each component of the last column vector is a constant  
the matrix determining the propagation of the error is

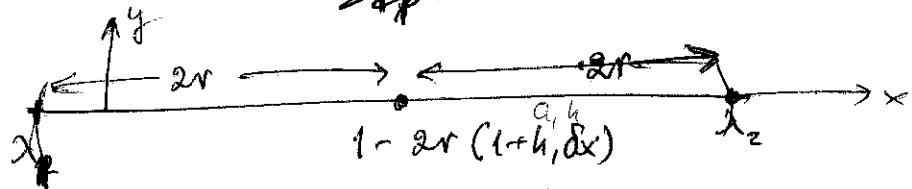
$$\begin{bmatrix} f_1 - 2r(1 + h_1 \frac{\partial}{\partial x}) & 2r \\ r & (1-2r) & r \\ \vdots & \vdots & \vdots \\ r & (1-2r) & r \\ 2r & \{f_1 - 2r(1 + h_2 \frac{\partial}{\partial x})\} \end{bmatrix}$$

Application of Brauner's Theorem to this matrix, with

$$Q_{ss} = 1 - 2r(1 + h_1 \frac{\partial}{\partial x}) \text{ and } P_s = 2r,$$

shows that its eigenvalues  $\lambda$  lie on or within the circle

$$|\lambda - \{1 - 2r(1 + h_1 \frac{\partial}{\partial x})\}| \leq 2r$$



$$-2r \leq \lambda - \{1 - 2r(1 + h, \delta x)\} \leq 2r$$

$$\begin{aligned} 1 - 2r - 2rh, \delta x - 2r &\leq \lambda \leq 2r + 1 - 2r - 2rh, \delta x \\ (1 - 4r - 2rh, \delta x) &\leq \lambda \leq 2r - 2rh, \delta x \\ \Rightarrow \lambda_1 = 1 - 2r(2 + h, \delta x), \quad \lambda_2 = 1 - 2r h, \delta x & \end{aligned}$$

and for stability,

$$|\lambda_1| \leq 1, \quad |\lambda_2| \leq 1.$$

hence

$$\begin{aligned} -1 \leq 1 - 2r(2 + h, \delta x) &\leq 1, \text{ giving} \\ 2r(2 + h, \delta x) \leq 2 &\Rightarrow r \leq \frac{1}{2 + h, \delta x}, \end{aligned}$$

and

$$\begin{aligned} |1 - 2rh, \delta x| \leq 1 &\Rightarrow -1 \leq 1 - 2rh, \delta x \\ 2rh, \delta x \leq 2 &\Rightarrow r_1 \leq \frac{1}{h, \delta x} \end{aligned}$$

The least of these is

$$r \leq \frac{1}{2 + h, \delta x}.$$

Similarly we require

$$r \leq \frac{1}{2 + h_2 \delta x}.$$

For row 2(1)(N-1)  $a_{ss} = 1 - 2r$ ,  $P_s = 2r$ , giving

$$|\lambda - (1 - 2r)| \leq 2r$$

$$\Rightarrow -2r \leq \lambda - (1 - 2r) \leq 2r \Rightarrow -2r \leq \lambda - 1 + 2r \leq 2r$$

$$1 - 2r - 2r \leq \lambda \leq 2r + 1 - 2r$$

$$\underbrace{\lambda_1}_{\lambda_1} \leq \lambda \leq \underbrace{\lambda_2}_{\lambda_2}$$

$$|\lambda_1| \leq 1 \rightarrow |1 - 4r| \leq 1$$

$$\begin{aligned} -1 \leq 1 - 4r &\leq 1 \\ 4r \leq 2 &\Rightarrow r \leq \frac{1}{2} \end{aligned}$$

108

For overall stability,

$$r \leq \min \left\{ \frac{1}{2 + h_1 \delta x}, \frac{1}{2 + h_2 \delta x} \right\}$$

Crank-Nicolson equations

$$\begin{aligned} -rU_{i-1,j+1} + (2+2r)U_{i,j+1} - rU_{i+1,j+1} &= rU_{i-1,j} + (2-2r)U_{i,j} + rU_{i+1,j} \\ \left\{ \begin{array}{l} U_{i-1,j} = U_{i,j} - \frac{a_1 h}{a_1 h + b_1} \delta x U_{0,j} + \frac{a_1 h}{a_1 h + b_1} \delta x V_1, \\ U_{i+1,j+1} = U_{i,j+1} - \frac{a_1 h}{a_1 h + b_1} \delta x U_{0,j+1} + \frac{a_1 h}{a_1 h + b_1} \delta x V_1, \\ U_{n+1,j} = U_{n-1,j} - \frac{a_2 h}{a_2 h + b_2} \delta x U_{n,j} + \frac{a_2 h}{a_2 h + b_2} \delta x V_2, \\ U_{n+1,j+1} = U_{n-1,j+1} - \frac{a_2 h}{a_2 h + b_2} \delta x U_{n,j+1} + \frac{a_2 h}{a_2 h + b_2} \delta x V_2 \end{array} \right. \end{aligned}$$

$$\text{for } i=0: -rU_{0,j+1} + (2+2r)U_{0,j+1} - rU_{1,j+1} = rU_{0,j} + (2-2r)U_{0,j} + rU_{1,j}$$

$$-r\underline{U_{1,j+1}} + 2rh_1 \delta x \underline{U_{0,j+1}} - 2rh_1 \delta x \underline{V_1} + (2+2r)\underline{U_{0,j+1}} - r\underline{U_{1,j+1}}$$

$$= rU_{0,j} - 2rh_1 \delta x \underline{U_{0,j}} + 2rh_1 \delta x \underline{V_1} + (2-2r)\underline{U_{0,j}} + r\underline{U_{1,j}}$$

$$\boxed{\{2+2r(1+h_1 \delta x)\}U_{0,j+1} - 2rU_{1,j+1} = \{2-2r(1+h_1 \delta x)\}U_{0,j} + 2rU_{1,j} + 4rh_1 \delta x v_1}$$

i=N:

$$-rU_{N-1,j+1} + (2+2r)U_{N,j+1} - rU_{N+1,j+1} = rU_{N-1,j} + (2-2r)U_{N,j} + rU_{N+1,j}$$

$$-r\underline{U_{N-1,j+1}} + (2+2r)\underline{U_{N,j+1}} - r\underline{U_{N+1,j+1}} + 2rh_2 \delta x \underline{U_{N,j+1}} - 2rh_2 \delta x \underline{V_2}$$

$$= rU_{N-1,j} + (2-2r)U_{N,j} + rU_{N+1,j} - 2rh_2 \delta x \underline{U_{N,j}} + 2rh_2 \delta x \underline{V_2}$$

$$\Rightarrow \boxed{-2rU_{N-1,j+1} + (2+2r)\underline{U_{N,j+1}} + \{2+2r(1+h_2 \delta x)\}U_{N,j} =}$$

$$= 2rU_{N-1,j} + \{2-2r(1+h_2 \delta x)\}U_{N,j} + 4rh_2 \delta x V_2$$

$$\begin{array}{c}
 \overrightarrow{\longrightarrow} \\
 \left[ \begin{array}{ccc}
 \{2 + 2r(1+h, \delta x)\} & \overset{a_1 h}{\longrightarrow} & -2r \\
 -r & (2+2r) & -r \\
 -r & (2+2r) & -r \\
 \vdots & \vdots & \vdots \\
 -r & (2+2r) & -r \\
 2r & \{2 + 2r(1+h_2, \delta x)\} & \overset{a_2 h}{\longrightarrow}
 \end{array} \right]
 \end{array}$$

$$\times \begin{bmatrix} u_{0,j+1} \\ u_{1,j+1} \\ \vdots \\ u_{n-1,j+1} \\ u_{n,j+1} \end{bmatrix} = \begin{bmatrix}
 \{2 - 2r(1+h, \delta x)\} & \overset{a_1 h}{\longrightarrow} & +2r \\
 r & (2-2r) & r \\
 r & (2-2r) & r \\
 \vdots & \vdots & \vdots \\
 r & (2-2r) & r \\
 2r & \{2 - 2r(1+h_2, \delta x)\} & \overset{a_2 h}{\longrightarrow}
 \end{bmatrix}$$

$$\times \begin{bmatrix} u_{0,j} \\ u_{1,j} \\ \vdots \\ u_{n-1,j} \\ u_{n,j} \end{bmatrix} + \begin{bmatrix} 4rh, \delta x \\ 0 \\ \vdots \\ 0 \\ 4rh, \delta x \end{bmatrix} \Rightarrow \begin{array}{l}
 B\vec{e}_{j+1} = (4I - B)\vec{e}_j \\
 \vec{e}_{j+1} = (4B^{-1} - I)\vec{e}_j \\
 B\vec{U}_{j+1} = (4I - B)\vec{U}_j \Rightarrow \vec{U}_{j+1} = (4B^{-1} - I)\vec{U}_j + B\vec{e}_j \\
 \vec{e}_{j+1} = \{4B^{-1} - I\}\vec{e}_j
 \end{array}$$

$\boxed{|\frac{4}{\lambda} - 1| \leq 1}$   
 $-1 \leq \frac{4}{\lambda} - 1 \leq 1$   
 $0 \leq \frac{4}{\lambda} \leq 2$   
 $\boxed{\lambda \geq 2}$

By Gershgorin's circle theorem

$$|\lambda - \{2 + 2r(1+h, \delta x)\}| \leq 2r$$

$$-2r + 2 + 2r(1+h, \delta x) \leq \lambda \leq 2r + 2 + 2r(1+h, \delta x)$$

$$2 + 2rh, \delta x \leq \lambda \leq 2 + 4r + 2rh, \delta x$$

$$\Rightarrow \lambda \geq 2 + 2rh, \delta x \Rightarrow \boxed{\lambda \geq 2}$$

Similarly  $\lambda \geq 2$  and  $\lambda \geq 2 + 2rh, \delta x$  for the remaining rows.

Hence the equations are unconditionally stable.

The stability of three or more time-level  
difference eqs

The following theorem is useful.

Theorem. If the matrix  $A$  can be written as

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,n} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,n} \\ \vdots & \vdots & & \vdots \\ A_{m,1} & A_{m,2} & \cdots & A_{m,n} \end{bmatrix},$$

where each  $A_{i,j}$  is an  $N \times N$  matrix, and all the  $A_{i,j}$  have a common set of  $n$  LI eigenvectors, then the eigenvalues of  $A$  are given by the eigenvalues of the matrices

$$\begin{bmatrix} \lambda_{1,1}^{(k)} & \lambda_{1,2}^{(k)} & \cdots & \lambda_{1,n}^{(k)} \\ \lambda_{2,1}^{(k)} & \lambda_{2,2}^{(k)} & \cdots & \lambda_{2,n}^{(k)} \\ \vdots & \vdots & & \vdots \\ \lambda_{m,1}^{(k)} & \lambda_{m,2}^{(k)} & \cdots & \lambda_{m,n}^{(k)} \end{bmatrix}, \quad k=1(1)n,$$

where  $\lambda_{i,j}^{(k)}$  is the  $k$ -th eigenvalue of  $A_{i,j}$  corresponding to the  $k$ -th eigenvector  $\vec{v}_k$  common to all the  $A_{i,j}$ 's.

-405-

Example: Investigate the stability of the \_\_\_\_\_ and Crank-Nicolson approximation to the equation

$$\frac{\partial U}{\partial t} = \frac{\partial^2 U}{\partial x^2}, 0 < x < 1,$$

given that  $U$  is known initially for  $0 \leq x \leq 1$  and is known on the boundaries  $x=0$  and  $x=1$ . The approximation at the  $(ih, jk)$  is

$$\frac{1}{2K} (U_{i,j+1} - U_{i,j-1}) = \frac{1}{h^2} (U_{i-1,j} - (U_{i,j+1} + U_{i,j-1})) + \text{---}$$

which may be written as

$$U_{i,j+1} - U_{i,j-1} = 2r (U_{i-1,j} + U_{i+1,j}) - 2r U_{i,j} - \text{---}$$

$$(1+2r)U_{i,j+1} = 2r (U_{i-1,j} + U_{i+1,j}) + (1-2r)U_{i,j} - \text{---}$$

where  $r = \frac{K}{h^2}$ . For known boundary values and equations in matrix form are

$$(1+2r) \begin{bmatrix} U_{1,j+1} \\ U_{2,j+1} \\ U_{3,j+1} \\ \vdots \\ U_{N-1,j+1} \end{bmatrix} = 2r \begin{bmatrix} 0 & 1 & & & & & \\ 1 & 0 & 1 & & & & \\ & 1 & 0 & 1 & & & \\ & & & & \ddots & & \\ & & & & & 1 & \\ & & & & & & 0 \end{bmatrix} \begin{bmatrix} U_{0,j} \\ U_{1,j} \\ U_{2,j} \\ U_{3,j} \\ \vdots \\ U_{N,j} \end{bmatrix} + \text{---}$$

$$+ (1-2r) \begin{bmatrix} U_{1,j-1} \\ U_{2,j-1} \\ U_{3,j-1} \\ \vdots \\ U_{N-1,j-1} \end{bmatrix} + 2r \begin{bmatrix} U_{0,j} \\ 0 \\ 0 \\ \vdots \\ U_{N,j} \end{bmatrix}$$

giving

$$\vec{U}_{j+1} = \frac{2r}{1+2r} A \vec{U}_j + \frac{1-2r}{1+2r} \vec{U}_{j-1} + \vec{G}_j$$

-105-

where  $A$  is as displayed and  $\vec{C}_j$  is a vector of known values.

Put

$$\vec{v}_j = \begin{bmatrix} \vec{U}_j \\ \vec{U}_{j-1} \end{bmatrix}.$$

Then eq. (1) and the identity  $\vec{U}_j = \vec{U}_j$  can be written as

$$\begin{bmatrix} \vec{U}_{j+1} \\ \vec{U}_j \end{bmatrix} = \begin{bmatrix} \frac{2r}{1+2r} A & \frac{(1-2r)}{1+2r} I_{N-1} \\ I_{N-1} & 0 \end{bmatrix} \begin{bmatrix} \vec{U}_j \\ \vec{U}_{j-1} \end{bmatrix} + \begin{bmatrix} \vec{C}_j \\ 0 \end{bmatrix}$$

where  $I_{N-1}$  is the unit matrix of order  $(N-1)$ ,

i.e., as

$$\vec{V}_{j+1} = P\vec{V}_j + d_j,$$

where  $P$  is the matrix shown and  $d_j$  a column vector of known constants. This technique has reduced a three-level difference equation to a two-level one. The equations will be stable when each eigenvalue of  $P$  has a modulus  $\leq 1$ .

The matrix  $A$  has  $(N-1)$  different eigenvalues so it has  $(N-1)$  LI eigenvectors  $\vec{v}_s$ ,  $s=1(1)(N-1)$ . Although the matrix  $I_{N-1}$  has  $(N-1)$  eigenvalues each equal to 1 it has  $(N-1)$  LI eigenvectors which may be taken as  $\vec{v}_s$ ,  $s=1(1)(N-1)$ , ( $I_{N-1}\vec{v}_s = 1\vec{v}_s$ ). Hence the eigenvalues  $\lambda$  of  $P$  are the eigenvalues of the matrix

$$\begin{bmatrix} \frac{2r\lambda_k}{1+2r} & \frac{1-2r}{1+2r} \\ 1 & 0 \end{bmatrix},$$

$$\boxed{\lambda_s = a + 2b \left(\frac{c}{b}\right)^{\frac{1}{2}} \cos \frac{s\pi}{N+1}}$$

where  $\lambda_k$  is the  $k$ -th eigenvalue of  $A$ . ( $\lambda_k = 2 \cos\left(\frac{k\pi}{N}\right)$ ,  $k=1(1)(N-1)$ )

$$\Rightarrow \det \begin{bmatrix} \left\{ \frac{2r\lambda_k}{1+2r} - \lambda \right\} & \frac{1-2r}{1+2r} \\ 1 & -\lambda \end{bmatrix} = 0$$

~~-107-~~

$$\Rightarrow \lambda^2 - \frac{2r\lambda_{ic}}{1+2r} \lambda - \frac{1-2r}{1+2r} = 0.$$

Hence

$$\lambda^2 - \frac{2r \cdot 2C_0 \frac{\kappa\theta}{N}}{1+2r} \lambda - \frac{1-2r}{1+2r} = 0,$$

$$\lambda = \frac{2r C_0 \frac{\kappa\theta}{N}}{1+2r} \pm \sqrt{\frac{4r^2 C_0^2 \frac{\kappa\theta^2}{N}}{(1+2r)^2} + \frac{1-2r}{1+2r}}$$

$$= \frac{2r C_0 \frac{\kappa\theta}{N}}{1+2r} \pm \sqrt{\frac{4r^2 C_0^2 \frac{\kappa\theta^2}{N} + 1-4r^2}{(1+2r)^2}}$$

$$= \frac{2r C_0 \frac{\kappa\theta}{N}}{1+2r} \pm \frac{\sqrt{4r^2 (C_0^2 \frac{\kappa\theta^2}{N} - 1)} + 1}{1+2r}$$

$$= \frac{2r C_0 \frac{\kappa\theta}{N}}{1+2r} \pm \frac{\sqrt{1-4r^2 \sin^2 \frac{\kappa\theta}{N}}}{1+2r}$$

$$\lambda = \frac{2r C_0 \frac{\kappa\theta}{N}}{1+2r} \pm \sqrt{1-4r^2 \sin^2 \frac{\kappa\theta}{N}}$$

Case (i)  $1 > 1-4r^2 \sin^2 \frac{\kappa\theta}{N} \geq 0$

Then

$$|\lambda| < \frac{2r+1}{1+2r} = 1$$

Case (ii)  $1-4r^2 \sin^2 \frac{\kappa\theta}{N} \leq 0$ .

Then

$$|\lambda|^2 = \frac{1}{(2r+1)^2} \left\{ \left( 2r C_0 \frac{\kappa\theta}{N} \right)^2 + 4r^2 \sin^2 \frac{\kappa\theta}{N} - 1 \right\}$$

$$= \frac{4r^2 - 1}{4r^2 + 4r + 1} < 1 \text{ since } r > 0.$$

Therefore the equations are unconditionally stable for all positive  $r$ .

## Stability by the Fourier series method (von Neumann's method)

This method expresses

The idea of this method consists of that, an initial line of errors in terms of a finite Fourier series, and considers the growth of a function that reduces to this series for  $t=0$  by a "variables separable" method identical with that commonly used for deriving analytical solutions of PDEs. The Fourier series can be formulated in terms of sines or cosines but the algebra is easier if the complex exponential form is used i.e., with

$$\sum a_n \cos \frac{n\pi x}{l} \text{ or } \sum b_n \sin \frac{n\pi x}{l}$$

replaced by the equivalent

$$\sum A_n e^{i \frac{n\pi x}{l}}, \quad i = \sqrt{-1}, \quad l \text{ is the interval throughout which the function is defined.}$$

We use the notation  $U(ph, qk) = U_{p,q}$  (instead of  $u_{ij}$ ).

Then  $A_n e^{i \frac{n\pi x}{l}} = A_n e^{i \frac{n\pi ph}{Nh}} = A_n e^{i \frac{\beta_n ph}{h}}$ , where

$$\beta_n = \frac{n\pi}{Nh} \text{ and } Nh = l$$

Denote the errors at the pivotal points along  $t=0$ , between  $x=0$  and  $Nh$ , by  $E(ph) = E_p, p=0, 1, \dots, N$ .

Then the  $(N+1)$  equations

$$E_p = \sum_{n=0}^N A_n e^{i \beta_n ph}, \quad (p=0, 1, \dots, N),$$

are sufficient to determine the  $(N+1)$  unknowns  $A_0, A_1, \dots, A_N$  uniquely, showing that an arbitrary distribution of initial errors can be expressed in this exponential form.

To investigate the propagation of this error as  $t$  increases we need to find a solution of the finite-difference eq. which reduces to  $e^{iBph}$  when  $t = qk = 0$ .

→ 109

Assume

$$E_{p,q} = e^{i\beta x} e^{\alpha t} = e^{i\beta ph} e^{qK} = e^{i\beta ph} \xi^q, \quad (1)$$

where  $\xi = e^{qK}$ , and  $\alpha$ , in general, is a complex constant.  
This obviously reduces to  $e^{i\beta ph}$  when  $q=0$ . The error will not increase as  $t$  increases provided

$$|\xi| \leq 1.$$

for the stability

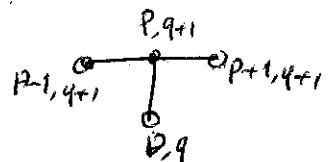
The criterion  $|\xi| \leq 1$  is necessary and sufficient for two time-level difference eqs. but is not always sufficient for three or more level eqs. although it is always necessary.

Example 1 Investigate the stability of the fully-implicit finite-difference equation

$$\frac{U_{p,q+1} - U_{p,q}}{k} = \frac{U_{p-1,q+1} - 2U_{p,q+1} + U_{p+1,q+1}}{h^2} \quad (2)$$

approximating the parabolic eq.

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



As the error function  $E_{p,q}$  satisfies the same difference equation as  $U_{p,q}$ , substitution of  $E_{p,q}$  from eq. (1) into eq. (2) gives

$$e^{i\beta ph} \xi^{q+1} - e^{i\beta ph} \xi^q = r \{ e^{i\beta(p-1)h} \xi^{q+1} - 2e^{i\beta ph} \xi^{q+1} + e^{i\beta(p+1)h} \xi^{q+1} \},$$

where  $r = k/h^2$ . Division by  $e^{i\beta ph} \xi^q$  leads to

$$\begin{aligned} \xi - 1 &= r \xi (e^{-i\beta h} - 2 + e^{i\beta h}) \\ &= r \xi (2 \cos \beta h - 2) = -4 r \xi \sin^2 \left( \frac{\beta h}{2} \right). \end{aligned}$$

Hence

$$\xi = \frac{1}{1 + 4r \sin^2 \left( \frac{\beta h}{2} \right)}.$$

For stability  $|\xi| \leq 1$ , which is so for all positive values of  $r$ .

Example 2. The hyperbolic equation

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

is approximated by the explicit scheme

$$\frac{U_{p,q+1} - 2U_{p,q} + U_{p,q-1}}{k^2} = \frac{U_{p+1,q} - 2U_{p,q} + U_{p-1,q}}{h^2}$$

(3)

Investigate its stability.

Putting  $E_{p,q} = e^{i\beta ph} \xi^q$ ,  $\xi = e^{ikx}$  in (3), we obtain

$$\xi^{-2} + \xi^{-1} = r^2 (e^{i\beta ph} - 2 + e^{-i\beta ph})$$

$$\Rightarrow \xi^2 - 2\xi + 1 = \xi r^2 (e^{i\beta ph} - 2) = -4\xi r^2 e^{i\beta ph}$$

$$\Rightarrow \underbrace{\xi^2 - 2(1 - 2r^2 \sin^2(\frac{\beta h}{2}))}_{A} \xi + 1 = 0, \quad r = \frac{k}{h},$$

$$\Rightarrow \xi^2 - 2A\xi + 1 = 0 \quad (4)$$

$$\Rightarrow \xi_1 = A + \sqrt{A^2 - 1} \quad \text{and} \quad \xi_2 = A - \sqrt{A^2 - 1}.$$

For stability

$$|\xi| \leq 1.$$

As  $r, k, \beta$  are real  $A \leq 1$

when  $A < -1$ ,  $|\xi_2| > 1$ , giving instability.

When  $-1 \leq A \leq 1$ ,  $A^2 \leq 1$ ,  $\xi_1 = A + i(1-A^2)^{\frac{1}{2}}$ ,  $\xi_2 = A - i(1-A^2)^{\frac{1}{2}}$ ,

Hence

$$|\xi_1| = |\xi_2| = \sqrt{A^2 + (1-A^2)} = 1,$$

providing that eq. (3) is stable for  $-1 \leq A \leq 1$ . By (4)

we then have  $-1 \leq 1 - 2r^2 \sin^2(\frac{\beta h}{2}) \leq 1 \Rightarrow -1 \leq 1 - 2r^2 \sin^2(\frac{\beta h}{2})$

$$\Rightarrow r^2 \leq \frac{1}{\sin^2(\frac{\beta h}{2})} \Rightarrow r \leq \frac{1}{\sin(\frac{\beta h}{2})} \Rightarrow r \leq 1$$

$$x_i^{(n+1)} = \frac{1}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\}, \quad i=1(1)m.$$

Jacobi method:

$$x_i^{(n+1)} = \frac{1}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\}, \quad i=1(1)m$$

If the G-S are written as

Gauss-Seidel method

$$\Rightarrow x_i^{(n+1)} = x_i^{(n)} + \frac{1}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\}$$

SOR (Successive over-relaxation method) :

$$x_i^{(n+1)} = x_i^{(n)} + \frac{\omega}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\}, \quad i=1(1)m$$

it may be written as

$$\begin{aligned} x_i^{(n+1)} &= x_i^{(n)} + \frac{\omega}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n+1)} - a_{ii} x_i^{(n)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\} \\ &= x_i^{(n)} - \omega x_i^{(n)} + \frac{\omega}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\} \end{aligned}$$

$$\Rightarrow \begin{aligned} x_i^{(n+1)} &= \frac{\omega}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(n+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(n)} \right\} + (1-\omega) x_i^{(n)} \\ &= \omega \left( \text{R.H.S. of the Gauss-Seidel iteration equation} \right) \\ &\quad - (1-\omega) x_i^{(n)}. \end{aligned}$$

The factor  $\omega$ , called the acceleration parameter or relaxation factor, generally lies in the range  $1 < \omega < 2$ .

$\omega=1 \rightarrow$  Gauss-Seidel iteration.

## Iterative methods in matrix form

we write A as

$$Ax = b \quad (1)$$

For  $4 \times 4$

$$A = D - L - U, \text{ where}$$

$$D = \begin{bmatrix} a_{11} & 0 & 0 & 0 \\ 0 & a_{22} & 0 & 0 \\ 0 & 0 & a_{33} & 0 \\ 0 & 0 & 0 & a_{44} \end{bmatrix}, \quad -L = \begin{bmatrix} 0 & 0 & 0 & 0 \\ a_{21} & 0 & 0 & 0 \\ a_{31} & a_{32} & 0 & 0 \\ a_{41} & a_{42} & a_{43} & 0 \end{bmatrix},$$

$$-U = \begin{bmatrix} 0 & a_{12} & a_{13} & a_{14} \\ 0 & 0 & a_{23} & a_{24} \\ 0 & 0 & 0 & a_{34} \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

The eq. (1) can then be written as

$$(D - L - U)x = b \Rightarrow Dx = (L + U)x + b$$

Jacobi iteration

$$\Rightarrow Dx^{(n+1)} = (L + U)x^{(n)} + b \quad (1)$$

$$\Rightarrow x^{(n+1)} = D^{-1}(L + U)x^{(n)} + D^{-1}b$$

The matrix  $D^{-1}(L + U)$  is called the ~~Jacobi~~ iteration matrix.

The Gauss-Seidel iteration is defined by

$$Dx^{(n+1)} = Lx^{(n+1)} + Ux^{(n)} + b \quad (2)$$

Hence

$$(D - L)x^{(n+1)} = Ux^{(n)} + b$$

giving that

$$x^{(n+1)} = (D - L)^{-1}Ux^{(n)} + (D - L)^{-1}b.$$

which shows that the Gauss-Seidel iteration matrix is  $(D - L)^{-1}U$ .

SOR :

$$x_i^{(u+1)} = x_i^{(u)} + \frac{\omega}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(u+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(u)} \right\}, \quad i=1(1)m.$$

or

$$x_i^{(u+1)} = \cancel{\frac{\omega}{a_{ii}}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(u+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(u)} \right\} - (\omega-1)x_i^{(u)}$$

$$= \omega \left( \text{R.H.S. of the Gauss-Seidel iteration eq.8} \right) \\ - (\omega-1)x_i^{(u)}.$$

The correction or displacement vector  $d^{(u)} = x^{(u+1)} - x^{(u)}$  of the SOR method is taken to be  $\omega$  times the displacement vector  $d_i^{(u)}$  defined by the Gauss-Seidel iteration.

$$\mathcal{D}d_i^{(u)} = \mathcal{D}(x^{(u+1)} - x^{(u)}) = Lx^{(u+1)} + Ux^{(u)} + b - \mathcal{D}x^{(u)}$$

Hence the SOR iteration, defined by

$$d^{(u)} = \omega d_i^{(u)}$$

can be written as

$$x^{(u+1)} - x^{(u)} = \omega \mathcal{D}^{-1} (Lx^{(u+1)} + Ux^{(u)} + b - \mathcal{D}x^{(u)}).$$

Therefore

$$(I - \omega \mathcal{D}^{-1} L)x^{(u+1)} = x^{(u)} + \omega \mathcal{D}^{-1} Ux^{(u)} + \omega \mathcal{D}^{-1} b - \omega x^{(u)} \\ = \{(1-\omega)I + \omega \mathcal{D}^{-1} U\}x^{(u)} + \omega \mathcal{D}^{-1} b.$$

Hence

$$x^{(u+1)} = (I - \omega \mathcal{D}^{-1} L)^{-1} \{(1-\omega)I + \omega \mathcal{D}^{-1} U\}x^{(u)} + (I - \omega \mathcal{D}^{-1} L)^{-1} \omega \mathcal{D}^{-1} b \quad (3)$$

showing the ~~the~~ SOR iteration matrix is

$$(I - \omega \mathcal{D}^{-1} L)^{-1} \{(1-\omega)I + \omega \mathcal{D}^{-1} U\}.$$

A necessary and sufficient condition  
for the convergence iterative methods

Each of the three iterative methods described  
above can be written as

$$\mathbf{x}^{(n+1)} = \mathbf{G}\mathbf{x}^{(n)} + \mathbf{c}, \quad (4)$$

where  $\mathbf{G}$  is the iteration matrix and  $\mathbf{c}$  a column  
vector of known values. This equation was derived from  
the original equations by rearranging them into the form

$$\mathbf{x} = \mathbf{f}\mathbf{x} + \mathbf{e}, \quad (5)$$

i.e., the unique solution of the  $m$  linear eq.s  $\mathbf{A}\mathbf{x} = \mathbf{b}$

is the solution of equation (5). The error  $\mathbf{e}^{(n)}$   
in the  $n$ th approximation to the exact solution is  
defined by  $\mathbf{e}^{(n)} = \mathbf{x} - \mathbf{x}^{(n)}$ :

$$\mathbf{e}^{(n+1)} = \mathbf{G}\mathbf{e}^{(n)}.$$

Therefore

$$\mathbf{e}^{(n)} = \mathbf{G}\mathbf{e}^{(n-1)} = \mathbf{G}^2\mathbf{e}^{(n-2)} = \dots = \mathbf{G}^n\mathbf{e}^{(0)}. \quad (6)$$

The sequence of iterative values  $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}, \dots$  will  
converge to  $\mathbf{x}$  as  $n$  tends to infinity if

$$\lim_{n \rightarrow \infty} \mathbf{e}^{(n)} = \mathbf{0}.$$

Since  $\mathbf{x}^{(0)}$  and therefore  $\mathbf{e}^{(0)}$  is arbitrary it follows that the  
iteration will converge if and only if

$$\lim_{n \rightarrow \infty} \mathbf{G}^n = \mathbf{0}.$$

Assume now that the matrix  $\mathbf{G}$  of order  $m$  has  $m$   
LI eigenvectors  $\mathbf{v}_s$ ,  $s=1(1)m$ . Then these eigenvectors  
can be used as a basis for our  $m$ -dimensional vector  
space and the arbitrary error vector  $\mathbf{e}^{(0)}$ , with its  $m$

components, can be expressed uniquely as a linear combination of them, namely

$$e^{(0)} = \sum_{s=1}^m c_s v_s ,$$

where the  $c_s$  are scalars. Hence

$$e^{(1)} = Ge^{(0)} = \sum_{s=1}^m c_s G v_s ,$$

But  $Gv_s = \lambda_s v_s$ , where  $\lambda_s$  is the eigenvalue corresponding to  $v_s$ . Hence

$$e^{(1)} = \sum_{s=1}^m c_s \lambda_s v_s , \quad e^{(n)} = Ge^{(n-1)} = G \left[ \sum_{s=1}^m c_s \lambda_s v_s \right] \\ = \sum_{s=1}^m c_s \lambda_s G v_s = \sum_{s=1}^m c_s \lambda_s^2 v_s$$

Similarly,  $e^{(n)} = Ge^{(n-1)} = \sum_{s=1}^m c_s \lambda_s^n v_s . \quad (7)$

Therefore  $e^{(n)}$  will tend to the null vector as  $n$  tends to infinity, for arbitrary  $e^{(0)}$ , if and only if  $|\lambda_s| < 1$  for all  $s$ . In other words, the iteration will converge for arbitrary  $x^{(0)}$  if and only if the spectral radius  $\rho(G)$  of  $G$  is less than one.

As a corollary to this result a sufficient condition for convergence is that  $\|G\| < 1$ . To prove this we have that  $Gv_s = \lambda_s v_s$ . Hence

$$\|Gv_s\| = \|\lambda_s v_s\| = |\lambda_s| \|v_s\| .$$

On the other hand  $\|Gv_s\| \leq \|G\| \|v_s\|$

$$\Rightarrow |\lambda_s| \|v_s\| \leq \|G\| \|v_s\| \Rightarrow |\lambda_s| \leq \|G\|, \quad s=1(1)m.$$

-16-

It follows from this that a sufficient condition for convergence is that  $\|G\| < 1$ . It is not a necessary condition because the norm of  $G$  can exceed one even when  $f(G) < 1$ .

As an example consider the Jacobi iteration matrix  $D^{-1}(L+U)$ . If the  $i$ -th equation of  $Ax=b$

is

$$a_{i1}x_1 + a_{i2}x_2 + \dots + a_{ii}x_i + \dots + a_{in}x_n = b_i,$$

then the  $i$ -th row of  $D^{-1}(L+U)$  is

$$\frac{a_{i1}}{a_{ii}} \quad \frac{a_{i2}}{a_{ii}} \quad \dots \quad \frac{a_{ii}}{a_{ii}} \quad 0 \quad \frac{a_{i,i+1}}{a_{ii}} \quad \dots \quad \frac{a_{in}}{a_{ii}}.$$

Then if we take our matrix norm as the infinity norm, the Jacobi iteration will converge if

$$|a_{i1}| + |a_{i2}| + \dots + |a_{i,i-1}| + 0 + |a_{i,i+1}| + \dots + |a_{in}| < |a_{ii}|.$$

$\Rightarrow$  The Jacobi method applied to the equation  $Ax=b$  will converge if  $A$  is a strictly diagonally dominant matrix.