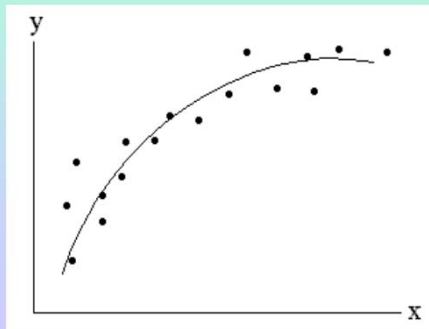


## 9. Approximation Theory (Spectral Methods)

### § Curve fitting

GIVEN:  $(x_i, y_i), i = 1, 2, \dots, N$

Find: a function which best approximates  
the unknown function  $y = f(x)$   
(not necessary passing the data points)



### Approximation Theory :

STEP1: choose a proper function from  $y = f(x; a, b, c, \dots)$   
with some adjustable parameters  $a, b, c, \dots$

STEP2: optimize the fitting in some way, i.e., minimize the  
error term

$$E \equiv \sum_{i=1}^N (y_i - f(x_i; a, b, c, \dots))^2 = E(a, b, c, \dots)$$

i.e. looking for a set of  $\{a, b, c, \dots\}$  such that

$$\frac{\partial E}{\partial a} = 0$$

$$\frac{\partial E}{\partial b} = 0$$

$$\frac{\partial E}{\partial c} = 0$$



Example: Given  $(x_i, y_i)$  and find a linear curve fitting.

Assume  $y = ax + b$ .

Define the derivation as

$$E \equiv \sum_{i=1}^N (y_i - ax_i - b)^2 = E(a, b)$$

To minimize the error  $E(a, b)$ , we look for  $(a, b)$  such that

$$\frac{\partial E}{\partial a} = 0 = \sum_{i=1}^N 2(y_i - ax_i - b)(-x_i)$$

$$\frac{\partial E}{\partial b} = 0 = \sum_{i=1}^N 2(y_i - ax_i - b)(-1)$$




$$0 = -\sum_{i=1}^N x_i y_i + a \sum_{i=1}^N x_i^2 + b \sum_{i=1}^N x_i$$

$$0 = -\sum_{i=1}^N y_i + a \sum_{i=1}^N x_i + b \sum_{i=1}^N 1$$

$$a = \frac{N \left( \sum_{i=1}^N x_i y_i \right) - \left( \sum_{i=1}^N x_i \right) \left( \sum_{i=1}^N y_i \right)}{N \left( \sum_{i=1}^N x_i^2 \right) - \left( \sum_{i=1}^N x_i \right)^2}$$

$$b = \frac{\left( \sum_{i=1}^N x_i^2 \right) \left( \sum_{i=1}^N y_i \right) - \left( \sum_{i=1}^N x_i y_i \right) \left( \sum_{i=1}^N x_i \right)}{N \left( \sum_{i=1}^N x_i^2 \right) - \left( \sum_{i=1}^N x_i \right)^2}$$



### § Approximation Theory

GIVEN: some function  $f(x)$  and some domain  $[a,b]$

STEP1: select a proper set of linearly independent functions  $\{\phi_i(x)\}$

STEP2: approximate the function from  $y = f(x)$  of the form

$$f(x) \approx S(x) = \sum_{i=0}^n c_i \phi_i(x)$$

which is expected to be convergent to the exact function  $f(x)$  as  $n \rightarrow \infty$

STEP3: adjust the values of parameters  $\{c_i\}$

so that the resulting approximation  $S(x)$  is the best.

### Definition (**linearly Independence**)

The set of functions  $\{\phi_0(x), \phi_1(x), \phi_2(x), \dots, \phi_n(x)\}$  is said to be linearly independent on  $[a, b]$  if whenever

$c_0\phi_0(x) + c_1\phi_1(x) + \dots + c_n\phi_n(x) = 0$  for  $c_i's \in R$  and all  $x \in [a, b]$ , then  $c_0 = c_1 = \dots = c_n = 0$ .



### Definition (Orthogonal)

$\{\phi_0, \phi_1, \phi_2, \dots, \phi_n\}$  is said to be an orthogonal set of functions on the interval  $[a, b]$  with respect to the weight function  $\omega(x)$  if

$$\int_a^b \omega(x) \phi_i(x) \phi_j(x) dx = \begin{cases} 0 & \text{if } i \neq j \\ \alpha_i \neq 0 & \text{if } i = j \end{cases} = \alpha_i \delta_{ij}$$

where  $\omega(x) \geq 0$  for all  $x \in [a, b]$

but  $\omega(x) \neq 0$  on any subinterval of  $[a, b]$ .

**Definition (Ortho-Normal)::** orthogonal with  $\alpha_i=1$



$$f(x) \approx S(x) = \sum_{i=0}^n c_i \phi_i(x)$$

- Best approximation:  
define the  $\omega$ -weighted error term as follows

$$\begin{aligned} E &\equiv \int_a^b \omega(x) (f(x) - S(x))^2 dx \\ &= \int_a^b \omega(x) \left( f(x) - \sum_{i=0}^n c_i \phi_i(x) \right)^2 dx \end{aligned}$$

The best approximation has a set of  $\{c_j\}$  which minimizes the error  $E$ .

$$\frac{\partial E}{\partial c_j} = 0 \text{ for all } j = 0, 1, 2, \dots, n$$

~  $(n+1)$  equations solve  $(n+1)$  unknowns  $\{c_j\}$  ~



Solution:

$$\begin{aligned} 0 &= \frac{\partial E}{\partial c_j} = \frac{\partial}{\partial c_j} \int_a^b \omega(x) \left( f(x) - \sum_{i=0}^n c_i \phi_i(x) \right)^2 dx \\ &= \int_a^b \omega(x) 2 \left( f(x) - \sum_{i=0}^n c_i \phi_i(x) \right) (-\phi_j(x)) dx \\ &= -2 \left\{ \int_a^b \omega(x) f(x) \phi_j(x) dx - \sum_{i=0}^n c_i \int_a^b \omega(x) \phi_i(x) \phi_j(x) dx \right\} \\ &\quad \sum_{i=0}^n c_i \int_a^b \omega(x) \phi_i(x) \phi_j(x) dx = \int_a^b \omega(x) f(x) \phi_j(x) dx \end{aligned}$$

~ (n+1) equations (j=0,1,2,...,n) for (n+1) unknowns ~

(solutions exists as long as  $\{\phi_i\}$  is L.I.)



With orthogonality (not necessary), we have

$$\int_a^b \omega(x) \phi_i(x) \phi_j(x) dx = \begin{cases} 0 & \text{if } i \neq j \\ \alpha_i \neq 0 & \text{if } i = j \end{cases} = \alpha_i \delta_{ij}$$

$$\sum_{i=0}^n c_i \alpha_i \delta_{ij} = \int_a^b \omega(x) f(x) \phi_j(x) dx$$

$$\alpha_j c_j = \int_a^b \omega(x) f(x) \phi_j(x) dx$$

$$c_j = \int_a^b \omega(x) f(x) \phi_j(x) dx / \alpha_j$$

$$= \int_a^b \omega(x) f(x) \phi_j(x) dx / \int_a^b \omega(x) \phi_j(x) \phi_j(x) dx$$



§ Commonly used functions  $\{\phi_i\}$

(i) **Fourier series** (orthogonal)

choose  $[a, b] = [-\pi, \pi]$  and  $w(x) = 1$  and

$$\phi_0(x) = 1$$

$$\phi_k(x) = \cos(kx) \text{ for } k = 1, 2, \dots, n$$

$$\phi_{k+n}(x) = \sin(kx) \text{ for } k = 1, 2, \dots, n$$

$$\int_{-\pi}^{\pi} \sin^2(kx) dx = \int_{-\pi}^{\pi} \cos^2(kx) dx = \pi$$

$$\phi_k(x) = \exp(\sqrt{-1} kx), \text{ for } k = 0, 1, 2, \dots, n$$



(ii) **Chebychev polynomials** (orthogonal)

choose  $[a, b] = [-1, 1]$  and  $w(x) = 1/\sqrt{1-x^2}$  and

$$\phi_k(x) = T_k(x) = \cos(k \cos^{-1} x) \text{ for } k = 0, 1, 2, \dots, n$$


$$\int_{-1}^1 T_k^2(x) / \sqrt{1-x^2} dx = \begin{cases} \pi/2 & \text{for } k \geq 1 \\ \pi & \text{for } k = 0 \end{cases}$$

(iii) **Legendre polynomials** (orthogonal)

choose  $[a, b] = [-1, 1]$  and  $w(x) = 1$  and

$$\phi_k(x) = L_k(x) \text{ for } k = 0, 1, 2, \dots, n$$

$$\int_{-1}^1 L_k^2(x) dx = \frac{2}{2k+1}$$




§ Application to ODEs ~ Rayleigh-Ritz method

Consider a linear 2nd - order ODE :  $y = y(x)$

$$-\frac{d}{dx} \left( p(x) \frac{dy}{dx} \right) + q(x)y = f(x)$$

for  $0 \leq x \leq 1$  and  $y(0) = y(1) = 0$

where  $p(x) \in C^1[0,1]$ ,  $q, f \in C^0[0,1]$   
 and  $\exists \delta > 0 \ni p(x) \geq \delta \geq 0$  and  $q(x) \geq 0$  for  $0 \leq x \leq 1$   
 (sufficient conditions for a unique solution)




**PRINCIPLE OF VARIATION**

The unique solution is the function which

- $\in C^2[0,1]$
- satisfies BCs
- minimizes the integral


$$I(u(x)) \equiv \int_0^1 \left\{ p(x) \left( \frac{du}{dx} \right)^2 + q(x)u(x)^2 - 2f(x)u(x) \right\} dx$$



§ Approximation Theory ~ ONE IDEA

- select a set of LI  $\{\phi_k(x)\}_{k=0}^n$   
each of which satisfies the given BCs:  
 $\phi_k(0) = \phi_k(1) = 0$

Define  $\Phi_n \equiv$  the set of functions spanned by  $\{\phi_k(x)\}_{k=0}^n$

$$= \left\{ y_n(x) \mid y_n(x) = \sum_{k=0}^n c_k \phi_k(x) \right\}$$


- select a function from  $\Phi_n$   
to approximate the exact solution  $y(x)$

$$y(x) \approx y_n(x) = \sum_{k=0}^n c_k \phi_k(x) \in \Phi_n$$


The approximation  $y_n(x) \in \Phi_n \subseteq C^2[0,1]$  and  $\ni$  B.C.s

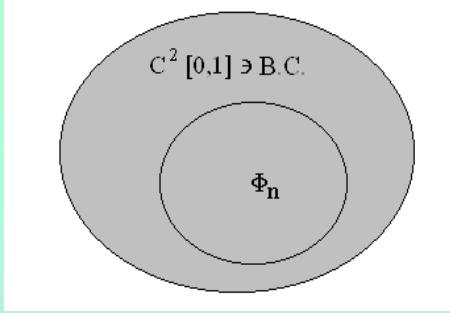
$$y_n(0) = \sum_{k=0}^n c_k \phi_k(0) = 0$$

$$y_n(1) = \sum_{k=0}^n c_k \phi_k(1) = 0$$

The exact solution  $y(x) \in C^2[0,1]$  and  $\ni$  B.C.s








$\Phi_n \rightarrow C^2[0,1]$  as  $n \rightarrow \infty$

- Which  $y_n$  is the best? The one having a minimum value of  $I$

Convergence:  $y_n \rightarrow y(x)$  as  $n \rightarrow \infty$



**SOLUTION::**

$$I(u(x)) \equiv \int_0^1 \left\{ p(x) \left( \frac{du}{dx} \right)^2 + q(x)u(x)^2 - 2f(x)u(x) \right\} dx$$

$$I(y_n(x) \in \Phi_n) = I \left( \sum_{k=0}^n c_k \phi_k(x) \right)$$

$$= \int_0^1 \left\{ p(x) \left( \sum_{k=0}^n c_k \phi_k'(x) \right)^2 + q(x) \left( \sum_{k=0}^n c_k \phi_k(x) \right)^2 - 2f(x) \sum_{k=0}^n c_k \phi_k(x) \right\} dx$$

**MINIMIAZATION:**  $\frac{\partial I(\Phi)}{\partial c_j} = 0$ , for  $j = 0, 1, 2, \dots, n$

$$\begin{aligned}
0 &= \frac{\partial I(\Phi)}{\partial c_j} \\
&= \frac{\partial}{\partial c_j} \int_0^1 \left\{ p(x) \left( \sum_{k=0}^n c_k \phi'_k(x) \right)^2 + q(x) \left( \sum_{k=0}^n c_k \phi_k(x) \right)^2 - 2f(x) \sum_{k=0}^n c_k \phi_k(x) \right\} dx \\
&= \int_0^1 \left\{ 2p(x) \phi'_j(x) \left( \sum_{k=0}^n c_k \phi'_k(x) \right) + 2q(x) \phi_j(x) \left( \sum_{k=0}^n c_k \phi_k(x) \right) - 2f(x) \phi_j(x) \right\} dx \\
&= 2 \sum_{k=0}^n \left\{ c_k \int_0^1 \left[ p(x) \phi'_j(x) \phi'_k(x) + q(x) \phi_j(x) \phi_k(x) \right] dx \right\} - 2 \int_0^1 f(x) \phi_j(x) dx
\end{aligned}$$

$a_{ij} = a_{jk}$ 
 $b_j$

Thus we obtain  $\sum_{k=0}^n a_{jk} c_k = b_j$ , for  $j = 0, 1, 2, \dots, n$

~ Solving ODE becomes solving matrix (algebraic) equations ~

§ Application to PDE ~ spectral methods, finite element methods, etc.

Given:  $\frac{\partial u}{\partial t} = L(u) + f(x, t)$

where  $L$  is some spatial operator, e.g. Laplacian  $L = \nabla^2$

$f(x, t)$  is a known external source term.

I.C.:  $u(x, 0) = g(x), a \leq x \leq b$

FIND:  $u(x, t)$  for  $t \geq 0$



§ Global approximation theory:

STEP1: find a set of LI  $\{\phi_j(x)\}$  and

the associated weight function  $w(x)$

$$u(x,t) \approx u_n(x,t) = \sum_{j=0}^n c_j(t) \phi_j(x)$$

P.S. The coefficients  $\{c_j\}$  is a function of  $t$  now.

for all  $x \in \Omega$

STEP2 : Substitute the above formula into PDE :  $\frac{\partial u}{\partial t} = L(u) + f$

$$\sum_{j=0}^n \frac{dc_j}{dt} \phi_j(x) \stackrel{?}{=} L\left(\sum_{j=0}^n c_j \phi_j(x)\right) + f(x,t)$$



§ Global approximation theory:

$$\sum_{j=0}^n \frac{dc_j}{dt} \phi_j(x) \stackrel{?}{=} L\left(\sum_{j=0}^n c_j \phi_j(x)\right) + f(x,t)$$

Let  $\Phi_n \equiv$  the space spanned by  $\{\phi_j^k(x)\}_{j=0}^n$ .

$$LHS \in \Phi_n$$

$$RHS \stackrel{?}{\in} \Phi_n$$

**not necessary!**

• Several ways to obtain an approximation:

(i) best approximation: look for a set of  $\{c_j(t)\}$  which minimizes

$$\left\| \sum_{j=0}^n \frac{dc_j}{dt} \phi_j(x) - L \left( \sum_{j=0}^n c_j \phi_j(x) \right) - f(x,t) \right\|$$

(ii) Collocation methods: choose a set of nodes  $\{x_i\}_{i=0}^n$  at which

$$\sum_{j=0}^n \frac{dc_j}{dt} \phi_j(x_i) = L \left( \sum_{j=0}^n c_j \phi_j(x_i) \right) + f(x_i, t)$$

(iii) Galerkin approximation

$$\begin{aligned} & \int_a^b \omega(x) \phi_i(x) \cdot \sum_{j=0}^n \frac{dc_j}{dt} \phi_j(x) dx \\ &= \int_a^b \omega(x) \phi_i(x) L \left( \sum_{j=0}^n c_j \phi_j(x) \right) dx + \int_a^b \omega(x) \phi_i(x) f(x,t) dx \end{aligned}$$

§ Galerkin approximation – linear  $L$

$$\begin{aligned} & \sum_{j=0}^n \frac{dc_j}{dt} \cdot \int_a^b \omega(x) \phi_i(x) \phi_j(x) dx \\ &= \sum_{j=0}^n c_j \int_a^b \omega(x) \phi_i(x) L \phi_j(x) dx + \int_a^b \omega(x) \phi_i(x) f(x,t) dx \end{aligned}$$

define  $a_{ij} = \int_a^b \omega(x) \phi_i(x) \phi_j(x) dx = \delta_{ij}$  if  $\omega$  – orthonormal

$$b_{ij} = \int_a^b \omega(x) \phi_i(x) L \phi_j(x) dx$$

$$r_i(t) = \int_a^b \omega(x) \phi_i(x) f(x,t) dx$$

$$\sum_{j=0}^n a_{ij} \frac{dc_j}{dt} = \sum_{j=0}^n b_{ij} c_j(t) + r_i(t)$$

~ PDE is reduced to a system of coupled ODEs ~

~ go time-marching ~



- special case:  $\omega$ -orthonormal  $\{\phi_j(x)\} : a_{ij} = \delta_{ij}$

$$\frac{dc_j}{dt} = \sum_{j=0}^n b_{ij}c_j(t) + r_i(t) \sim \text{decoupled on LHS} \sim$$

### § Initial Conditions, i.e. $c_j(0)$

$$u(x,0) = g(x) \approx u_n(x,0) = \sum_{j=0}^n c_j(0)\phi_j(x)$$

- best approximation : minimize  $\left\| g(x) - \sum_{j=0}^n c_j(0)\phi_j(x) \right\|$

- Collocation:  $g(x_i) = \sum_{j=0}^n c_j(0)\phi_j(x_i)$

- Galerkin:  $\int_a^b \omega(x)\phi_i(x)g(x)dx = \sum_{j=0}^n c_j(0) \int_a^b \omega(x)\phi_i(x)\phi_j(x)dx$



### § Boundary Conditions:

e.g.  $Bu(x,t) = 0$  for some linear spatial operator  $B$

- (a) select  $\{\phi_j(x)\}$  with the additional restriction

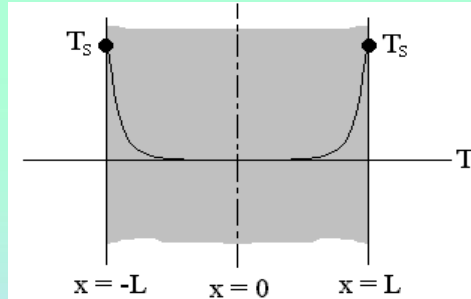
that each  $\phi_j(x) \in BCs$ . i.e.  $B\phi_j(x) = 0$

$$\Rightarrow Bu_n(x,t) = B \sum_{j=0}^n c_j(t)\phi_j(x) = \sum_{j=0}^n c_j(t)B\phi_j(x) = 0$$

- (b) (Tau method) Discard as many ODEs as the number of BCs and impose

$$Bu_n(x,t) = \sum_{j=0}^n c_j(t)B\phi_j(x) = 0$$

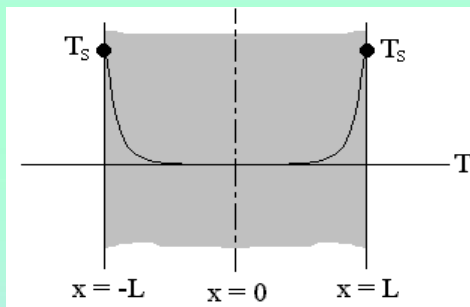
§ Example: thermal boundary layer



Governing equation :  $\frac{\partial T}{\partial t} = D \frac{\partial^2 T}{\partial x^2}$

I.C.:  $T(x, t = 0) = T_i$

B.C.s:  $T(x = \pm L, t \geq 0^+) = T_s$



define  $\Theta(x, t) = \frac{T - T_i}{T_s - T_i}$

Governing equation :  $\frac{\partial T}{\partial t} = D \frac{\partial^2 T}{\partial x^2}$


I.C.:  $T(x, t = 0) = T_i$

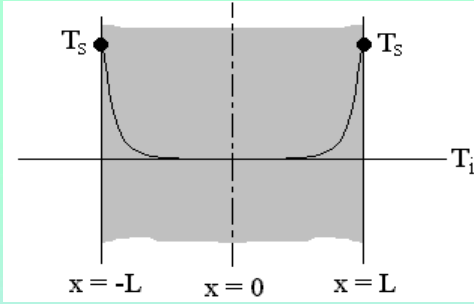
B.C.s:  $T(x = \pm L, t \geq 0^+) = T_s$

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

$\Theta(x, 0) = 0$

$\Theta(L, t) = 1$






- A thermal boundary layer is expected to develop near  $x = \pm 1$   
 $\Rightarrow$  better denser grid points near  $x = \pm 1$
- Expect the solution is symmetric.  
 $\Rightarrow$  choose Chebychev polynomials of even degrees

$$\{T_{2j}(x) = \cos(2j \cos^{-1} x)\}$$

### § collocation + $\tau$ methods

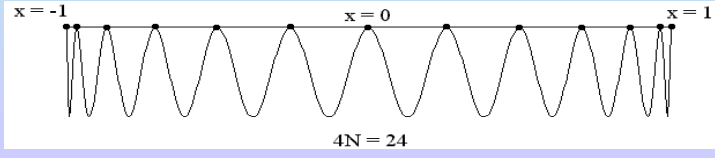


- choose  $\phi_j(x) = T_{2j}(x)$

$$\Theta(x, t) \approx \Theta_n(x, t) \equiv \sum_{j=0}^n c_j(t) T_{2j}(x)$$

- collocation grid points: choose the locations of peaks of Chebychev,

$$x_i \equiv \text{locations of peak of } T_{4n}(x) = \cos\left(\frac{2(n-i)\pi}{4n}\right) \quad i = 0, 1, 2, \dots, n$$



### § collocation + $\tau$ methods

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

$$\Theta(x, t) \approx \Theta_n(x, t) \equiv \sum_{j=0}^n c_j(t) T_{2j}(x)$$

$$\sum_{j=0}^n \frac{dc_j}{dt} T_{2j}(x) \stackrel{?}{=} D \sum_{j=0}^n c_j \frac{d^2 T_{2j}(x)}{dx^2}$$

$$\sum_{j=0}^n \frac{dc_j}{dt} T_{2j}(x_i) = D \sum_{j=0}^n c_j(t) \frac{d^2 T_{2j}(x)}{dx^2}(x_i)$$

for  $i = 0, 1, 2, \dots, n-1$  ( $i = n$  discarded)

$$\text{B.C.: } \Theta(1, t) = \sum_{j=0}^n c_j(t) T_{2j}(1) = \sum_{j=0}^n c_j(t) = 1$$

### • Initial Conditions

$$\Theta(x, t) \approx \Theta_n(x, t) \equiv \sum_{j=0}^n c_j(t) T_{2j}(x)$$

$$\Theta(x_i, 0) = 0 = \sum_{j=0}^n c_j(0) T_{2j}(x_i)$$

for  $i = 0, 1, 2, \dots, n-1$

$$\Theta(\pm 1, 0) = 1 = \sum_{j=0}^n c_j(0)$$

~ (n+1) equations for (n+1) unknowns ~



### § Galerkin method + $\{\phi_j(x)\}$ $\hat{=}$ B.C.s

- Since  $T_0(x) = 1$  and  $T_{2j}(\pm 1) = 1$

Let  $\phi_j(x) \equiv T_{2j}(x) - T_0(x)$ . Thus  $\phi_j(\pm 1) = 0$

$$\begin{aligned}\Theta(x, t) &\approx \Theta_n(x, t) \equiv 1 + \sum_{j=1}^n c_j(t) \phi_j(x) \\ &= 1 + \sum_{j=1}^n c_j(t) (T_{2j}(x) - T_0(x))\end{aligned}$$

$$\Theta_n(\pm 1, t) = 1 + \sum_{j=1}^n c_j(t) \phi_j(\pm 1) = 1 + \sum_{j=1}^n c_j(t) \cdot 0 = 1 \quad \hat{=} \text{B.C.s}$$

- Substitute into governing equation:  $\frac{\partial \Theta}{\partial t} = D \frac{\partial^2 \Theta}{\partial x^2}$


$$\Theta_n(x, t) = 1 + \sum_{j=1}^n c_j(t) \phi_j(x)$$

$$\sum_{j=1}^n \frac{dc_j}{dt} \phi_j(x) \stackrel{?}{=} D \sum_{j=1}^n c_j(t) \frac{d^2 \phi_j(x)}{dx^2}$$


- Use orthogonal property:

$$\sum_{j=1}^n \frac{dc_j}{dt} \cdot \int_{-1}^1 \frac{\phi_i(x) \phi_j(x)}{\sqrt{1-x^2}} dx = D \sum_{j=1}^n c_j(t) \int_{-1}^1 \frac{\phi_i(x) (d^2 \phi_j(x) / dx^2)}{\sqrt{1-x^2}} dx$$

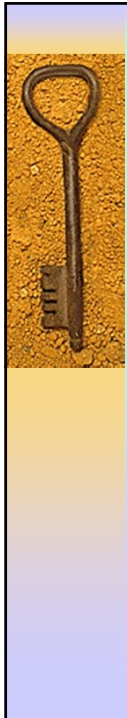
$$\sum_{j=1}^n a_{ij} \frac{dc_j}{dt} = D \sum_{j=1}^n b_{ij} c_j(t) \quad \text{for } i = 1, 2, \dots, n$$



$$\begin{aligned}
 a_{ij} &= \int_{-1}^1 \frac{\phi_i(x)\phi_j(x)}{\sqrt{1-x^2}} dx, \quad i, j = 1, 2, \dots, n \\
 &= \int_{-1}^1 (T_{2i}(x) - T_0(x))(T_{2j}(x) - T_0(x)) / \sqrt{1-x^2} dx \\
 &= \int_{-1}^1 (T_{2i}T_{2j} - T_0T_{2i} - T_0T_{2j} + T_0^2) / \sqrt{1-x^2} dx \\
 &= \int_{-1}^1 (T_{2i}T_{2j} + T_0^2) / \sqrt{1-x^2} dx \\
 &= \begin{cases} \pi & \text{if } i \neq j \\ \frac{3\pi}{2} & \text{if } i = j \geq 1 \end{cases}
 \end{aligned}$$



$$\begin{aligned}
 b_{ij} &\equiv \int_{-1}^1 \phi_i(x) \frac{d^2\phi_j(x)}{dx^2} / \sqrt{1-x^2} dx \\
 &= \int_{-1}^1 \left\{ (T_{2i}(x) - T_0(x)) \cdot \frac{d^2}{dx^2} (T_{2j}(x) - T_0(x)) \right\} / \sqrt{1-x^2} dx \\
 &= \int_{-1}^1 \left\{ (T_{2i} - T_0) \cdot \frac{d^2 T_{2j}}{dx^2} \right\} / \sqrt{1-x^2} dx \\
 \frac{d^2 T_{2j}(x)}{dx^2} &= P_{2j-2}(x) = \sum_{k=0}^{j-1} \gamma_k T_{2k}(x) \\
 b_{ij} &= \sum_{k=0}^{j-1} \gamma_k \cdot \int_{-1}^1 (T_{2i}T_{2k} - T_0T_{2k}) / \sqrt{1-x^2} dx
 \end{aligned}$$

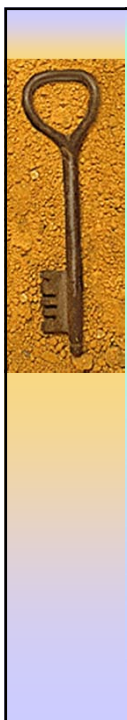


- Initial Conditions:  $\Theta(x, 0) = 0$

$$\Theta_n(x, t) = 1 + \sum_{j=1}^n c_j(t) \phi_j(x)$$

$$\Theta_n(x, 0) = 1 + \sum_{j=1}^n c_j(0) \phi_j(x)$$

$$\int_{-1}^1 \frac{(\Theta(x, 0) - 1) \phi_i(x)}{\sqrt{1-x^2}} dx = \sum_{j=1}^n c_j(0) \int_{-1}^1 \frac{\phi_i(x) \phi_j(x)}{\sqrt{1-x^2}} dx$$

$$\pi = \int_{-1}^1 \frac{(1 - T_{2i}(x))}{\sqrt{1-x^2}} dx = \sum_{j=1}^n a_{ij} c_j(0) \quad \text{for } i = 1, 2, \dots, n$$


- ♣ **Local approximation theory**
- divide the interested domain  $\Omega \equiv [a, b]$  into several subdomins  $\Omega = \bigcup_{k=1}^K \Omega_k$

$$u(x, t) \approx u_n^k(x, t) \equiv \sum_{j=0}^n c_j^k(t) \phi_j^k(x) \quad \text{for all } x \in \Omega_k$$

P.S.  $\{\phi_j^k(x)\}_{j=0}^n$  is different from one subdomain to another.

i.e. require  $K$  sets of L.I. functions  $\{\phi_j^k(x)\}_{j=0}^n$  and will have  $K$  sets of corresponding coefficients  $\{c_j^k(t)\}_{j=0}^n$

§ **Finite Element Methods** ~ piecewise approximation

Test problem : (domain  $\Omega$  with boundary  $\partial\Omega$ )

$$\frac{\partial}{\partial x} \left( p(x, y) \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( q(x, y) \frac{\partial u}{\partial y} \right) + r(x, y)u(x, y) = f(x, y)$$

$$u(x, y) = g(x, y) \quad \text{on } \partial\Omega$$

PRINCIPLE OF VARIATION ::

The exact solution (if exists and unique) is the function  $\in C^2(\Omega)$  which satisfies BCs and minimizes the value of

$$I(u) = \int_{\Omega} \left\{ \frac{1}{2} \left[ p(x, y) \left( \frac{\partial u}{\partial x} \right)^2 + q(x, y) \left( \frac{\partial u}{\partial y} \right)^2 - r(x, y)u^2 \right] + f(x, y)u(x, y) \right\} dx dy$$


STEP1: Divide the interested domain into several parts, i.e.  $\Omega = \bigcup_{k=1}^K \Omega_k$

Each subdomain  $\Omega_k$  is called an "element".

STEP2: within each subdomain  $\Omega_k$ , choose a proper set of  $\left\{ \phi_j^k(x, y) \right\}_{j=0}^n$

$$u(x, y) \approx u_n^k(x, y) \equiv \sum_{j=0}^n c_j^k \phi_j^k(x, y) \quad \text{for } (x, y) \in \Omega_k$$

P.S.  $\left\{ \phi_j^k(x, y) \right\}_{j=0}^n$  are different from one element to another in general.



$$\text{Put } u(x, y) \approx \bigcup_{k=1}^K u_n^k(x, y) = \bigcup_{k=1}^K \sum_{j=0}^n c_j^k \phi_j^k(x, y)$$

STEP3: Look for the values of  $\{c_j^k\}$ ,  $j = 0, 1, \dots, n$  and  $k = 0, 1, \dots, K$  such that  $\{u_n^k(x, y)\}$  minimizes  $I(u)$

$$I(u) = \int_{\Omega} \left\{ \frac{1}{2} \left[ p(x, y) \left( \frac{\partial u}{\partial x} \right)^2 + q(x, y) \left( \frac{\partial u}{\partial y} \right)^2 - r(x, y) u^2 \right] + f(x, y) u(x, y) \right\} dx dy$$

$$\text{Thus } I(u) = \int_{\Omega} \dots = \sum_{k=1}^K \int_{\Omega_k} \dots = I(c_j^k)$$

Therefore for all  $i \in \{0, 1, \dots, n\}$  and all  $m \in \{0, 1, \dots, K\}$


$$0 = \frac{\partial I(u)}{\partial c_i^m} \quad \text{required (minimized)}$$

$$= \frac{\partial}{\partial c_i^m} \sum_{k=1}^K \int_{\Omega_k} \frac{1}{2} p(x, y) \left( \sum_{j=0}^n c_j^k \frac{\partial \phi_j^k}{\partial x} \right)^2 + \dots$$

$$= \frac{\partial}{\partial c_i^m} \int_{\Omega_m} \frac{1}{2} p(x, y) \left( \sum_{j=0}^n c_j^m \frac{\partial \phi_j^m}{\partial x} \right)^2 + \dots$$

$$= \int_{\Omega_m} \frac{\partial}{\partial c_i^m} \left\{ \frac{1}{2} p(x, y) \left( \sum_{j=0}^n c_j^m \frac{\partial \phi_j^m}{\partial x} \right)^2 + \dots \right\} dx dy \quad \sum_{j=0}^n \frac{\partial \phi_j^m}{\partial x} \delta_{ij} = \frac{\partial \phi_i^m}{\partial x}$$

$$= \int_{\Omega_m} \left\{ p(x, y) \left( \sum_{j=0}^n c_j^m \frac{\partial \phi_j^m}{\partial x} \right) \frac{\partial}{\partial c_i^m} \left( \sum_{j=0}^n c_j^m \frac{\partial \phi_j^m}{\partial x} \right) + \dots \right\} dx dy$$




$$0 = \frac{\partial I(u)}{\partial c_i^m} \quad \text{required (minimized)}$$

$$= \int_{\Omega_m} \left\{ p(x, y) \sum_{j=0}^n c_j^m \frac{\partial \phi_j^m}{\partial x} \frac{\partial \phi_i^m}{\partial x} + q(x, y) \sum_{j=0}^n c_j^m \frac{\partial \phi_j^m}{\partial y} \frac{\partial \phi_i^m}{\partial y} \right. \\ \left. - r(x, y) \sum_{j=0}^n c_j^m \phi_j^m \phi_i^m + f(x, y) \phi_i^m \right\} dx dy$$

After rearrangement : (change variable  $m \rightarrow k$ )

$$0 = \sum_{j=0}^n c_j^k \cdot \int_{\Omega_k} \left\{ p(x, y) \frac{\partial \phi_j^k}{\partial x} \frac{\partial \phi_i^k}{\partial x} + q(x, y) \frac{\partial \phi_j^k}{\partial y} \frac{\partial \phi_i^k}{\partial y} - r(x, y) \phi_j^k \phi_i^k \right\} dx dy$$

$$- \int_{\Omega_k} f(x, y) \phi_i^k(x, y) dx dy$$


$$a_{ij}^k = \int_{\Omega_k} \left\{ p(x, y) \frac{\partial \phi_j^k}{\partial x} \frac{\partial \phi_i^k}{\partial x} + q(x, y) \frac{\partial \phi_j^k}{\partial y} \frac{\partial \phi_i^k}{\partial y} - r(x, y) \phi_i^k \phi_j^k \right\} dx dy$$


$$= a_{ji}^k$$

$$b_i^k = - \int_{\Omega_k} f(x, y) \phi_i^k(x, y) dx dy$$

$$\sum_{j=0}^n a_{ij}^k c_j^k = b_i^k \quad \text{or} \quad A^k c^k = b^k \sim \text{elemental matrix equation}$$

$$\text{for } k = 1, 2, \dots, K$$

- $A^k$  is singular! i.e. Each element matrix equation is underminate.
- ▶ too many degrees of freedom!



• Continuity requirement : if  $(x, y) \in \Omega_{k_1}$ , and also  $\in \Omega_{k_2}$ , then


$$u(x, y) \approx u_n^{k_1} = \sum_{j=0}^n c_j^{k_1} \phi_j^{k_1}(x, y)$$

$$u(x, y) \approx u_n^{k_2} = \sum_{j=0}^n c_j^{k_2} \phi_j^{k_2}(x, y)$$

Thus  $C^0$ -requirement needs

$$\sum_{j=0}^n c_j^{k_1} \phi_j^{k_1}(x, y) = \sum_{j=0}^n c_j^{k_2} \phi_j^{k_2}(x, y)$$

$\Rightarrow$  The minimization problem should be redone with the  $C^0$ -requirements



§ Example: use collocation points


STEP1: within element  $\Omega_k$ , choose  $(n+1)$  points denoted by  $(x_i^k, y_i^k)$  for  $i = 0, 1, 2, \dots, n$

STEP2 : within each element, choose an interpoint function, say two - dimensional Lagrangian polynomials,  $\phi_i^k(x, y)$

$$\phi_i^k(x, y) \equiv \prod_{\substack{l=0 \\ l \neq i}}^n \frac{\sqrt{(x - x_l^k)^2 + (y - y_l^k)^2}}{\sqrt{(x_i^k - x_l^k)^2 + (y_i^k - y_l^k)^2}} \quad \ni \phi_i^k(x_j^k, y_j^k) = \delta_{ij}$$

and approximate the desired function as

$$u(x, y) \approx u_n^k \equiv \sum_{i=0}^n c_i^k \phi_i^k(x, y), \quad \text{whenever } (x, y) \in \Omega_k$$




$$u(x, y) \approx u_n^k \equiv \sum_{i=0}^n c_i^k \phi_i^k(x, y), \quad \text{whenever } (x, y) \in \Omega_k$$

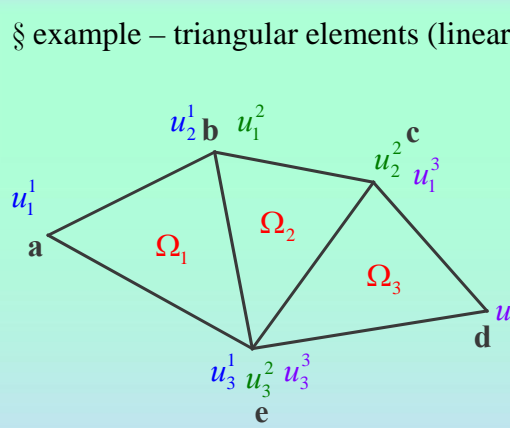
Notice:

$$u_n^k(x_i, y_i) \equiv u_i^k = \sum_{j=0}^n c_j^k \phi_j^k(x_i, y_i)$$

$$u_i^k = \sum_{j=0}^n c_j^k \delta_{ij} = c_i^k$$

$$u_n^k(x, y) = \sum_{j=0}^n u_j^k \phi_j^k(x, y)$$


§ example – triangular elements (linear interpolation)




$$C^0 : \begin{cases} u_1^1 = u_a \\ u_2^1 = u_1^2 = u_b \\ u_2^2 = u_1^3 = u_c \\ u_3^2 = u_d \\ u_3^1 = u_3^2 = u_3^3 = u_e \end{cases}$$

If no  $C^0$  constraint: d.o.f.s =  $u_1^1, u_2^1, u_3^1, u_1^2, u_2^2, u_3^2, u_1^3, u_2^3, u_3^3$

with  $C^0$  - constraint: d.o.f.s =  $u_a, u_b, u_c, u_d, u_e$





\*d.o.f.  $u_a (= u_1^1)$ :

$$0 = \frac{\partial I(u)}{\partial u_a}$$

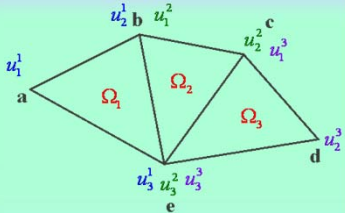

$$= \frac{\partial}{\partial u_a} \sum_{k=1}^K \int_{\Omega_k} \dots \sum_{j=0}^n u_j^k \phi_j^k(x, y) \dots$$

$$= \frac{\partial}{\partial u_a} \int_{\Omega_1} \dots \sum_{j=0}^n u_j^1 \phi_j^1(x, y) \dots$$

$$= \frac{\partial}{\partial u_1^1} \int_{\Omega_1} \dots \sum_{j=0}^n u_j^1 \phi_j^1(x, y) \dots$$

$$0 = \{a_{11}^1 u_1^1 + a_{12}^1 u_2^1 + a_{13}^1 u_3^1 - b_1^1\}$$

\*d.o.f.  $u_d (= u_2^3)$ :  $0 = \{a_{21}^3 u_1^3 + a_{22}^3 u_2^3 + a_{23}^3 u_3^3 - b_2^3\}$

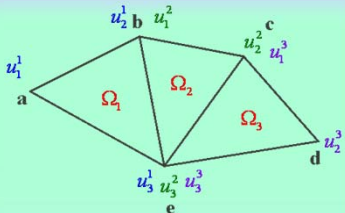




\*d.o.f.  $u_b (= u_2^1 = u_1^2)$ :

$$0 = \frac{\partial I(u)}{\partial u_b}$$

$$= \frac{\partial}{\partial u_b} \left( \int_{\Omega_1} \dots \sum_{j=0}^n u_j^1 \phi_j^1(x, y) \dots + \int_{\Omega_2} \dots \sum_{j=0}^n u_j^2 \phi_j^2(x, y) \dots \right)$$

$$= \frac{\partial}{\partial u_2^1} \int_{\Omega_1} \dots \sum_{j=0}^n u_j^1 \phi_j^1(x, y) \dots + \frac{\partial}{\partial u_1^2} \int_{\Omega_2} \dots \sum_{j=0}^n u_j^2 \phi_j^2(x, y) \dots$$

$$0 = \{a_{21}^1 u_1^1 + a_{22}^1 u_2^1 + a_{23}^1 u_3^1 - b_2^1\} + \{a_{11}^2 u_1^2 + a_{12}^2 u_2^2 + a_{13}^2 u_3^2 - b_1^2\}$$


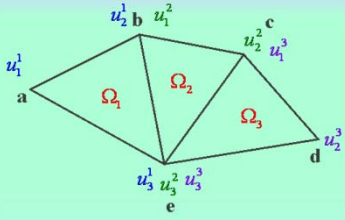



\*d.o.f.  $u_c (= u_2^2 = u_1^3)$ :

$$0 = \frac{\partial I(u)}{\partial u_c}$$

$$= \frac{\partial}{\partial u_c} \left( \int_{\Omega_2} \cdots \sum_{j=0}^n u_j^2 \phi_j^2(x, y) \cdots + \int_{\Omega_3} \cdots \sum_{j=0}^n u_j^3 \phi_j^3(x, y) \cdots \right)$$

$$= \frac{\partial}{\partial u_2^2} \int_{\Omega_2} \cdots \sum_{j=0}^n u_j^2 \phi_j^2(x, y) \cdots + \frac{\partial}{\partial u_1^3} \int_{\Omega_3} \cdots \sum_{j=0}^n u_j^3 \phi_j^3(x, y) \cdots$$

$$0 = \{a_{21}^2 u_1^2 + a_{22}^2 u_2^2 + a_{23}^2 u_3^2 - b_2^2\} + \{a_{11}^3 u_1^3 + a_{12}^3 u_2^3 + a_{13}^3 u_3^3 - b_1^3\}$$



\*d.o.f.  $u_e (= u_3^1 = u_3^2 = u_3^3)$ :

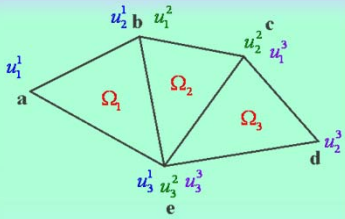
$$0 = \frac{\partial I(u)}{\partial u_e}$$

$$= \frac{\partial}{\partial u_c} \left( \int_{\Omega_1 + \Omega_2 + \Omega_3} \cdots \right)$$

$$= \frac{\partial}{\partial u_3^1} \int_{\Omega_1} \cdots \sum_{j=0}^n u_j^1 \phi_j^1(x, y) \cdots + \frac{\partial}{\partial u_3^2} \int_{\Omega_2} \cdots \sum_{j=0}^n u_j^2 \phi_j^2(x, y) \cdots$$

$$+ \frac{\partial}{\partial u_3^3} \int_{\Omega_3} \cdots \sum_{j=0}^n u_j^3 \phi_j^3(x, y) \cdots$$

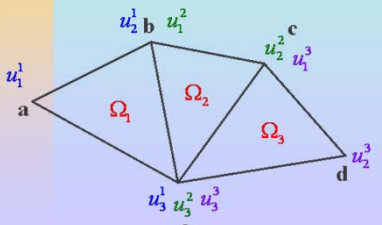
$$0 = \{a_{31}^1 u_1^1 + a_{32}^1 u_2^1 + a_{33}^1 u_3^1 - b_3^1\} + \{a_{31}^2 u_1^2 + a_{32}^2 u_2^2 + a_{33}^2 u_3^2 - b_2^2\}$$

$$+ \{a_{31}^3 u_1^3 + a_{32}^3 u_2^3 + a_{33}^3 u_3^3 - b_3^3\}$$


• Global matrix equation  $Au=b$

$a$	$b$	$c$	$d$	$e$	
$a_{11}^1$	$a_{12}^1$	0	0	$a_{13}^1$	$\begin{pmatrix} u_a = u_1^1 \\ u_b = u_2^1 = u_1^2 \\ u_c = u_2^2 = u_1^3 \\ u_d = u_2^3 \\ u_e = u_3^1 = u_2^2 = u_3^3 \end{pmatrix}$
$a_{21}^1$	$a_{22}^1 + a_{11}^2$	$a_{12}^2$	0	$a_{23}^1 + a_{13}^2$	
0	$a_{21}^2$	$a_{22}^2 + a_{11}^3$	$a_{12}^3$	$a_{23}^2 + a_{13}^3$	
0	0	$a_{21}^3$	$a_{22}^3$	$a_{23}^3$	
$a_{31}^1$	$a_{32}^1 + a_{31}^2$	$a_{32}^2 + a_{31}^3$	$a_{32}^3$	$a_{33}^1 + a_{33}^2 + a_{33}^3$	



$$= \begin{pmatrix} b_1^1 \\ b_2^1 + b_1^2 \\ b_2^2 + b_1^3 \\ b_2^3 \\ b_3^1 + b_3^2 + b_3^3 \end{pmatrix}$$

$$\begin{pmatrix} a_{11}^1 & a_{12}^1 & 0 & 0 & a_{13}^1 \\ a_{21}^1 & a_{22}^1 & 0 & 0 & a_{23}^1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ a_{31}^1 & a_{32}^1 & 0 & 0 & a_{33}^1 \end{pmatrix} \begin{pmatrix} u_1^1 \\ u_2^1 \\ 0 \\ 0 \\ u_3^1 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & a_{11}^2 & a_{12}^2 & 0 & a_{13}^2 \\ 0 & a_{21}^2 & a_{22}^2 & 0 & a_{23}^2 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & a_{31}^2 & a_{32}^2 & 0 & a_{33}^2 \end{pmatrix} \begin{pmatrix} 0 \\ u_1^2 \\ u_2^2 \\ 0 \\ u_3^2 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & a_{11}^3 & a_{12}^3 & a_{13}^3 \\ 0 & 0 & a_{21}^3 & a_{22}^3 & a_{23}^3 \\ 0 & 0 & a_{31}^3 & a_{32}^3 & a_{33}^3 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ u_1^3 \\ u_2^3 \\ u_3^3 \end{pmatrix} = \begin{pmatrix} b_1^1 \\ b_2^1 \\ 0 \\ 0 \\ b_3^1 \end{pmatrix} + \begin{pmatrix} 0 \\ b_1^2 \\ b_2^2 \\ 0 \\ b_3^2 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ b_1^3 \\ b_2^3 \\ b_3^3 \end{pmatrix}$$

stiffness summation

• local matrix equation

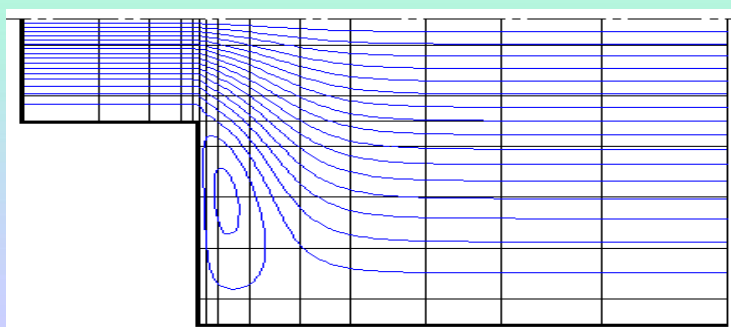
$$k = 1: \begin{pmatrix} a_{11}^1 & a_{12}^1 & a_{13}^1 \\ a_{21}^1 & a_{22}^1 & a_{23}^1 \\ a_{31}^1 & a_{32}^1 & a_{33}^1 \end{pmatrix} \begin{pmatrix} u_1^1 \\ u_2^1 \\ u_3^1 \end{pmatrix} = \begin{pmatrix} a_{11}^1 & a_{12}^1 & a_{13}^1 \\ a_{21}^1 & a_{22}^1 & a_{23}^1 \\ a_{31}^1 & a_{32}^1 & a_{33}^1 \end{pmatrix} \begin{pmatrix} u_a \\ u_b \\ u_e \end{pmatrix} = \begin{pmatrix} b_1^1 \\ b_2^1 \\ b_3^1 \end{pmatrix}$$


$$k = 2: \begin{pmatrix} a_{11}^2 & a_{12}^2 & a_{13}^2 \\ a_{21}^2 & a_{22}^2 & a_{23}^2 \\ a_{31}^2 & a_{32}^2 & a_{33}^2 \end{pmatrix} \begin{pmatrix} u_1^2 \\ u_2^2 \\ u_3^2 \end{pmatrix} = \begin{pmatrix} a_{11}^2 & a_{12}^2 & a_{13}^2 \\ a_{21}^2 & a_{22}^2 & a_{23}^2 \\ a_{31}^2 & a_{32}^2 & a_{33}^2 \end{pmatrix} \begin{pmatrix} u_b \\ u_c \\ u_e \end{pmatrix} = \begin{pmatrix} b_1^2 \\ b_2^2 \\ b_3^2 \end{pmatrix}$$

$$k = 3: \begin{pmatrix} a_{11}^3 & a_{12}^3 & a_{13}^3 \\ a_{21}^3 & a_{22}^3 & a_{23}^3 \\ a_{31}^3 & a_{32}^3 & a_{33}^3 \end{pmatrix} \begin{pmatrix} u_1^3 \\ u_2^3 \\ u_3^3 \end{pmatrix} = \begin{pmatrix} a_{11}^3 & a_{12}^3 & a_{13}^3 \\ a_{21}^3 & a_{22}^3 & a_{23}^3 \\ a_{31}^3 & a_{32}^3 & a_{33}^3 \end{pmatrix} \begin{pmatrix} u_c \\ u_d \\ u_e \end{pmatrix} = \begin{pmatrix} b_1^3 \\ b_2^3 \\ b_3^3 \end{pmatrix}$$

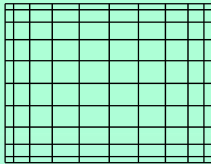
§ Example – 2D Poisson Solver

$$\nabla^2 P = \frac{\partial^2 P}{\partial x^2} + \frac{\partial^2 P}{\partial y^2} = f(x, y, t)$$





• Let



$$h_i^k(x) \equiv \prod_{\substack{m=0 \\ m \neq i}}^N \frac{(x - x_m^k)}{(x_i^k - x_m^k)}$$

$$h_j^k(y) \equiv \prod_{\substack{m=0 \\ m \neq j}}^N \frac{(y - y_m^k)}{(y_j^k - y_m^k)}$$

where  $\{x_i^k\}_{i=0}^N$  and  $\{y_j^k\}_{j=0}^N$  are Gauss Quadrature Intergration nodes

$$P(x, y, t) = \sum_{k=1}^K \sum_{i=0}^N \sum_{j=0}^N P_{ij}^k h_i^k(x) h_j^k(y)$$

$$f(x, y, t) = \sum_{k=1}^K \sum_{i=0}^N \sum_{j=0}^N f_{ij}^k h_i^k(x) h_j^k(y)$$


## Numerical Analysis

- 挑戰人類結合數學及物理知識能力
- 激發人類的創造力與邏輯思考
- 數位化儀器的秘密武器
- 提供經濟的設計分析工具

- 數值誤差 – truncation error & rounding error
- 數值現象 – numerical diffusion & dispersion
- 穩定性、準確性、效率

