Computational Fluid Dynamics — Prof. Suman Chakraborty

Lecture 1: Introduction to CFD

CFD: Computational Fluid Dynamics

Not just limited to Fluid Dynamics, but to any general transport phenomena. (like heat transfer, mass transfer)

CFD applications:

- Aerospace — Interior: ventilation system, combustor engine
                Exterior: flow over A/c.
- automobile
- biomedical
- chemical — (mixing of chemicals ; bubbles at the interface, separation & mixing, injection of streams).
- Electronics — efficient cooling strategies.
- Energy -
                                                        & dynamic coupling.
- Fluid structure interaction (two way coupling — fluid interacts with structure & structure interacts with fluid)
- Marine (eg: flow past ships & boats).
                                                        CFD.
- Materials processing (eg: grain growth ; need not be a deterministic approach; can be stochastic like Monte Carlo as well, but welding ?, mold filling ?
                                                        ⤷ helps choose the correct filling process to avoid casting defects ?.

find out if the distribution of impurities via their convection
                                is okay
across the material & viable ? for the material itself)

- Micro fluidics — microscale fluid flows. (micron or sub-micron scales)
        (eg: mix two streams — good mixing via pulsating flows, droplet dynamics — micro reactor studies,
                Fluid structure interaction at small scales, flap placed
                        reacting
                on micro flows — need to optimise the flaps movement,
flap acting like a mixer & a pump)

- sports: (racing cars, golf balls, running motion etc.)
- Turbomachines: (flow over blade passages

○ Is CFD inevitable?

Numerical vs Analytical vs Experimental.

- Experimental investigation:
    - full scale                    (no substitute for this: seeing is believing)
        · expensive & often impossible
        · measurement errors.
    - on a scaled model
        · simplified             (should upscale/downscale in such a
        · difficult to extrapolate results       manner that the flow
                                                physics doesn't changes
        · measurement errors.                    eg: micro capillary
                                                    scaled up:
  · Theoretical calculation:                       at small scales,
                                                    surface tension
    - analytical solutions:                        forces are
        · if exists, gives us exact answers}        significant)
          but exists only for a few cases.  (need to maintain
                                            all similarities - kinematic,
                                                    dynamic etc···
  ⎡ CFD cannot stand on its own without experimental  or  analytical}
  ⎣ solutions. B/c we need to benchmark our solution.

        · Sometimes complex

    - numerical solutions:
                exists
        · for almost any problems.
        · Continuous nature of the problem is compromised, but
          at its expense we get answers to complex problems

# Modeling vs Experimentation.

## • Advantages of modeling:

- cheaper
- more complete information (all details of all variables can be obtained)
- can handle any degree of complexity as long as·

## • Disadvantages of modeling:

- deals with a mathematical description not a reality.
  Numerical solution is as good as the input fed to the problem
- Mathematical description can be inadequate
  (Governing $eq^{ns}$ may / may not capture the correct physics).
- multiple solutions can exist·
  (non-linear problems may have multiple different solutions depending on the $IC$).

In conclusion: no real substitute for experimentation, but experimentation is limited by many restrictions & cannot handle multiple trials).

Usual plan of action: try analytical solution → do numerical simulation → create good experimental design & validate results.

→ Cross validate the goodness of numerical solution. from the numerical simulation.

$$\boxed{\text{Lecture 2: Classification of PDEs}}$$

$$\frac{\partial}{\partial t}(\rho \phi) + \nabla \cdot (\rho \vec{V} \phi) = \nabla \cdot (\Gamma \nabla \phi) + S$$

— 2$^{nd}$ order PDE  — from $\nabla \cdot (\Gamma \nabla \phi)$.

$\mathsf{C}$onsider a 2$^{nd}$ order PDE of the form:

$$A\phi_{xx} + B\phi_{xy} + C\phi_{yy} + \underbrace{D\phi_x + E\phi_y + F\phi + G}_{H} = 0 \cdot \quad —①$$

$\phi_x \rightarrow \dfrac{\partial \phi}{\partial x}$

$\phi_y \rightarrow \dfrac{\partial \phi}{\partial y}$

$\phi_{xx} \rightarrow \dfrac{\partial^2 \phi}{\partial x^2}$

$\phi_{xy} \rightarrow \dfrac{\partial^2 \phi}{\partial x \partial y}$

$\phi_{yy} \rightarrow \dfrac{\partial^2 \phi}{\partial y^2}$

- $A, B, C$ need not be consts. They can be fns / lower order partial derivatives / independent variables.

- $\phi = \phi(x, y)$.

- One classification —
  1. Linear : $A, B, C \rightarrow$ fns of $x, y$. ~~D, E, F, G fns~~   remaining terms $DE, F, G$ linear fns of $\phi, \phi_x, \phi_y$.
  2. ~~Non Linear~~ Quasi linear:   $A, B, C \rightarrow$ fns of $x, y, \phi, \phi_x, \phi_y$.

- $\mathsf{C}$haracteristics of the PDE:

  Highest order derivatives in a PDE may be continuous or discontinuous in the domain of consideration. There may be lines $a$ to $b$ $\nearrow\!\!\!\!\!\nwarrow$ along which these highest order derivatives could be have are discontinuous. Such lines are called characteristics lines of the PDE.

  Why They're important? B/c in our numerical method, we have to account for such discontinuities apriori.

  Objective: to get characteristics of the PDE
  ~~Get~~ $\phi_x = \phi_x(x, y) \Rightarrow d\phi_x = \dfrac{\partial \phi_x}{\partial x} dx + \dfrac{\partial \phi_x}{\partial y} dy.$

  $$d\phi_x = \phi_{xx}\, dx + \phi_{xy}\, dy \quad —②$$

  $\amalg^{rly}$ $\phi_y = \phi_y(x, y) \Rightarrow d\phi_y = \phi_{yx}\, dx + \phi_{yy}\, dy.$

  $\therefore \phi_{xy} = \phi_{yx}, \qquad d\phi_y = \phi_{xy}\, dx + \phi_{yy}\, dy. \quad —③$

From ①, ②, ③:

$$\begin{bmatrix} A & B & C \\ dx & dy & 0 \\ 0 & dx & dy \end{bmatrix} \begin{bmatrix} \phi_{xx} \\ \phi_{xy} \\ \phi_{yy} \end{bmatrix} = \begin{bmatrix} -H \\ d\phi_x \\ d\phi_y \end{bmatrix}$$

We are interested in the case where the solution of this

matrix expression → $\begin{bmatrix} \phi_{xx} \\ \phi_{xy} \\ \phi_{yy} \end{bmatrix}$ doesn't exist.

Analogy with algebraic eq$^{ns}$:

$$\begin{matrix} x+y = 2 \\ 2x+2y = 5 \end{matrix} \Big\} \text{ sol}^n \text{ doesn't exist.}$$

$$\begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

$$\hookrightarrow \Delta = 0.$$

We want to find locus of pts across which we have discontinuites

in $\phi_{xx}$, $\phi_{yy}$, & $\phi_{xy}$. That is possible when $\det\left(\begin{bmatrix} A & B & C \\ dx & dy & 0 \\ 0 & dx & dy \end{bmatrix}\right) = 0$

For $\phi_{xx}$, $\phi_{yy}$, $\phi_{xy}$ to be discontinuous,

$\Delta = 0.$

$$\begin{vmatrix} A & B & C \\ dx & dy & 0 \\ 0 & dx & dy \end{vmatrix} = 0 \Rightarrow$$

{only coeffs of highest order derivatives matters}

$$A (dy)^2 - B (dx\,dy) + C((dx)^2) = 0$$

$$A \left(\frac{dy}{dx}\right)^2 - B\left(\frac{dy}{dx}\right) + C = 0$$

$$\frac{dy}{dx} = \frac{B \pm \sqrt{B^2 - 4AC}}{2A}.$$

Number of real characteristics existing will depend on whether

$B^2 - 4AC \geq 0$ or $< 0$.

If $B^2 - 4AC = 0 \rightarrow$ only one real characteristic.
(parabolic).

$B^2 - 4AC < 0 \rightarrow$ no real characteristic (elliptic PDE)

$B^2 - 4AC > 0 \rightarrow$ 2 real characteristics (hyperbolic PDE)

# Lecture 3 : Examples of PDEs

**Ex 1:** $\nabla^2 \phi = 0$   (Laplace eq$^n$ — most commonly encountered simple PDE)

$\nabla^2 T = 0$   (T: Temperature, uniform heat conductivity, steady state, no heat source).
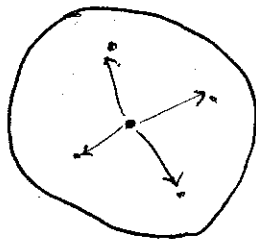
2D example:

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = 0$$

$A = 1, \ B = 0, \ C = 1.$

$B^2 - 4AC = -4 \longrightarrow$ elliptic equation.

Say, we have a domain with uniform temperature throughout. Put a heat source at a pt $\rightarrow$ acts as a disturbance that propagates in all directions in the domain at infinite speed. i.e. Thermal disturbance propagates in all directions with infinite speed.



This disturbance tries to nullify the temperature differences at different points in the domain. to make the temp distribution homogeneous everywhere.

— So while numerically formulating the problems, a point on the domain under our consideration will be influenced by all the other points in our domain.

— B/c has to be specified at the boundary of the domain. The B/c can be discontinuous. For eg: half the boundary may be in contact with steam & the other half in contact with ice, so discontinuities are possible on the boundary. But since the disturbance travel at infinite speeds (message propagation is fast) means that there will not be discontinuities within the domain.

— This type of problem $\rightarrow$ BVP.

Ex - 2

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

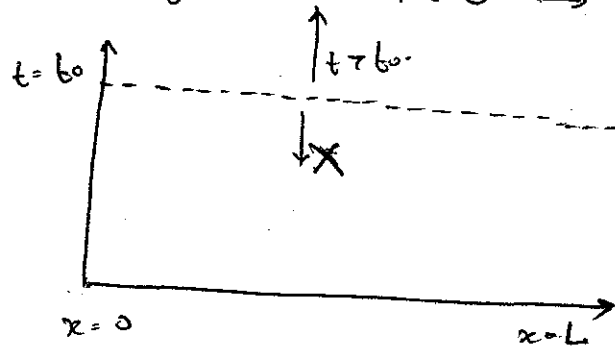Eg of 1D unsteady heat conduction.

A = α.

B = 0

C = 0.

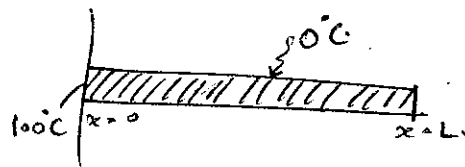B² - 4AC = ~~4αc~~ 0 → parabolic : one characteristic.



t = t₀

t > t₀.

↓X

x = 0          x = L

Say we have a 1D rod at uniform temp. At t = t₀, create a sudden disturbance in T at one end of the rod

Abrupt disturbance at t = t₀ will propagate in a direction forward in time.



50°C

100°C | x = 0          x = L.

The disturbance at t = t₀ will influence what will happen for t > t₀. It cannot influence back what has already happened sometime back.

~~Adiab~~ → time marching problems.

— will have only one type of discontinuity. That discontinuity at reference time at which the disturbance is imposed.

- At time t > t₀, the abrupt disturbance originated at t = t₀ may have made its presence known throughout some part of the domain. That part is called domain of influence.

- The total region in the domain where the presence of disturbance can potentially make its presence known is called domain of disturbance.

t > t₀ → domain of influence

~~to ellip~~ t ≤ t₀ → domain of disturbance

In elliptic case → entire domain is the domain of influence.

Initial-boundary value problem:
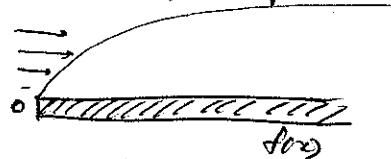~~IVP~~ since $t=t_0$ disturbance introduced

As $t \to t_0$, steady state reached $\Rightarrow \frac{\partial}{\partial t} = 0$ or it becomes elliptic

So it has some elliptic nature build into it.

So more correct way of saying would be that this eq$^n$ is parabolic in time and elliptic in space.

. Another e.g. B/L over a flat plate



- Space marching problems — disturbance at $x=0$. Whatever happens before $x=0$ is not influenced by the disturbance at $x=0$.

Its possible b/c of high Re. High Re $\Rightarrow$ high inertial forces. Inertial force are predominantly uni-directional compared to viscous forces which spread out in all directions. High Re $\Rightarrow$ disturbances predominantly carried uni-directionally.

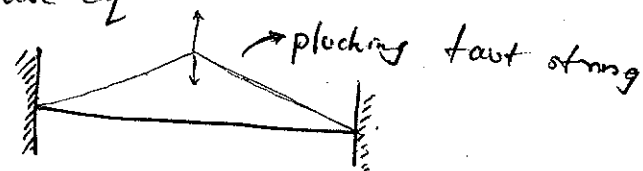(~~that any~~ until & unless ~~not~~ B/L separation occurs).

Ex-3: $\quad \frac{\partial^2 \phi}{\partial t^2} = c^2 \frac{\partial^2 \phi}{\partial x^2} \qquad$ wave eq$^n$



plucking taut string

Use

$x = X$
$t = Y$ $\Big\}$ to formulate in the form $A\phi_{xx} + B\phi_{xy} + C\phi_{yy} = -H$.

$A = +c^2$      |  $B^2 - 4AC = 4c^2$
$B = ~~0~~ 0$     |  $\frac{dy}{dx} \quad \frac{dY}{dX} = \frac{\pm 2c}{2c^2} = \pm \frac{1}{c}$
$C = -1$      |  
            |  $\frac{dt}{dx} = \pm \frac{1}{c}$
            |          $\Rightarrow$
            |  $\therefore \frac{dx}{dt} = \pm c$
            |              $\Rightarrow$

Integrating,

→ Hyperbolic eq$^{ns}$.

$x = \pm ct + c_1$.

Forget $c_1$ (It just shifts. The sol$^n$ by a const. amount everywhere).

main characteristics, $x - ct = \xi$.

$x + ct = \eta$.

Effect is combined spatio-temporal effect.

It's possible to write the entire eq$^n$ in terms of characteristic variables $\xi$ and $\eta$.

---

**Lecture 4: Examples of partial differential equations (contd).**

$$\frac{\partial^2 \phi}{\partial t^2} = c^2 \frac{\partial^2 \phi}{\partial x^2} \qquad : \begin{cases} \xi = x - ct \\ \eta = x + ct \end{cases} \text{two characteristic variables.}$$

$\phi(x,t) \to \phi(\xi, \eta)$.

Sol$^n$ can be written in terms of characteristic variables

$$\frac{\partial \phi}{\partial t} = \frac{\partial \phi}{\partial \xi} \overset{-c}{\frac{\partial \xi}{\partial t}} + \frac{\partial \phi}{\partial \eta} \overset{c}{\frac{\partial \eta}{\partial t}}$$

$$\frac{\partial^2 \phi}{\partial t^2} = \frac{\partial}{\partial \xi}\left[ -c\frac{\partial \phi}{\partial \xi} + c\frac{\partial \phi}{\partial \eta} \right] \overset{-c}{\frac{\partial \xi}{\partial t}} + \frac{\partial}{\partial \eta}\left[ -c\frac{\partial \phi}{\partial \xi} + c\frac{\partial \phi}{\partial \eta} \right] \overset{c}{\frac{\partial \eta}{\partial t}}$$

$$= c^2 \frac{\partial^2 \phi}{\partial \xi^2} + c^2 \frac{\partial^2 \phi}{\partial \eta^2} - 2c^2 \frac{\partial^2 \phi}{\partial \xi \partial \eta} \qquad \cdots ①$$

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial \xi} \overset{1}{\frac{\partial \xi}{\partial x}} + \frac{\partial \phi}{\partial \eta} \overset{1}{\frac{\partial \eta}{\partial x}}$$

$$\frac{\partial^2 \phi}{\partial x^2} = \frac{\partial}{\partial \xi}\left( \frac{\partial \phi}{\partial \xi} + \frac{\partial \phi}{\partial \eta} \right)\frac{\partial \xi}{\partial x} + \frac{\partial}{\partial \eta}\left( \frac{\partial \phi}{\partial \xi} + \frac{\partial \phi}{\partial \eta} \right)\frac{\partial \eta}{\partial x}$$

$$= \frac{\partial^2 \phi}{\partial \xi^2} + \frac{\partial^2 \phi}{\partial \eta^2} + 2\frac{\partial^2 \phi}{\partial \xi \partial \eta} \qquad \cdots ②$$

Apply ① and ② in original wave eq$^n$.

gives, $4\frac{\partial^2 \phi}{\partial \xi \partial \eta} = 0$ or $\frac{\partial^2 \phi}{\partial \xi \partial \eta} = 0$. → $\frac{\partial}{\partial \xi}\left( \frac{\partial \phi}{\partial \eta} \right) = 0$.

$\Rightarrow \frac{\partial}{\partial n}\phi = f_1(n).$

Integrating, $\phi = F(n) + G(\xi).$

Conclusion: general sol$^n$ can be written in terms of characteristic variables.

$IC \to$ At $t = 0$, $\phi = f(x)$.

$\qquad t = 0, \frac{\partial \phi}{\partial t} = g(x)$

$F(x) + G(x) = f(x)$

$cF'(x) - cG'(x) = g(x)$

$\Rightarrow F(x) - G(x) = \frac{1}{c}\int_0^{x} g(\tau)\, d\tau.$

$F(x) = \frac{1}{2c}\int_0^{\frac{1}{2}f(x) + x} g(\tau)\, d\tau + C_k$

$G(x) = \frac{1}{2}f(x) - \frac{1}{2c}\int_0^{x} g(\tau)\, d\tau.$

$\phi = F(n) + G(\xi).$

$\quad = F(x+ct) + G(x-ct)$

$\quad = \frac{1}{2}\left[ f(x+ct) + f(x-ct) \right] + \frac{1}{2c}\int_{x-ct}^{x+ct} g(\tau)\, d\tau$

. Say a disturbance propagates in a fluid medium.

disturbance speed = Sonic speed.

$\qquad c =$ sonic speed.

$\qquad \upsilon =$ speed of the source of disturbance.
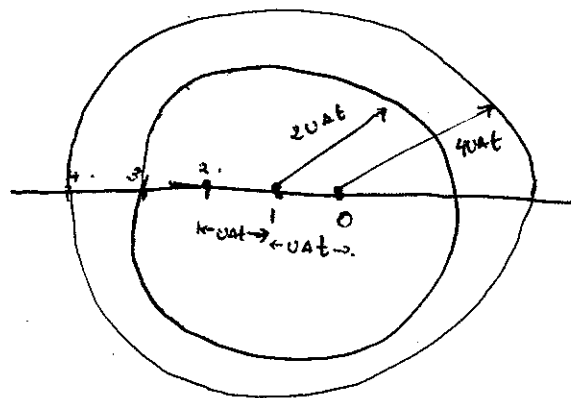
$Ex-1 \quad M_a = 0$ (source of disturbance $\{M_a = \upsilon/c.\}$ doesn't move).



$\qquad \leftarrow$ Take $\Delta t, 2\Delta t, 3\Delta t$ etc....
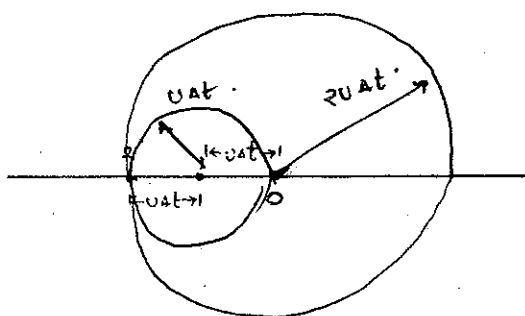
Ex-2    $M_a = \frac{1}{2} \longrightarrow \frac{U}{c} = \frac{1}{2}$
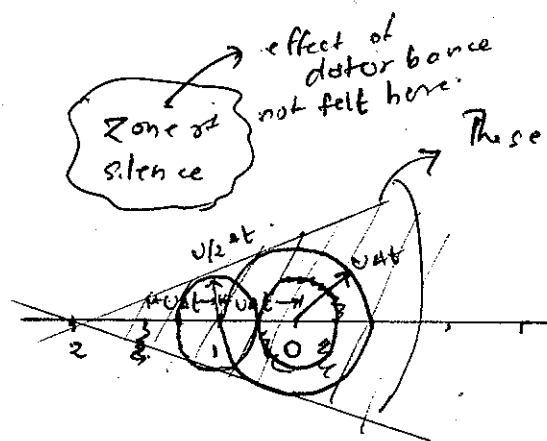


Ex-3    $M_a = 1 \longrightarrow U = c$



Disturbance wave doesn't propagate more than where the source is located.

Ex-4    $M_a > 1$   $M_a = 2 \longrightarrow U = 2c$.



Zone of silence

effect of disturbance not felt here.

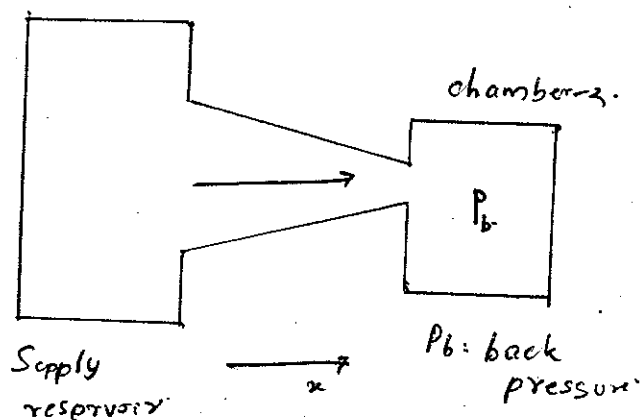These two straight lines are characteristic lines (weak discontinuity)

Mach cone - imaginary cone within which the effects of disturbance is felt.

Highly incompressible flows — disturbance can propagate only with finite speed (unlike in incompressible case where speed is infinite).

As such the effect of disturbance gets accumulated

Say we have a converging nozzle.



chamber-2.

$P_b$.

$P_b$: back pressure.

Supply reservoir

To increase flow rate, $P_b \downarrow$.

$\Rightarrow \dot{m} \uparrow$

$P_b$ regulation essentially creates a disturbance in chamber-2 which is propagated upstream & gives message to supply reservoir to respond to that & send more mass flow.

But what actually happens is That the mass flow rate can be increased upto $Ma = 1$ by decreasing $P_b$. Beyond that limit it cannot be increased.

Explanation:

$\overrightarrow{V_{Df}}$ = velocity of disturbance relative to flow

$= \overrightarrow{V_b} - \overrightarrow{V_f}$

$\overrightarrow{V_f} = +c$

$\overrightarrow{V_{bf}}' = -c$

$\Rightarrow \overrightarrow{V_b} \sim 0$

In hyperbolic cases, when source of disturbance moves faster than the disturbance itself, it leads to discontinuities in the flow medium which has to be taken into account while designing the numerical simulation.

Lecture 5: Nature of the Characteristics of partial differential equation

General $2^{nd}$ order pde. form:

$$\sum_i \sum_j A_{ij} \frac{\partial^2 \phi}{\partial x_i \partial x_j} + B = 0.$$

→ lower order terms

Nature of

Characteristics depends on eigen value of A.
(coeff matrix)

To get eigs, use $\det |A - \lambda I| = 0 \longrightarrow \lambda's.$

. If any $\lambda$ is zero → parabolic

. If none is zero and all $\lambda$s are of the same sign → elliptic

. If none is zero and all but one $\lambda$ is opposite sign → hyperbolic.

Ex $(1 - M_0^2) \frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = 0.$ → eqⁿ relating velocity potential for isentropic inviscid compressible flow over ste bodies with slender shapes.

$A_{11} \frac{\partial^2 \phi}{\partial x_1^2} + \cancel{A_{12} \frac{\partial^2 \phi}{\partial x_1^2}}$

$A_{12} \frac{\partial^2 \phi}{\partial x_1 \partial x_2} + A_{21} \frac{\partial^2 \phi}{\partial x_2 \partial x_1} + A_{22} \frac{\partial^2 \phi}{\partial x_2^2} + B = 0.$

$\begin{aligned} x_1 &= x & A_{11} &= 1 - M_\infty^2 \\ x_2 &= y & A_{12} &= A_{21} = B = 0. \\ & & A_{22} &= 1. \end{aligned}$ $\Bigg\}$ $\begin{vmatrix} (1 - M_\infty^2) - \lambda & 0 \\ 0 & 1 - \lambda \end{vmatrix} = 0$

$(1 - \lambda - M_\infty^2)(1 - \lambda) = 0.$

$\lambda = 1,$
$\lambda = 1 - M_\infty^2$

If $M_\infty = 1 \longrightarrow$ parabolic

$M_\infty < 1 \longrightarrow$ elliptic

$M_\infty > 1 \longrightarrow$ hyperbolic

So the same stream eqⁿ depending on $M_\infty$ can be parabolic, hyperbolic or elliptic.

**Ex.**
- Heat transfer
- Unsteady
- 1D
- low $\dfrac{k}{\rho C_p} \to 0$
- $U = U_\infty$

$$\frac{\partial}{\partial t}(\rho T) + \nabla(\rho \vec{V} T) = \nabla\left(\frac{k}{C\rho}\nabla T\right) + S$$

Using the conditions given,

$$\nabla(\rho \vec{V} T) = \frac{d}{dx}(\cdots)$$

$$\nabla\left(\frac{k}{C\rho}\nabla T\right) \to 0$$

$$\Rightarrow \frac{\partial T}{\partial t} + U_\infty \frac{\partial T}{\partial x} = S$$

$T = T(x,t)$

$$dT = \frac{\partial T}{\partial x}dx + \frac{\partial T}{\partial t}dt$$

$$\begin{bmatrix} 1 & U_\infty \\ dt & dx \end{bmatrix}\begin{bmatrix} \frac{\partial T}{\partial t} \\ \frac{\partial T}{\partial x} \end{bmatrix} = \begin{bmatrix} S \\ dT \end{bmatrix}$$

$\left|\ \ \right| = 0 \longrightarrow dx - U_\infty dt = 0 \Rightarrow \frac{dx}{dt} = U_\infty \leftarrow$ characteristics.

What is the nature of characterstics here? Not parabolic!

(2nd order form of governing) ↑

$$\frac{\partial T}{\partial t} + U_\infty \frac{\partial T}{\partial x} = 0$$

$$\frac{\partial^2 T}{\partial t \partial x} + U_\infty \frac{\partial^2 T}{\partial x^2} = 0 \cdots ①$$

$$\frac{\partial T}{\partial t} + U_\infty \frac{\partial T}{\partial x} = 0$$

$$\frac{\partial^2 T}{\partial t^2} + U_\infty \frac{\partial^2 T}{\partial x \partial t} = 0 \cdots ②$$

$① \times U_\infty - ② \longrightarrow$

$$U_\infty^2 \frac{\partial^2 T}{\partial x^2} - \frac{\partial^2 T}{\partial t^2} = 0$$

$$\frac{\partial^2 T}{\partial t^2} = U_\infty^2 \frac{\partial^2 T}{\partial x^2} \longrightarrow \text{hyperbolic!}$$

**HW:**
$$\left.\begin{array}{l} \dfrac{\partial u}{\partial x} = \dfrac{\partial v}{\partial y} \\[2mm] \dfrac{\partial v}{\partial y} = 4V \end{array}\right\} . \text{ Find the nature of pde:}$$
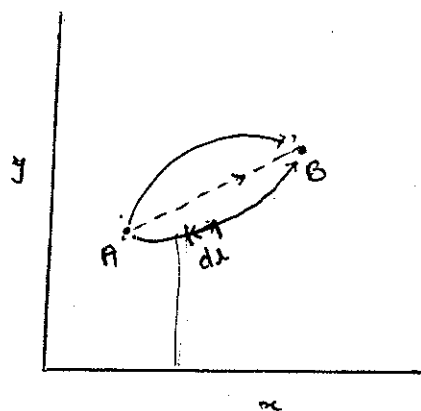
**HW:** Make a chart for the features of parabolic, hyperbolic, elliptic eq⁸

Nature of the characterstics, how many characters, zone of influence, zone of disturbance, speed of propagation of disturbance.

| Lecture 6 : Euler- Lagrangian Equation |
|---|

• Error minimization — key principle with which many numerical methods are founded.

  – B/c a good approximate sol$^n$ is the one which is curs least error.

  – includes many considerations one such is variations.

• Calculus of variations in brief:

  Say we have two pts. Objective: find the path with least distance b/w them.



$$dl = \sqrt{dx^2 + dy^2}$$

$$= \sqrt{1+\left(\frac{dy}{dx}\right)^2} \cdot dx$$

$$y' \equiv \frac{dy}{dx} \quad , \quad y'' \equiv \frac{d^2y}{dx^2}$$

$$l = \int dl = \int \sqrt{1+y'^2} \, dx$$

Problem statement becomes:

  Find the path AB that minimizes $l$

  $\Rightarrow$ minimizes $I = \int \boxed{\sqrt{1+y'^2}} \, dx$

$$I = \int F(x, y, y') \, dx$$

To minimize $I$, take $\delta I$.

$\delta I$ — arbitrarily small virtual change in $I$.

$\delta I$ will only involve changes in dependent variables, not on independent variable (x here).

Why? We are trying to find what the $dy$, should $dy'$ etc.

$I_c$ to minimize $I$, so that our answer fall on the desired path i.e. the straight line.

$F(x, y, y')$

independent variables fns of x ⟶ themselves fns of x

So $F$ is a fn of functions

∴ $F$ is a functional.

$$\delta I = \int \left[ \frac{\partial F}{\partial y} \delta y + \frac{\partial F}{\partial y'} \delta y' \right] dx.$$

$y$ is fixed at $A$ and $B$. $\therefore$ $\delta y = 0$ at $A, B$.

Simplify $2^{nd}$ term $\frac{using}{by}$ integration by parts:

$$\delta I = \int \left[ \frac{\partial F}{\partial y} \right] \delta y + \left[ \frac{\partial F}{\partial y'} \delta y \right]_A^B - \int \frac{d}{dx} \left( \frac{\partial F}{\partial y'} \right) \delta y \, dx.$$

Minimize $I = \int F(x, y, y') dx$ subject to constraint

$y_{@A}$, $y_{@B}$ specified.

$$\therefore 2^{nd} \text{ term } \left[ \frac{\partial F}{\partial y'} \delta y \right]_A^B = 0 \qquad (\delta y_{@A} = \delta y_{@B} = 0).$$

$$\therefore \delta I = \int \left[ \frac{\partial F}{\partial y} - \frac{d}{dx} \left( \frac{\partial F}{\partial y'} \right) \right] \delta y \, dx \}$$

For min $I$, $\delta I = 0$ for any arbitrary $\delta y$.

This is possible when integrand $= 0$.

$$\Rightarrow \frac{\partial F}{\partial y} - \frac{d}{dx} \left( \frac{\partial F}{\partial y'} \right) = 0 \quad \leftarrow \text{Euler - Lagrange equation.}$$

for minimization of dist b/w two pts in this example.

$$I = \int \sqrt{1+y'^2} \, dx \quad \Rightarrow F(x, y, y') = \sqrt{1+y'^2}$$

$$\frac{\partial F}{\partial y} = 0 \quad , \quad \frac{\partial F}{\partial y'} = \frac{2y'}{2\sqrt{1+y'^2}} = \frac{y'}{\sqrt{1+y'^2}}$$

$$\therefore \text{E-L eq}^n: \quad 0 - \frac{d}{dx} \left( \frac{y'}{\sqrt{1+y'^2}} \right) = 0$$

$$\text{or} \quad \frac{d}{dx} \left( \frac{y'}{\sqrt{1+y'^2}} \right) = 0. \quad \Rightarrow \frac{y'}{\sqrt{1+y'^2}} = const$$

$$\Rightarrow y' = const = C$$
$$\text{or} \quad \frac{dy}{dx} = C \quad \rightarrow \text{path is a straight line.}$$

$$\boxed{\text{Lecture 7: Approximate Solutions of Differential Equations}}$$

(Prob) Show that an alternative form of the Euler-Lagrange eq$^n$ is given

by $\quad \dfrac{\partial F}{\partial x} - \dfrac{d}{dx}\left(y'\dfrac{\partial F}{\partial y'}\right) - \dfrac{\partial F}{\partial x} = 0$, where $F(x, y, y')$. ---①

(Ans) $\quad dF = \dfrac{\partial F}{\partial x}dx + \dfrac{\partial F}{\partial y}dy + \dfrac{\partial F}{\partial y'}dy'$

$\Rightarrow \quad \dfrac{dF}{dx} = \dfrac{\partial F}{\partial x} + \dfrac{\partial F}{\partial y}\dfrac{dy}{dx} + \dfrac{\partial F}{\partial y'}\dfrac{dy'}{dx}$

$\dfrac{d}{dx}\left(y'\dfrac{\partial F}{\partial y'}\right) = y'\dfrac{d}{dx}\left(\dfrac{\partial F}{\partial y'}\right) + \dfrac{dy'}{dx}\left(\dfrac{\partial F}{\partial y'}\right)$

$\therefore \quad \dfrac{dy'}{dx}\left(\dfrac{\partial F}{\partial y'}\right) = \dfrac{d}{dx}\left(y'\dfrac{\partial F}{\partial y'}\right) - y'\dfrac{d}{dx}\left(\dfrac{\partial F}{\partial y'}\right)$.

Sub into main expression $\dfrac{dF}{dx}$

$\Rightarrow \quad \dfrac{dF}{dx} = \dfrac{\partial F}{\partial x} + \dfrac{\partial F}{\partial y}\dfrac{dy}{dx}^{y'} + \dfrac{d}{dx}\left(y'\dfrac{\partial F}{\partial y'}\right) - y'\dfrac{d}{dx}\left(\dfrac{\partial F}{\partial y'}\right)$ ----②

$\Rightarrow \quad \dfrac{\partial F}{\partial x} + \dfrac{d}{dx}\left(y'\dfrac{\partial F}{\partial y'}\right) + y'\left[\dfrac{\partial F}{\partial y} - \dfrac{d}{dx}\left(\dfrac{\partial F}{\partial y'}\right)\right]$

$\rightarrow \quad \dfrac{dF}{dx} = \dfrac{d}{dx}\left(y'\dfrac{\partial F}{\partial y'}\right) - \dfrac{\partial F}{\partial x} = 0$

$\circ$ (since this is same as original E-L eq$^n$).

(Prob-2) Hence show that the closed curve that minimizes the perimeter for a given area is a circle.

$P = \displaystyle\int dl = \int\sqrt{dx^2 + dy^2} = \int\sqrt{1+y'^2}\,dx , \quad y' = \dfrac{dy}{dx}.$

$A = A^* = \displaystyle\int y\,dx.$

Obj. minimize $P$ & $A$ is a const.

Do this via Lagrange multiplier.

Introduce $\quad F = P + dA$, where $d$ is Lagrange multiplier.

$$I = \int (P + \lambda A) = \int \left[ \sqrt{1+y'^2} + \lambda y \right] dx$$

$$\underbrace{\qquad\qquad}_{F(x,y,y')}$$

Use alternate form of E-L eq$^n$.

$$\frac{d}{dx}\left[ F - y'\frac{\partial F}{\partial y'} \right] - \cancel{\frac{\partial F}{\partial x}}^{0} = 0$$

$$\Rightarrow F - y'\frac{\partial F}{\partial y'} = \text{const } c$$

$$\sqrt{1+y'^2} + \lambda y - \frac{(y')^2}{\sqrt{1+(y')^2}} = c$$

$$\Rightarrow 1 + \cancel{y'^2} + \lambda y \sqrt{1+y'^2} - \cancel{(y')^2} = c\sqrt{1+(y')^2}$$

$$1 = (c - \lambda y)\sqrt{1+y'^2}$$

Use $y' = \tan\theta = \frac{dy}{dx}$.

$$\cos\theta = c - \lambda y.$$

$$\frac{dy}{dx} = \frac{dy/d\theta}{dx/d\theta}$$

$$\frac{d}{d\theta}(\cos\theta) = \frac{d}{d\theta}(c - \lambda y)$$

$$-\sin\theta = -\lambda\frac{dy}{d\theta}.$$

$$\text{So, } \frac{dy}{d\theta} = \frac{1}{\lambda}\sin\theta$$

$$\text{So, } \frac{dy}{dx} = \frac{\frac{1}{\lambda}\sin\theta}{\left(\frac{dx}{d\theta}\right)}.$$

$$\tan\theta = \frac{\frac{1}{\lambda}\sin\theta}{\left(\frac{dx}{d\theta}\right)}$$

$$\therefore \frac{dx}{d\theta} = \frac{1}{\lambda}\cos\theta.$$

$$\therefore x = \frac{\sin\theta}{\lambda} + c_1.$$

$$\sin\theta = \lambda(x - c_1)$$

$$\cos\theta = c - \lambda y$$

$$\sin^2\theta + \cos^2\theta = 1$$

$$\Rightarrow [\lambda(x - c_1)]^2 + [c - \lambda y]^2 = 1.$$

$$\lambda^2(x - c_1)^2 + \lambda^2(y - c_2)^2 = 1.$$

This is of the form.

$$(x - a)^2 + (y - b)^2 = r^2$$

$\longrightarrow$ eq$^n$ of circle

Thus circle minimizes perimeter for const area shapes!

Ex/ Functionals involving higher order derivatives:

Say, $I = \int F(x, y, y', y'') \, dx$

$\longrightarrow \left[\frac{\partial F}{\partial y''} * \delta y'\right]_A^B - \int \frac{d}{dx}\left(\frac{\partial F}{\partial y''}\right)\delta y' \, dx$

$\delta I = \int\left[\frac{\partial F}{\partial y}\delta y + \frac{\partial F}{\partial y'}\delta y' + \frac{\partial F}{\partial y''}\delta y''\right] dx$

$- \left[\frac{d}{dx}\left(\frac{\partial F}{\partial y''}\right)\delta y\right]_A^B + \int \frac{d^2}{dx^2}\left(\frac{\partial F}{\partial y''}\right)\delta y \, dx$

[ to denote first order variation (usually omitted) ]

$\longrightarrow \left[\frac{\partial F}{\partial y'}\delta y\right]_A^B - \int \frac{d}{dx}\left(\frac{\partial F}{\partial y'}\right)\delta y \, dx$

$\left[\frac{\partial F}{\partial y''}\delta y'\right]_A^B - \int$

Boundary terms:

$$\left[\frac{\partial F}{\partial y'}\delta y\right]_A^B, \quad \left[\frac{\partial F}{\partial y''}\delta y'\right]_A^B, \quad -\left[\frac{d}{dx}\left(\frac{\partial F}{\partial y'}\right)\delta y\right]_A^B$$

Boundary terms $= 0$ if $\delta y, \delta y' = 0$.

So assume $y, y'$ are specified at A & B :

$\delta y, \delta y' = 0$.

$\Rightarrow \delta I = \int\left[\frac{\partial F}{\partial y} - \frac{d}{dx}\left(\frac{\partial F}{\partial y'}\right) + \frac{d^2}{dx^2}\left(\frac{\partial F}{\partial y''}\right)\right]\delta y \, dx = 0$

$\delta I = \int\left[\frac{\partial F}{\partial y} + \frac{d^2}{dx^2}\left(\frac{\partial F}{\partial y''}\right)\right]\delta y \, dx = 0$

$\Rightarrow \quad \frac{\partial F}{\partial y} - \frac{d}{dx}\left(\frac{\partial F}{\partial y'}\right) + \frac{d^2}{dx^2}\left(\frac{\partial F}{\partial y''}\right) = 0 \quad$ to minimize $\underline{I}$.

For more higher order eq?, we have the form:

$$\left\{\frac{\partial F}{\partial y} - \frac{d}{dx}\left(\frac{\partial F}{\partial y'}\right) + \frac{d^2}{dx^2}\left(\frac{\partial F}{\partial y''}\right) - \frac{d^3}{dx^3}\left(\frac{\partial F}{\partial y'''}\right) + \frac{d^4}{dx^4}\left(\frac{\partial F}{\partial y''''}\right) - \cdots = 0\right\}$$

$\hookrightarrow$ Euler-Poisson form.

Approximate solutions of differential equations through variational formulation:

Eg: 1D steady state heat conduction with constant heat source:

$$\frac{\partial}{\partial t}(\rho T) + \nabla(\ell \vec{V} T) = \nabla\left(\frac{k}{C_p} \nabla T\right) + \frac{S}{C_p}$$

Also assume: const. Thermal properties $k, C_p$ etc const.

Steady state $\Rightarrow \frac{\partial}{\partial t}(\cdots) = 0$

Conduction problem $\Rightarrow$ no flow velocity involved.

$$\therefore \nabla(\ell \vec{V} T) = 0$$

1D Problems $\Rightarrow \nabla() = \frac{d}{dx}()$

Final form: $\frac{d}{dx}\left(k \frac{dT}{dx}\right) + S = 0. \longrightarrow$ "D" form (Differential form) $\hookrightarrow$ Strong form.

To make it a variational form, multiply with a variational parameter & integrate over the domain.

$$\int \left[\frac{d}{dx}\left(k\frac{dT}{dx}\right) + S\right] v \, dx = 0$$

$\hookrightarrow$ carries the meaning of $\delta T$.

Use integration by parts:

$$\left[v k \frac{dT}{dx}\right]_1^2 - \int k \frac{dv}{dx} \frac{dT}{dx} dx + \int S v \, dx = 0 \quad \text{? should ... }$$

Boundary conditions possible:

$\longrightarrow T$ specified.
$\quad \rightarrow \delta T = 0$ or $v = 0$ $\Big\}$ either of the two.
$\longrightarrow \frac{dT}{dx}$ specified.

Ex: $T$ specified at both boundaries.

$$\Rightarrow \boxed{\int k \frac{dT}{dx} \frac{dv}{dx} dx = \int S v \, dx} \longrightarrow \text{Weak form.}$$
$$\text{"V" form.}$$

## Lecture 8: Variational formulation

In the weak form, it requires continuity only upto the first order derivative, while in the strong form, continuity upto the second order derivative is required. Hence why the name — weak form

Generic form of "V" form:

$$a(T, v) = l(v)$$

bilinear operator    linear operator

why this form?

Make some observations:

obs-1

$$a(\alpha_1 T + \alpha_2 v, \beta_1 T + \beta_2 v) = \int k \frac{d}{dx}(\alpha_1 T + \alpha_2 v) \frac{d}{dx}(\beta_1 T + \beta_2 v) \, dx$$

$$= \alpha_1 \beta_1 \int \frac{k \, dT}{dx} \frac{dT}{dx} \, dx + \alpha_1 \beta_2 \int k \frac{dT}{dx} \frac{dv}{dx} \, dx + \alpha_2 \beta_1 \int k \frac{dv}{dx} \frac{dT}{dx} \, dx +$$

$$\alpha_2 \beta_2 \int k \frac{dv}{dx} \frac{dv}{dx} \, dx$$

$$= \alpha_1 \beta_1 \, a(T, T) + \alpha_1 \beta_2 \, a(T, v) + \alpha_2 \beta_1 \, a(v, T) + \alpha_2 \beta_2 \, a(v, v).$$

$\Rightarrow$ $a(,)$ is bilinear (i.e. linear in each slot)

obs-2

$$l(\alpha T + \beta v) = \int s (\alpha T + v) \, dx$$

$$= \alpha l(T) + \beta l(v), \text{ if such a property is satisfied,}$$

$l$ is a linear operator

obs-3.

$$a(T, v) = a(v, T) \Rightarrow a \text{ is symmetric (self adjoint)}$$

obs-4

$$a(v, v) = \int k\left(\frac{dv}{dx}\right)^2 dx. \text{ is a +ve definite operator}$$

$\downarrow$ integral over the domain is +ve
$\geq 0$.

$a$ is a scalar product on V.

Say we have a function $g(\varepsilon) = \frac{1}{2} a(T+\varepsilon v, T+\varepsilon v) - l(T+\varepsilon v)$

M) problem : Minimize $g$ at $\varepsilon = 0 \longrightarrow \frac{1}{2} a(T,T) - l(T)$

Assumptions:

$\longrightarrow$ $a$ is bilinear & $l$ is linear

$$g(\varepsilon) = \frac{1}{2} a(T,T) + \frac{\varepsilon}{2} a(T,v) + \frac{\varepsilon}{2} a(v,T) + \frac{\varepsilon^2}{2} a(v,v)$$
$$- l(T) - \varepsilon l(v)$$

To minimize '$g$' , $g'(\varepsilon)\big/_{\varepsilon=0} = 0$

$$g'(\varepsilon)\big|_{\varepsilon=0} = \frac{1}{2} a(T,v) + \frac{1}{2} a(v,T) + \varepsilon a(v,v)^{0} - l(v) = 0.$$

$$\frac{1}{2} (a(T,v) + a(v,T)) = l(v) \cdot$$

From obs-g , $\rightarrow$ $a(T,v) = l(v)$ which is same as '$v$' form.

i.e. $a$ is symmetric $\Rightarrow$ $a(T,v) = a(v,T)$

$$\Rightarrow a(T,v) = l(v) \rightarrow \text{``}v\text{''} \text{ form.}$$

Conclusion: 'M' form $\Rightarrow$ 'v' form , provided $a$ is bilinear, $l$ is linear & $a$ is symmetric.

Question B: starting from 'v' form, is it possible to reach 'M' form?

Consider $a$ as symmetric.

$$g(\varepsilon) = \frac{1}{2} a(T,T) - l(T) + \varepsilon [a(T,v) - l(v)] + \frac{\varepsilon^2}{2} a(v,v)$$

If 'v' form is true, $\varepsilon [a(T,v) - l(v)] = 0$

Then $g(\varepsilon) \geq \frac{1}{2} a(T,T) - l(T)$; provided $a(v,v)$ is +ve.

It is true when $a$ is +ve definite +ve

$$\therefore \quad g(\varepsilon) = \frac{1}{2} a(T,T) - l(T) \quad \text{is the maximum of } g(\varepsilon)$$

for v to M, we require additional constraint re: $a$ is +ve definite .

$\hookrightarrow$ 'M' form

at $\varepsilon = 0$.

$$\int k \frac{dT}{dx} \frac{dv}{dx} \, dx = \int S v \, dx$$

sub $\int T$ in place of $v$,

$$\delta \int \left[ \frac{1}{2} k \cdot \left(\frac{dT}{dx}\right)^2 dx \mp \int S T \right] dx = 0.$$

$$\underbrace{\hspace{5cm}}_{\pi}$$

Essentially, we are minimizing $\pi$.

Q) How to get 'D'-form from 'V' form?

'V' form: $\int k \underset{1}{\underline{\frac{dT}{dx}}} \underset{2}{\underline{\frac{d(\delta T)}{dx}}} \, dx = \int S \, \delta T \, dx.$

Integrate by parts:

$$k \cdot \frac{dT}{dx} \delta T \Big]^2 - \int k \frac{d}{dx} \, k \frac{d}{dx} \left( k \frac{dT}{dx} \right) v \, dx = \int_{3}^{2} S v \, dx.$$

$$\int \left[ \frac{d}{dx} \left( k \frac{dT}{dx} \right) \varphi + S \right] v \, dx = 0$$

$$\to \quad \frac{d}{dx} \left( k \frac{dT}{dx} \right) + S = 0 \longrightarrow \text{'D' form}.$$

So far we took BC as $T$ is specified. we can do the same

with $\frac{dT}{dx}$ specified as a instead.

- Boundary conditions in the variational formulation:

$T$ specified : variable for which variation appears in the boundary terms

B/c $v$ is variation of $T$.  (primary variable)

$k \frac{dT}{dx}$ specified : $\longrightarrow$ coeff. of variation in the boundary term

(secondary variable)

Specifying the primary variable at boundary is called as

essential BC.

Specifying the secondary variable at boundary is called as

natural BC. Natural BC that terms automatically.

appears in the eqⁿ ( -ve of heat flux).

**Lecture-9:** Example of Variational formulation and introduction to Weighted Residual Method.

**Eg.:**

$$\frac{d^2}{dx^2}\left[a(x)\frac{d^2y}{dx^2}\right] + b(x) = 0.$$

Obj → cast in variational formulation.

$$\int \left[\frac{d^2}{dx^2}\left[a(x)\frac{d^2y}{dx^2}\right] + b(x)\right] v\, dx = 0.$$

Integrate by parts:

$$v \oint \frac{d}{dx} \quad \left. v\,\frac{d}{dx}\left[a(x)\frac{d^2y}{dx^2}\right]\right|_{x=0}^{x=L} - \int \frac{dv}{dx}\frac{d}{dx}\left[a(x)\frac{d^2y}{dx^2}\right] dx$$

$$+ \int b(x)\, v\, dx = 0.$$

$$\left.v\frac{d}{dx}\left[a(x)\frac{d^2y}{dx^2}\right]\right|_{x=0}^{x=L} \quad - \quad \left.\frac{dv}{dx}\left(a(x)\frac{d^2y}{dx^2}\right)\right|_{x=0}^{x=L} \quad + \quad \int_{x=0}^{x=L} \frac{d^2v}{dx^2}\,a(x)\frac{d^2y}{dx^2}\, dx. \qquad a(y,v)$$

$$+ \int_{x=0}^{x=L} b(x)\, v\, dx = 0.$$

$$\ell(v)$$

Primary variable → $y$.  (E.B.C)

Secondary variable → $\dfrac{d}{dx}\left[a(x)\frac{d^2y}{dx^2}\right]$

(N.B.C)

Primary variable : $\dfrac{dy}{dx}$ (E.B.C)

Secondary variable : $a(x)\dfrac{d^2y}{dx^2}$

(N.B.C)

Boundary terms need not be zero after applying B/Cs.

The terms that we get after BC application can be dubbed with $\ell(v)$ term.

Specifying B/Cs.

Let $y = 0$ at $x = 0$.

$\dfrac{d}{dx}\left[a(x)\dfrac{d^2y}{dx^2}\right] = c_1$ at $x = L$.

$\dfrac{dy}{dx} = 0$ at $x = 0$

$a(x)\dfrac{d^2y}{dx^2} = c_2$ at $x = L$.

$$\Rightarrow \quad V_L C_1 - \frac{dv}{dx}\Big|_L C_2 + \int_{x=0}^{x=L} a_{(x)}\frac{d^2v}{dx^2}\frac{d^2y}{dx^2}\,dx + \int_{x=0}^{x=L} b_{(x)}v\,dx = 0.$$

a) Write this in the form $\overset{A}{\Phi}(y,v) = L(v)$

$$A(y,v) = \int_{x=0}^{x=L} a_{(x)}\frac{d^2y}{dx^2}\frac{d^2v}{dx^2}\,dx.$$

$$L(v) = -\int_{x=0}^{x=L} b_{(x)}v\,dx - V_L C_1 + \frac{dv}{dx}\Big|_L C_2$$

$A(y,v) = L(v)$  is the required variational formulation.

## Approximate solutions of differential equations:

Weighted residual approach:

: gets a clue from the variational form.

$$\int(\ )v = 0 \quad \text{form.} \quad v \text{ is an arbitrarily small variation.}$$

Say we wish to solve $\dfrac{d^2y}{dx^2} = 0$

Call linear operator $L(y) = 0 \longrightarrow$ of $L = \dfrac{d^2}{dx^2}$.

$v$ till now is an abstract variational parameter. We try make it non-abstract by looking at the possibilities of functions we can use in place of $v$.

In the problem $\dfrac{d^2y}{dx^2} = 0$, to convert it into an algebraic eq$^n$ ( b/c algebraic are easier to solve), we replace $y$ with an approximate $y_{approx}$ polynomial & solve for it.

But $\dfrac{d^2}{dx^2}(y_{approx}) \neq 0$, in general.

Then, $L(y) - L(y_{approx}) = R$    $R \rightarrow$ residual

Our objective is to minimize $R$ in an integral sense over the domain.

$$\int_\Omega R\omega \, d\Omega = 0$$

Try to minimize the error or the residual in a weighted integral sense.

$\rightarrow y^*$ $(y_{approx})$

[trial function]

$\rightarrow \omega$

[weighting function]

Note : these functions need not be as rigourous as the variation. formulations. Restrictions like $a(y,v)$ symmetry & positive definiteness not needed.

$$\boxed{\text{Lecture 10: Weighted Residual Method.}}$$

Say, governing differential eq? $L(y) = 0$.

Substituting $y_{approx}$, $L(y_{approx}) \neq 0$.

$$L(y_{approx}) = R.$$

$$\int_\Omega R\omega \, d\Omega = 0.$$

In 1D problem, say $y = f(x)$,

$$\int R\omega \, dx = 0.$$    try to make sure $y_{approx}$ is appropriate to minimise the R.

Trial function $\longrightarrow y_{approx}$,

- polynomial — most convenient form
- Should satisfy the essential BC (key requirement)
- should be continuous
- derivatives of trial function must be square integrable.
$\int \left(\frac{d y_{approx}}{dx}\right)^2 dx < \infty$ — shows integral is not unbounded.

+ $H'$ fn: $1^{st}$ derivative is square integrable.

Requirements for weighting function: $\longrightarrow \omega$

- should satisfy homogeneous part of the EBC.

  {Say if $y = 5$ is EBC, then $\omega = 0$ is the homogeneous part}

  why? B/c if $y$ is specified, then variation is $y = 0$. $\omega$ has similar meaning as that of variation in $y$. $\therefore \omega = 0$. at Boundary

- Should be continuous.

## Some specific examples:

Prob: 1D, steady state heat transfer with uniform Thermal conductivity $k$, source $S$.

Governing DE: $\quad k\dfrac{d^2 T}{dx^2} + S = 0$.

$$\frac{d^2 T}{dx^2} + \left(\frac{S}{k}\right) = 0 . \quad \text{100, say}$$

BCs: $\quad x = 0, \; T = 0$
$\quad\quad\quad x = 10, \; T = 0$.

$\longrightarrow$

| $\dfrac{d^2 y}{dx^2} + 100 = 0$ | At $x = 0, y = 0$ |
|---|---|
| | $x = 10, y = 0$ |

Obj: find approx. sol$^n$.

## Ex-1 Least squared method

$y \longrightarrow y_{approx}$

$$\frac{d^2 y_{approx}}{dx^2} + 100 = R$$

Interested in minimizing $R^2$

$\longrightarrow \int R^2 dx \longrightarrow$ minimized.

$\equiv$ sum of square of errors.

$y_{approx} \longrightarrow 2^{nd}$ order polynomial with $1$ parameter.

$\left\{ \begin{array}{l} \int R \, dx \quad\text{—— not correct as} \\ \text{some errors can be +ve \&} \\ \text{-ve \& sums to } 0, \text{ can be} \\ \text{misleading} \\ \cdot \int |R| \, dx \text{— not used b/c of} \\ \text{tedious algebra.} \end{array} \right\}$

General form: $ax(10-x) = y_{approx}$

$\longrightarrow$ polynomial, $2^{nd}$ order

Find out $a$ s $\int R^2 dx$ is minimized.

$y_{approx} = ax(10-x)$

$\qquad = 10ax - ax^2$

$\dfrac{dy_{approx}}{dx} = 10a - 2ax$

$\dfrac{d^2 y_{approx}}{dx^2} = -2a$

$R = y - y_{approx}$

$\qquad = -2a + 100$

$\dfrac{\partial}{\partial a}\left(\int R^2 dx\right) = 0.$

$\Rightarrow \int 2R\,\dfrac{\partial R}{\partial a}\,dx = 0.$

$\Rightarrow \int R\,\dfrac{\partial R}{\partial a}\,dx = 0.$

So $\dfrac{\partial R}{\partial a}$ is the weighing function.

$\dfrac{\partial R}{\partial a} = \dfrac{\partial}{\partial a}(-2a+100) = -2$

$\therefore \displaystyle\int_0^{10} (-2a+100)(-2)\,dx = 0$

$\qquad \Rightarrow -2a+100 = 0$

$\qquad\qquad$ or $a = 50$

| Lecture 11: Point Collocation method, Galerkin's method & The 'M' form. |
| --- |

Ex-2. Point Collocation method:

$$\omega = \delta(x - x_i) \qquad \delta : \text{Dirac-Delta function}.$$

Idea: you try to satisfy the value of the function at chosen points $x_i$

$$\int_x R \omega \, dx = 0.$$

Keep trial function same: $ax(10-x)$

$$R_{x=x_i} = 0 \qquad \text{Consider only 1 collocation point, say } x=5.$$

$$R_{x=5} = 0$$

$$-2a + 100 = 0 \Rightarrow a = 50$$

Ex-3 Galerkin's method:

It considers the weighting function as The trial function.

$$\omega = x(10-x) \qquad \text{Putting } a \text{ doesnt matter as its just a const & it will go away in The expression } \int R\omega \, dx = 0.$$

$$\int R\omega \, dx = 0.$$

$$\int_0^{10} (-2a+100)\, x(10-x)\, dx = 0$$

$$\longrightarrow \quad -2a+100 = 0.$$

$$\text{or} \quad a = 50$$

Going through routes of the 'M' form:

$$\int_0^{100} \left(\frac{d^2y}{dx^2} + 100\right) v \, dx = 0$$

$$\left[ v\frac{dy}{dx} \right]_0^{100} - \int_0^{100} \frac{dv}{dx}\frac{dy}{dx}\, dx + \int 100\, v\, dx = 0$$

PV: $y$     $y=0$ at $x=0$    $\Rightarrow \int_0^{100} \frac{dv}{dx}\frac{dy}{dx}\, dx = \int_0^{100} 100\, v\, dx.$

SV: $\frac{dy}{dx}$     (given)

$$a(y,v) = \ell(v)$$

$a(y,y) > 0$. here $\Rightarrow$ 'M' form exist here!

N) form:

Minimize $\Pi = \frac{1}{2} a(y, y) - \ell(y)$.

$$= \frac{1}{2} \int_0^{10} \left(\frac{dy}{dx}\right)^2 dx - \int_0^{10} 100y\, dx$$

Use $y_{approx} \longrightarrow$ in place of $y$ & minimize $\Pi$, & can also have reduced requirement for continuity.

Here we don't require any weighting $fn$. Only $fnal$ function is necessary. This convenience comes at the cost of additional requirement of positive definiteness of $a(y,y)$.

$\dfrac{\partial \Pi}{\partial a} = 0$     Use $fnal$ $fn$. $y_{approx} = ax(10-x)$.

— Rayleigh – Ritz method.

$$\frac{dy_{approx}}{dx} = 10a - 2ax$$

$$\Pi = \frac{1}{2} \int_0^{10} \left(\frac{dy_{approx}}{dx}\right)^2 dx - \int_0^{10} 100y_{approx}\, dx$$

$$= \frac{1}{2} \int_0^{10} (10-2ax)^2 dx - \int_0^{10} 100(10a-2ax)\, dx$$

$$\approx \frac{d}{da} \oint_{0}^{10} \frac{\pi}{\partial a} \left(500 + \frac{2a^2}{3}\times1000 - \frac{100}{x} \right) + 10000a + 4000a = 0$$

$$= \frac{4a}{3}\times 500 = 100$$

$\dfrac{\partial \Pi}{\partial a} = 0 \implies a = 100 \int_0^{}$

$$\frac{dy}{dx} = -100x + C_1$$

$$y = -\frac{100x^2}{2} + C_1 x + C_2.$$

$C_2 = 0 \longrightarrow 0 = -\frac{100\times1000}{2} + 10C_1 \longrightarrow C_1 = 500.$   exact!

$$y = 50(10x - x^2)$$
$$= 50 x (10-x) \checkmark.$$

How to reduce calculations associated with higher order polynomials?

Divide the domain into smaller & smaller subdomains & use lower order polynomials in those subdomains. B/c even though across the entire domain our $f_n$ may be quite complex, in smaller subdomains, such functions may be simpler. Now the solution will be fitted b/w discrete points rather than over the entire domain ⟶ basic idea behind discretization.

↗ here we lose the continuous nature of the domain.

<u>Discretization</u>:

- Divide the domain into a number of discrete subdomains.
  (element, control volume, ⋯)

- Each subdomain is represented by a discrete set of points.
  (grid points, nodes, ⋯)

- Objective is to convert the governing DE into a system of algebraic equations valid at each of these discrete points.

---

| Lecture 12: Finite Element Method (FEM) of discretization |
| --- |

Discretization principles:

→ Divide the domain into a number of discrete subdomains, each subdomain being ~~characterized~~ represented by a number of discrete points.

→ Derive algebraic equations from the governing diff. eq^ns, valid at these discrete points.

→ Solve the system of algebraic equations to obtain values of the dependent variables at the discrete points.

---

→ Broad steps in overall analysis:
  → Pre-processing: set-up geometry, discretized eq^ns, input data (property data), initial cond, BCs.
  → Solution: algebraic eq^ns
  → Post-processing: Graphical representation of the obtained results.

# Finite Element Method (FEM)

$$Ex: \frac{d}{dx}\left(k\frac{dT}{dS}\right) + S = 0$$

1, 2, 3, 4 → nodes

①, ②, ③ → elements.

Prepare node-element connectivity chart:

| Element | node i | node j |
|---------|--------|--------|
| 1 | 1 (,) | 2 (,) |
| 2 | 2 (,) | 3. (,) |
| 3 | 3 (,) | 4. (,) |

Consider any isolated element

Writing an algebraic eq? corresponding to the governing diff eq?.

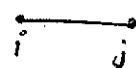$$\int_{x=0}^{x=L}\left[\frac{d}{dx}\left(k\frac{dT}{dx}\right) + S\right]w\,dx = 0$$

Integrate by parts:

$$\omega\, k\frac{dT}{dx}\Big]_{0}^{L} - \int_{0}^{L}\frac{d\omega}{dx}k\frac{dT}{dx}\,dx + \int_{0}^{L}Sw\,dx = 0.$$

for the isolated element, this would become.

$$\omega k\frac{dT}{dx}\Big]_{x_i}^{x_j} - \int_{x_i}^{x_j}k\frac{d\omega}{dx}\frac{dT}{dx}\,dx + \int_{x_i}^{x_j}Sw\,dx = 0.$$

If it were not a two-noded element, we'd require a higher order polynomial to approximate T.

$$T = a_0 + a_1 x \quad \text{(trial function)}$$

At $x = x_i$, $T = T_i$

$x = x_j$, $T = T_j$

$$\left.\begin{array}{c}T_i = a_0 + a_1 x_i \\ T_j = a_0 + a_1 x_j\end{array}\right\} \rightarrow$$

$$a_1 = \frac{T_j - T_i}{x_j - x_i}$$

$$a_0 = T_i - a_1 x_i = T_i - \left(\frac{T_j-T_i}{x_j-x_i}\right)x_i$$

We choose trial function for each element, not for the whole domain.

We finally get piecewise continuous function for the whole domain.

$$a_0 = T_i - \left(\frac{T_j - T_i}{x_j - x_i}\right) x_i \qquad \Big| \qquad a_1 = \frac{T_j - T_i}{x_j - x_i}$$

$$= \frac{T_i x_j - T_j x_i}{x_j - x_i}$$

$$T = \left(\frac{T_i x_j - T_j x_i}{x_j - x_i}\right) + \left(\frac{T_j - T_i}{x_j - x_i}\right) x$$

$$= \underbrace{\frac{(x_j - x)}{(x_j - x_i)}}_{} T_i + \underbrace{\left(\frac{x - x_i}{x_j - x_i}\right)}_{} T_j =$$

$$T = \underbrace{N_i T_i + N_j T_j}_{}$$

$\longrightarrow$ Interpolation functions / shape functions.

Property of shape functions:

$$N_i = 1 \text{ at node } i \quad, \quad = 0 \text{ at node } j$$
$$N_j = 0 \text{ at node } j \quad, \quad = 1 \text{ at node } i.$$

Writing in matrix form

$$T = \underbrace{[N_i \quad N_j]}_{[N]} \underbrace{\begin{bmatrix} T_i \\ T_j \end{bmatrix}}_{[T]} =.$$

$$\omega = [N] \overset{[\omega]}{\underset{\hat{}}{\cancel{\omega}}} \quad \text{is Galerkin form.}$$
$$\longrightarrow \begin{bmatrix} \omega_i \\ \omega_j \end{bmatrix}$$

$$\Rightarrow \omega^T = [\omega]^T [N]^T$$

$$\therefore \quad \omega k \frac{dT}{dx}\Big]_{x_i}^{x_j} - \int_{x_i}^{x_j} k \frac{d\omega}{dx} \frac{dT}{dx} dx + \int_{x_i}^{x_j} S\omega \, dx = 0.$$

$$\rightarrow [\omega]^T [N]^T q''\Big]_i^j - [\omega]^T \int_{x_i}^{x_j} \left[\frac{dN^T}{dx}\right] k \left[\frac{dN}{dx}\right] \overset{[T]}{\cancel{\text{\tiny }}} dx$$

$$+ \cancel{\text{\tiny }} + [\omega]^T \int S [N]^T dx = 0.$$

Since $[w]^T$ is arbitrary,

$$\underbrace{\left[-[N]^T q^u\right]_i^j}_{Term-1} - \underbrace{\left[\int_i^j \left[\frac{dN}{dx}\right]^T k \left[\frac{dN}{dx}\right]\right][T]}_{Term-2.} + \int_i^j S[N]^T dx = 0$$

$Term.1 = \left[-\begin{bmatrix} N_i \\ N_j \end{bmatrix} q^u\right]_i^j$

$= -\begin{bmatrix} 0 \\ q^u_j \end{bmatrix} + \begin{bmatrix} q^u_i \\ 0 \end{bmatrix} = \begin{bmatrix} +q^u_i \\ -q^u_j \end{bmatrix}$

$Term-2 = \int_i^j \begin{bmatrix} \frac{dN_i}{dx} \\ \frac{dN_j}{dx} \end{bmatrix} k \begin{bmatrix} \frac{dN_i}{dx} & \frac{dN_j}{dx} \end{bmatrix} dx$

$\frac{dN_i}{dx} = \frac{-1}{\underbrace{x_j - x_i}_{l_e}} \qquad \frac{dN_j}{dx} = \frac{1}{\underbrace{x_j - x_i}_{l_e}}$

$= \frac{-1}{l_e} \qquad\qquad = \frac{+1}{l_e}$

$\Rightarrow \quad Term-2 = \int_i^j \frac{k}{l_e^2} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \end{bmatrix} dx$

$= \int_i^j \frac{k}{l_e^2} \begin{bmatrix} +1 & -1 \\ -1 & +1 \end{bmatrix} dx$

$= \frac{k}{l_e^2} (x_j - x_i)$

$= \frac{k}{l_e} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$

$Term-3 : \int_i^j S[N]^T dx = S \int_i^j \begin{bmatrix} \frac{x_j - x}{l_e} \\ \frac{x - x_i}{l_e} \end{bmatrix} dx = \begin{Bmatrix} \frac{S l_e}{2} \\ \frac{S l_e}{2} \end{Bmatrix}$

→ entire effect of the element is manifested by the behavior of nodes.
Total is shared equally b/w nodes.

Assembling all terms together,

$$(\text{Term-2})\cdot[T] = (\text{Term-1}) + (\text{Term-3})$$

$$\Rightarrow \quad \frac{k}{le}\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}\begin{bmatrix} T_i \\ T_j \end{bmatrix} = \begin{bmatrix} q_i'' \\ -q_j'' \end{bmatrix} + Sle\begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}.$$



$$T = T_L$$

This is of the form:

→ Acts like a force term.

$$[\bar{k}][T] = [F]$$ → Similar to spring mass system.

Stiffness matrix.

as if like the stiffness of the system

as if behaves like displacement

for 2 nodes we get 2×2 form. Similarly for 4 nodes we get 4×4 form.

This part of coeff. matrix activated for 1st element.



$$\begin{bmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{bmatrix} = \begin{bmatrix} q_1'' \\ -q_2'' + q_2'' \\ -q_6'' + q_3'' \\ -q_4'' \end{bmatrix} + Sle\begin{bmatrix} \frac{1}{2} \\ \frac{1}{2}+\frac{1}{2} \\ \frac{1}{2}+\frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

This part for the third element

→ This part activated for 2nd element.

Important assumptions: ① Thermal conductivity $k$ as constant.

② All elements have same length. It need not be the case in reality.

Final form:

$$\begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}\begin{bmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{bmatrix} = \begin{bmatrix} q_1'' \\ 0 \\ 0 \\ q_2'' \end{bmatrix} + Sle\begin{bmatrix} 1/2 \\ 1 \\ 1 \\ 1/2 \end{bmatrix}$$

Complete the problem by imposing BCs.

At $x=0$, $q_1'' = 0$ (heat flux 0 zero as insulated).

At $x=L$, $\cancel{\frac{q}{4}}$ $T_4$ is specified.

Even if $T_4$ is specified, computer may try to find $T_4$ $\cancel{\text{fr.}}$ with $\cancel{\text{r}}$ first BC. In that case, we require matching b/w specified $T_4$ & computed $T_4$.

$$k_{41} T_1 + k_{42} T_2 + k_{43} T_3 + k_{44} T_4 = R.$$

Let $T_{4,\text{specified}} = T_4^*$.

Use small computational trick.

Replace $R$ with $\overset{L}{\cancel{k_{44}}} T_4^*$ & replace $k_{44}$ with $k_{44} + L$, where $L$ is a large number.

Then $k_{41} T_1 + k_{42} T_2 + k_{43} T_3 + (k_{44}^{+L}) T_4 = L T_4^*$.

$$T_4 = \frac{-k_{41} T_1}{k_{44} + L} - \frac{k_{12} T_2}{k_{44}+L} - \frac{k_{43} T_3}{k_{44}+L} + \frac{L T_4^*}{k_{44}+L}$$

If $L$ is very large,

$$\frac{-k_{41}}{k_{44}+L} \rightarrow 0, \quad \frac{-k_{42}}{k_{44}+L} \rightarrow 0, \quad \frac{-k_{43}}{k_{44}+L} \rightarrow 0, \quad \&$$

$$\frac{L}{k_{44}+L} \rightarrow 1.$$

Then we numerically get,

$$T_4 = T_4^*.$$

Lecture 14: Finite Difference Method (FDM) of discretization.

(H/w)

(Prob-1) Consider the DE: $\dfrac{d^2 v}{dx^2} + v + x = 0$.

with BC $v(0) = v(1) = 0$.

Solve the above eqⁿ using
(1) Least square
(2) Point collocation.
(3) Galerkin
(4) Rayleigh Ritz method.

Choose trial function $\longrightarrow v = a \sin \pi x$

(Prob-2) Consider a heat conduction problem with the following governing DE:

$$\frac{d}{dx}\left( A k \frac{dT}{dx} \right) + Q = 0, \qquad A = 10 m^2, \ k = 5 J/kms,$$

$$Q = 100 \ J/sm.$$

Domain $2 cm \leq x \leq 8 cm.$

BCs: $T(x = 2 cm) = 0°C$

$q''(x = 8 cm) = 15 \ J/m^2 s$.

Obtain temperature distribution in the domain using FEM with Three linear elements, and compare with the analytical solution.

---

• If solving a structural mechanics problem, 'M' form is essentially a statement of minimization of potential energy of that system; which governs the stability of the system at equilibrium.

• In FDM, we deal with the 'D'-form directly.
   Express derivatives in terms of suitable algebraic differences by using Taylor series expansion.

Consider a 1D domain.
we represent the domain with a collection of discrete grid points.
(don't have the concept of discrete elements here).

$$f(x+b) = f(x) + b f'(x) + \frac{b^2}{2!} f''(x) + \cdots \quad \text{—①}$$

$$^{\text{and}} \quad f(x-b) = f(x) - b f'(x) + \frac{b^2}{2!} f''(x) - \cdots \quad \text{—②}$$

We are interested in an algebraic expression for $f'(x)$.

From ①

Thus, $f'(x) = \dfrac{f(x+b) - f(x)}{b} \mp \dfrac{b}{2!} f''(x) \cdots$

here continuous derivative is represented as discrete algebraic quantities.

Error incurred $-\dfrac{b}{2!} f''(x) \cdots \longrightarrow$ truncation error. (TE)

$\sim O(b)$ dictated by the leading order term.

From ② ;

$$f'(x) = \frac{f(x) - f(x-b)}{b} + \frac{b}{2!} f''(x) + \cdots \longrightarrow \text{truncation error (TE)}$$

$$\sim O(b)$$

$-② + ①$ gives

$$f'(x) = \frac{f(x+b) - f(x-b)}{2b} \mp \frac{b^2}{3!} f'''(x) + \cdots \longrightarrow TE \sim O(b^2)$$

$$f'(x) = \frac{f(x+b) - f(x)}{b} + O(b) \longrightarrow \text{Forward difference}$$

$$f'(x) = \frac{f(x) - f(x-b)}{b} + O(b) \longrightarrow \text{Backward difference}$$

$$f'(x) = \frac{f(x+b) - f(x-b)}{2b} + O(b^2) \longrightarrow \text{Central difference}$$

① $+②$ gives

$$f(x+b) + f(x-b) = 2f(x) + \frac{2 \cdot b^2}{2!} f''(x) + \frac{2b^4}{4!} f''''(x) + \cdots$$

$$f''(x) = \frac{f(x+b) + f(x-b) - 2f(x)}{b^2} - \frac{b^2}{12} f''''(x) \cdots \longrightarrow \begin{array}{l}\text{Central difference} \\ \text{for 2nd order} \\ \quad PDE.\end{array}$$

$$f''(x) = \frac{f'(x+h) - f'(x)}{h} \qquad \text{forward - difference}$$

$$= \frac{\dfrac{f(x+h) - f(x)}{h} \nearrow^{BD} - \dfrac{f(x) - f(x-h)}{h} \to^{BD}}{h}$$

$$= \frac{f(x+h) + f(x-h) - 2f(x)}{h^2}$$

**Ex** Consider 1D, steady-state heat conduction problems.

$$\frac{d}{dx}\left(k \frac{dT}{dx}\right) + S = 0$$

$q_i'' = 0$

$x=0$ _____ $x=L$ $(T_L \text{ given})$

Rod $\nearrow$

Assumption: $k, S$ both constants.

$$k\frac{T_{i+1} + T_{i-1} - 2T_i}{h^2} + S = 0$$

$$T_{i+1} + T_{i-1} - 2T_i + \frac{Sh^2}{k} = 0. \qquad \hookrightarrow \text{algebraic eq}^n.$$

Consider 4 grid pts:



1   2   3   4'

The above eq$^n$ is valid for internal grid pts not at the boundary. (B/c we don't have $T_{i-1}$ @ left boundary & $T_{i+1}$ at right boundary).

[ What determines the whether we should choose large $h$ or small $h$?

It depends upon the temperature gradients in the domain. It may so happen that at some regions, the gradients may be steep. In such regions, we use finer value of $h$. At other regions, where the gradient is less, we may opt for a larger $h$. It all depends on the physics of the problem & our understanding about it ]

Grid 2 →

$$T_3 - T_1 - 2T_2 + \frac{Sh^2}{k} = 0$$

Grid 3 →

$$T_4 - T_2 - 2T_3 + \frac{Sh^2}{k} = 0$$

Grid 4 →

BC: $T_4 = T_L$ (given)

Grid 1 →

BC: $q'' = 0$

$$\Rightarrow k\frac{dT}{dx}\Big|_1 = 0$$

So $\dfrac{T_2 - T_1}{h} = 0$   FD formula

$$\Rightarrow T_1 = T_2$$

(Assign the value at the boundary with the interior value)

### Lecture 15: Well Posed Boundary Value Problem.

* Well posed BVP problems requirements:

   ⟶ Existence of solution

   ⟶ Uniqueness of sol$^n$.

   ⟶ A small perturbation in BC shouldn't lead to large changes in
                                               The solution.

       (This is important b/c such a perturbation may be unwillingly.
       introduced Through round-off errors; B/c of it, it may
       lead to large change in sol$^n$ → oversensitive BC).

* Possible types of BCs: (2$^{nd}$ order problems)

  Dirichlet BC:
1. Value of the dependent variable is specified → EBC.

2. Neumann BC :

   Value of the gradient of the dependent variable is specified
                                                 − NBC.

3. Mixed BC:

   value of the dependent variable is expressed as a function of
   The grad.

    Eg: convective heat transfer BC.



$$-k\frac{dT}{dx}\bigg|_{x=L} = h\,(T_L - T_\infty)$$

                         The expression can be used even in
                         unsteady case. (replace $\frac{dT}{dx}$ with $\frac{\partial T}{\partial x}$)

               This BC says, Whatever is the heat flux coming at
                         via conduction
               $x=L$, The same is the heat flux
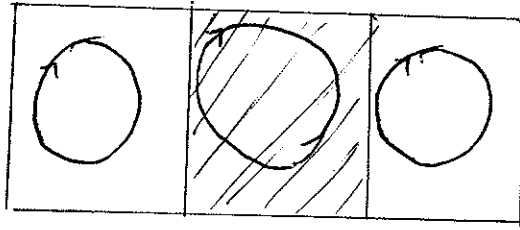               leaving $x=L$ via convection. That remains
                   true even if unsteady

             Here temperature at $x=L$ is expressed as a
             function of gradient of the temperature $\left(-k\frac{dT}{dx}\right)$.
             Hence its a mixed BC.

## 4. Periodic BC:



Say, in the domain. The physics of the problem is such that the solution is periodically repeated.

So, solve 1 part of the domain & extrapolate or extend the solution to other parts of the domain.



$$U_L = U_{\Delta x}$$
$$U_o = U_{L-\Delta x}$$

periodicity = $L - \Delta x$.

Here again we see boundary terms being represented in terms of interior terms, not the opposite.

---

Q) Is any condition specified at the boundary, a boundary condition?

Ans) Consider a simple 1D, steady-state, heat conduction problem, $S = 0$, $k = $ const.

$q_o'' = 1 w/m^2$

$q_L'' = 1 W/m^2$

$x = 0$     $x = L$.

→ Neumann BC at both the boundaries.

$$\frac{d^2 T}{dx^2} = 0$$

$$\frac{dT}{dx} = C_1$$

$$T = C_1 x + C_2 .$$

Say, $k = 1 W/mK$.

$x = 0$,   $-k\frac{dT}{dx} = 1$

$\Rightarrow \frac{dT}{dx} = -1 \Rightarrow C_1 = -1$

$111^{rly}$ at $x = L$,

$-k\frac{dT}{dx} = 1 \Rightarrow \frac{dT}{dx} = -1 \Rightarrow C_1 = -1$

Basically
$$T = -x + C_2 .$$

Cannot determine $C_2$.

⇒ Violates requirement for uniqueness of sol?

Plotting in $T-x$ plane gives all sol^{ns} to be parallel straight lines; no 1 sol?.

→ not legitimate BCs.

⇒ not well posed problem.

BV

**Lecture 16: Finite Volume Method (FVM) of Discretization.**

FDM → simple

Issues: ① Complex geometry.

One has to
tediously ~~create~~
~~#~~ handle the
boundary while
using cartesian grid.

Taylor series expansions:

$$f(x+h) = f(x) + h f'(x) + \boxed{\frac{h^2}{2}} f''(x) + \ldots$$

truncate from here

② But what if $f''(x)$ is large & enough to be non-negligible?
(B/c h cannot be tending to zero, it is still finite).

Ex: $f(x) = e^x \rightarrow f''(x)$ = non-negligible.

⇒ significant errors while truncating.
→ limitation of Taylor series based method.

Q) What do we expect from the discretization?

(ans) → ① Conservativeness → Discretized versions of conservation eqⁿs should exhibit that conservative nature.

→ ② Boundedness

→ ③ Transportiveness

Finite Difference

- ~~FD~~ discretization may not satisfy conservativeness b/c we haven't explicitly enforced that condition while expanding out the function in Taylor series form and ~~not~~ truncation. Conservativeness is not inbuilt while coming up with FDM.

· Boundedness : Say we have a rod.



$0°C$        $100°C$

We are interested in the temperature distribution in the rod.

We expect the values in the rod to lie b/w $0$ & $100°C$. The discretization should also ensure the physical nature of boundedness in the problem.

This boundedness is also not ensured while coming up with FDM.

· Transportiveness : If there is a predominant directionality to the flow involved in the problem, the ~~transport~~ transport properties should also have a predominant transport direction based on the flow direction.

(For high Re flows, for eg, enthalpy should predominantly be transported downstream).

_____

FEM — relatively more complicated, when compared to FDM.

   — Strong mathematical basis in error minimization.

   — not intuitive physically. all the V-formulation & M-formulation needs to have some physical meaning, which can be difficult to come up with.

   For fluid flow, mass flow, Conservation is important, while for Structural mechanics, minimization of potential energy is important.

   — Can handle complex geometries.

# Finite Volume Method (FVM):

Step-1: Divide the domain into a number of sub-domains (Control Volumes)
Each sub-domain is represented by a finite no. of gridpts. — finite size

Step-2: Integrate the governing differential equation over each subdomain

Step-3: Consider a profile assumption for the dependent variable for evaluating the above integrals & to express the result in terms of algebraic quantities at the grid points.
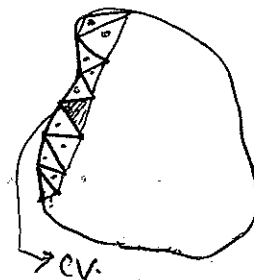
---

## Lecture 17: Illustrative Examples of Finite Volume Method.

Ex Steady state convection diffusion with $S=0$.

General transport equation:

$$\frac{\partial}{\partial t}(\rho\phi)^0 + \nabla\cdot(\rho\vec{V}\phi) = \nabla\cdot(\Gamma\nabla\phi) + \cancel{S}^{0}$$

$$\int_{cv} \nabla\cdot\underbrace{(\rho\vec{V}\phi - \Gamma\nabla\phi)}_{g.d.e.}\,d\forall = 0$$

$\overset{\vec{J}}{\overbrace{(\rho\vec{V}\phi - \Gamma\nabla\phi)}}$

$\rho\vec{V}\phi \to$ advection flux

$\Gamma\nabla\phi \to$ diffusion flux.

Using divergence theorem,

$$\int_{c.s} \vec{J}\cdot\hat{n}\,ds = 0.$$

Grid pts located at the center of each CV.

Q) Why the name 'finite' volume method?

Ans) While deriving the transport equations we considered an infinitesimally small control element. Now we're integrating back to apply to a finite volume. Hence the name finite volume method

Key step → step-2.

Similarity:

It can be thought of having $w=1$ in

$$\int_{CV} \nabla \cdot (\rho \vec{\nabla} \phi - \Gamma \nabla \phi) \, w \, dv = 1 \quad .$$
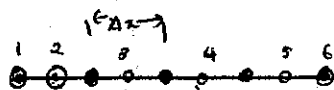
$\underset{r}{\downarrow}$

Galerkin method with $w=1$.

• Requirement of conservation when applied to each & control-volume will satisfy the conservation of across the whole domain .

∴ unlike FDM, FVM takes into account conservation requirement implicitly .

• Profile assumption method is used just for step-3. Afterwards in post analysis, we no longer require profile assumption. Here we have more flexibility in choosing profile assumption as compared to FEM.

Illustration:   1D steady state heat conduction, with $S$ const.

$$\frac{d}{dx}\left( k \frac{dT}{dx} \right) + S = 0$$

Divide 1-D domain into CVs.



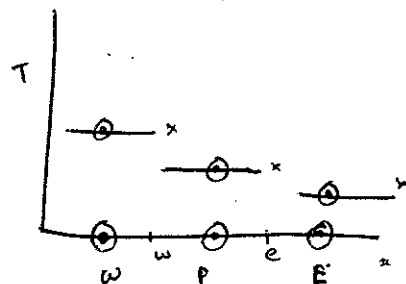In addition to centroids of CVs, we also consider grid pts at the boundary (just so we can impose BCs at those pts).



(adjacent grid pts)
West & East

$$\int_{\omega}^{e} \frac{d}{dx}\left( k \frac{dT}{dx} \right) dx + \int_{\omega}^{e} S \, dx = 0 .$$

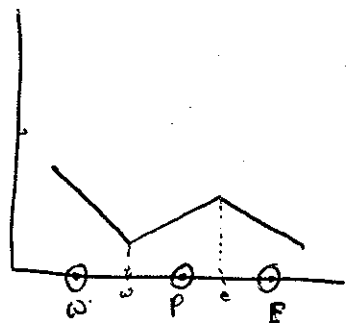$$\left. k \frac{dT}{dx} \right]_{e} - \left. k \frac{dT}{dx} \right]_{\omega} + S \, \Delta x = 0 .$$

• Choosing a profile assumption:



Can consider piecewise continuous functions for each CV.

But not here. B/c we need $\frac{dT}{dx}$ here. Since $T = const$ for each CV, This is not a valid profile assumption. (Here discontinuity is not a problem).
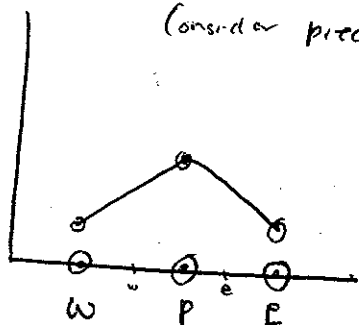
Piecewise linear non-const profile:



Problems with this case:

$k \frac{dT}{dx}$ is not continuous; i.e. physically it means heat flux is not continuous, which is incorrect.

⇒ not an acceptable profile.

Consider piecewise linear profile b/w the grid pts



This will work b/c $\frac{dT}{dx}$ is not evaluated @ P, but rather at the faces of each control volume w, e.

⇒ valid profile assumption

Profile assumption: Piecewise linear T b/w grid pts.

$$\Rightarrow \quad k\frac{dT}{dx}\Big]_e - k\frac{dT}{dx}\Big]_w + S\Delta x = 0$$

$$\rightarrow \quad k_e \frac{T_E - T_P}{\delta x_e} - k_w \frac{T_P - T_w}{\delta x_w} + S\Delta x = 0$$

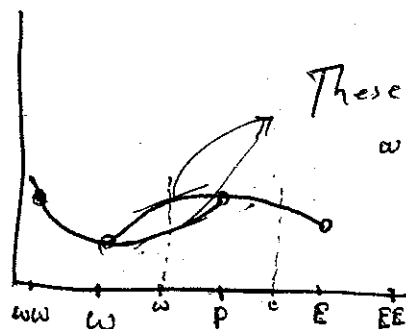$$\Rightarrow -\left(\frac{k_e}{\delta x_e} + \frac{k_w}{\delta x_w}\right) T_P \cancel{-\left(\frac{k_e}{\delta x_e} + \frac{k_w}{\delta x_w}\right)P} + \frac{k_e}{\delta x_e} T_E + \frac{k_w}{\delta x_w} T_w + S\Delta x = 0.$$

$$a_P T_P = a_E T_E + a_w T_w + b.$$
$$\;(i) \qquad\quad (i+1) \qquad (i-1)$$

where $a_E = \dfrac{k_e}{\delta x_e} \qquad a_w = \dfrac{k_w}{\delta x_w}$ , $a_P = a_E + a_w$ , $b = S\Delta x$

$\frac{k}{\delta x} \longrightarrow$ conductance (physical meaning)

Would it have been better to consider a higher order interpolation function?



These slopes need not necessarily be equal. It happens only when $\omega$ lies exactly in b/w $W$ and $P$. From Rolle's thm, a chord b/w two pts has a slope that will be attained by the curve b/w those two pts at some pt b/w them, given the curve is differentiable at all points. From MVT, this will lie at the midpoint of the curve. ~~It is only happens when all the~~

So higher order interpolate functions need not necessarily given more accurate results in FV M* (counter-intuitive).

---

## Lecture -18 : Illustrative Examples of Finite Volume Method. (Contd)

1-D steady state heat conduction equation.
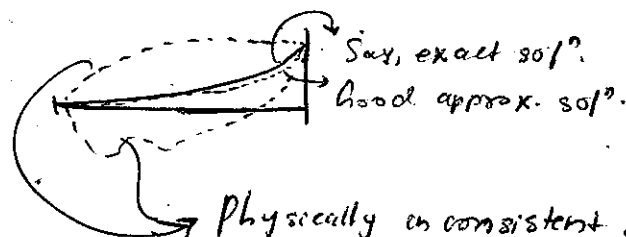
$$\frac{d}{dx}\left(k\frac{dT}{dx}\right) + S = 0$$

If $k, S \longrightarrow fn(T) \longrightarrow$ non linear eq" $\longrightarrow$ solved via iterative process
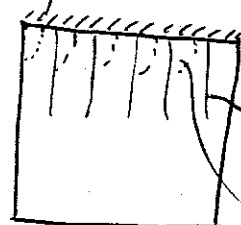
- Requirements of discretization:
    (1) Physical consistency
    (2) Overall balance.

Eg:



$T = 0°c$      $T = 100°C$      Assume $S = 0$
                                 $k = const.$

Say, exact sol".

Good approx. sol".

insulated $\longrightarrow$ Physically inconsistent solutions. (Boundedness invalidated)
                                                              for eg,.

Iso Thermal lines in 2D domain

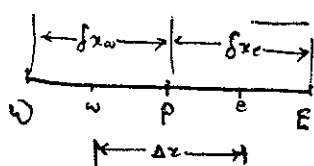$\longrightarrow$ These isotherms are physically consistent (normal to insulated surface)

$\longrightarrow$ These isotherms are physically inconsistent.

Variable $S$ : $\longrightarrow$ let $S(T)$

Ex: $\underline{S \text{ is a linear fn of } T.}$
i.e. $S = a + bT.$



$$k \frac{dT}{dx}\Big]_e - k\frac{dT}{dx}\Big]_w + \int_w^e (a+bT)\, dx = 0.$$

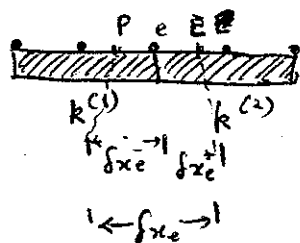$\hookrightarrow$ piecewise const $T$ profile within each CV.

Profile assumpts: piece wise linear $T$ b/w grid pts.

$\Rightarrow \quad k_e \dfrac{(T_E - T_P)}{\delta x_e} - k_w \dfrac{(T_P - T_w)}{\delta x_w} + (a + b\,T_P)\Delta x = 0.$

$a_P T_P = a_E T_E + a_w T_w + b.$

where $a_E = \dfrac{k_e}{\delta x_e}$ , $a_w = \dfrac{k_w}{\delta x_w}$ , $a_E + a_w - b\Delta x = a_P$ , $b = a\Delta x.$

$\underline{C}$omposite material with position dependent $k$:
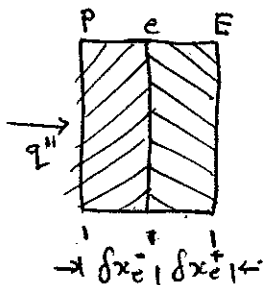


$\rightarrow 4$ CVs

$\rightarrow$ Requirement: Thermal conductivity at the interface.

Since interface shared by both materials, an equivalent thermal conductivity needs to be described.

If $k_P$ & $k_E$ is known, intuitive: $k_e = \dfrac{k_E + k_P}{2}$ $\longrightarrow$ linear interpolation

Physical assessment of this $k_e$ needs to be performed: $\rightarrow$ AN formulation.



$q''_{left} = \left(\dfrac{T_e - T_P}{\delta x_e^-}\right)(k_P)$

$q''_{left} = q''_{right}$

$= \dfrac{T_e - T_P}{\left(\dfrac{-\delta x_e^-}{k_P}\right)}$

$q''_{right} = \dfrac{T_E - T_e}{\delta x_e^+}(-k_P)$

$= \dfrac{T_E - T_e}{\left(\dfrac{-\delta x_e^+}{k_P}\right)}$

$\dfrac{T_P - T_e}{\dfrac{\delta x_e^-}{k_P}} = \dfrac{T_e - T_E}{\dfrac{\delta x_e^+}{k_E}} = \dfrac{T_P - T_E}{\dfrac{\delta x_e}{k_e}} = \dfrac{T_P - T_E}{\dfrac{\delta x_e^-}{k_P} + \dfrac{\delta x_e^+}{k_E}}$

$\hookrightarrow$ where $k_e$ is the equivalent thermal conductivity.

If $\int x_e^- = \int x_e^+$,

$$\frac{2}{k_e} = \frac{1}{k_p} + \frac{1}{k_E}$$

$\longrightarrow$ H.M. formulation

or $k_e = \dfrac{2 k_p k_E}{k_p + k_E} = \dfrac{2}{\dfrac{1}{k_p} + \dfrac{1}{k_E}}$

• $\underline{L}$imiting cases:

$$K_p \gg k_E$$

Then  AM $\longrightarrow$ $k_e \approx \dfrac{k_p}{2}$

HM $\longrightarrow$ $k_e \approx 2 k_E$

For interfacial conductivity variation, HM formulation is physically much more appealing. why?

Say $k_E = 0$. That is equivalent to highly insulated material-E. In that case $k_e$ should also be $0$ at the interface. This is reflected in HM, while not at all in AM formulation.

---

**Lecture 19 · Basic rules of finite Volume Discretization.**

4 Basic rules (of 1D steady state diffusion type problem):

(1) Physical consistency of fluxes at Control Volume faces.

→ Profile should be chosen in such a way that there is no discontinuity of flux at control volume faces.

(2) All coefficients in the discretized equation must be of the same sign.

Eg; say, $b = 0$

$$10 \, T_p = 15 \, T_E - 5 \, T_w.$$

Say $T_E = 10$

$T_w = 100$

$\therefore T_p = \dfrac{15 \times 10 - 5 \times 100}{10}$

$= -35$

$T_p$ is not bounded b/w $T_E$ & $T_w$.

⇒ Physically inconsistent.

This inconsistency has originated b/c of the -ve sign in the discretization eqn ($10 \, T_p = 15 \, T_E - 5 \, T_w$)

4, $10 T_p = 5 T_E + 5 T_w$

$T_p = 5 \times 10 + 5 \times 100 \over 10$

$= 55$

Here $T_E \leq T_p \leq T_w$.

Hence consistent.

(Same sign coefficients)

By sign convention, we will consider that sign to be +ve.

(3) If the source term is linearized as: $S = S_c + S_p T_p$,

Then $S_p$ must be $-ve \longrightarrow$ Extension of consideration-#2.

(4) If a linear governing DE is discretized, its discretized version should satisfy the following requirement:

If $T$ is a sol$^n$, Then $T+c$ is also a sol$^n$.

we are interested to see

$$a_p \cdot (T_p + c) \overset{?}{=} a_E (T_E + c) + a_w (T_w + c) \cdots ①$$

We already know

$$a_p T_p = a_E T_E + a_w T_w \cdots ②$$

① − ②

$$\Rightarrow \quad a_p \, c = (a_E + a_w) \, c$$

or $a_p = a_E + a_w$

This linearity is satisfied.

Q) What if the source term is non-linear?

Method to linearize a non-linear source term $\longrightarrow$

Source term linearization

Ex-1  Say, $S = 3 + 4T$       $S = S_c + S_p T_p$

$S_c = ?$   $S_p = ?$

from the form, it may appear that $S_c = 3$, $S_p = 4$.

Just now we've seen that $S_p$ should be taken $-ve$. So the above form is not valid.

Use iterative process. initially take:

$$S_c = 3 + 4 T_p^*, \quad S_p = 0.$$

If the iteration has a tendency to diverge fast, considering appropriate initial value for $S_p$ may be required.
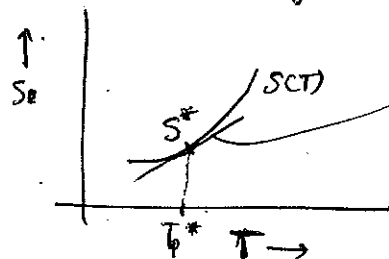
__Ex-2__   $S = 3 - 4T$

Here $S_c = 3, \quad S_p = -4$ is valid & the best choice.

To slow down the convergence, we may choose $S_c = 3 + 4 T_p^*$
$$S_p = -8.$$

__Ex-3__   $S = 3 - 4 T^3.$

Say we have any arbitrary $S(T)$.



Best linear function for this $\to$ tangent at $T_p^*$.
(why? Same first order derivative for tangent as that of the non-linear curve).

$$\frac{S - S^*}{T - T_p^*} = \frac{dS}{dT}\Big|_{T_p^*}$$

$$S - S^* = \frac{dS}{dT}\Big|_{T_p^*} (T - T_p^*)$$

$$S^* = 3 - 4 T_p^{*\,3}$$

$$\frac{dS}{dT}\Big|_{T_p^*} = -12 T_p^{*\,2}$$

$$S - S^* = S - (3 - 4 T_p^{*\,3})$$

$$= -12 T_p^{*\,2} (T - T_p^*)$$

$$\to \quad S = 3 + 8 T_p^{*\,3} - 12 T_p^{*\,2} T.$$

$$\therefore S_c = (3 + 8 T_p^{*\,3})$$

$$S_p = (-12 T_p^{*\,2}) \longrightarrow S_p \text{ -ve satisfied as } T_p^{*\,2} \geqslant 0$$

Ex-4

$S = 3 + 4T^3$.

$S^* = 3 + 4T_p^{*3}$ $\left\{ \frac{dS}{dT}\Big|_{T_p^*} = \right.$ $\wedge 12 T_p^{*2} (T - T_p^*)$

$S - [3 + 4T_p^{*3}] = 12 T_p^{*2} (T - T_p^*)$

$S = 8 - 8 T_p^{*3} + 12 T_p^{*2} T$.

$\Rightarrow S_c = 3 - 8 T_p^{*3}$

$S_p = 12 T_p^{*2}$ $\longrightarrow$ Here +ve. $S_o$ cannot work!
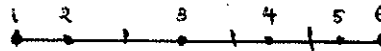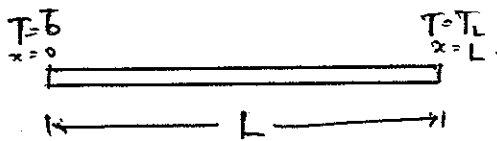
What to do then?

Dump the entirety to $S_c$.

$S_c = 3 + 4 T_p^{*3}$, $\quad S_p = 0$

This linearization is mathematically correct, but will give physically inconsistent solution

---

| Lecture 20: Implementation of boundary conditions in FVM. |
| --- |

Ex-1

$T = T_b$
$x = 0$

$T = T_L$
$x = L$

$\longleftarrow \quad L \quad \longrightarrow$|

1   2   3   4   5  6

$a_p T_p = a_E T_E + a_w T_w + b$

For grid point-2,

$a_2 T_2 = a_3 T_3 + a_1 T_1 + b$

But at point-1, $T_1 = T_o$ (given)

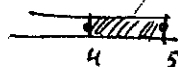(Use penalty approach in FEM for further analysis)

Ex-2.   $T = T_b$

$\longrightarrow 2'' = 2_L''$ (given)

$\longleftarrow \quad L \quad \longrightarrow$

To implement BC, consider CV at the boundary.

Take $\frac{1}{2}$ of a CV. (Smaller length can capture sharper gradients accurately)

G.de: $\frac{d}{dx}\left(k \frac{dT}{dx}\right) + S = 0$.

Integrating gde over the half CV →

$$\int_4^5 \frac{d}{dx}\left(k\frac{dT}{dx}\right)dx + \int_4^5 s\,dx = 0$$

Total length of $\frac{1}{2}$ CV = $\frac{\Delta x}{2}$.

$$\Rightarrow \quad \underbrace{k\frac{dT}{dx}\Big]_5}_{\downarrow} - \underbrace{k\frac{dT}{dx}\Big]_4}_{} + S\frac{\Delta x}{2} = 0.$$

$$-q_L''$$

Use profile assumption of
piecewise linear noñconstant T.

$$\Rightarrow \quad k\frac{T_5 - T_4}{\left(\frac{\Delta x}{2}\right)}$$

$$\Rightarrow \quad -q_L'' - \frac{2k}{\Delta x}(T_5 - T_4) + S\frac{\Delta x}{2} = 0.$$

$$\therefore \quad \frac{2k}{\Delta x}(T_5 - T_4) = -q_L'' + \frac{S\Delta x}{2}$$

$$T_5 = T_4 + \frac{S\Delta x^2}{4k} - \frac{q_L''\Delta x}{2k} \longrightarrow$$ Eq$^n$ of the form:

$$a_5 T_5 = \cancel{a_6 T_6} + a_4 T_4 + b$$

$\hookrightarrow$ ( Expression written as
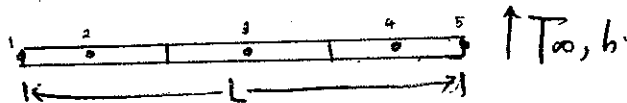boundary term as a function of
interior terms)

where $a_5 = 1$
$a_4 = 1$,
$b = \left(\frac{S\Delta x^2}{4k} - \frac{q_L''\Delta x}{2k}\right)$.

FD

$$q_L'' = -k\left(\frac{T_5 - T_4}{\frac{\Delta x}{2}}\right)$$

$$T_5 = T_4 - \frac{q_L''\cancel{2}\Delta x}{2k}.$$

Ex-3  (Mixed type BC)

 $\uparrow T_\infty, h$

$$-k\frac{dT}{dx}\Big|_{x=L} = h\left(T_{x=L} - T_\infty\right)$$

$$\underbrace{\phantom{xxxx}}_{q_L''}$$

$$-q_L'' - \frac{2k}{\Delta x}(T_5 - T_4) + \frac{S\Delta x}{2} = 0$$

$$\underbrace{\phantom{x}}_{}$$

$$h(T_5 - T_\infty) \quad \leadsto \quad () T_5 = () T_4 + () \quad \longrightarrow \text{very much analogous to}$$
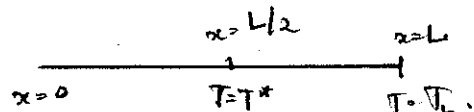$$\text{Neumann BC case.}$$

I → Heat flux enters through left boundary.

II → Source at 2. Requirement: equivalent heat flux passing through face b/w 2 & 3 same as in I. Boundary-1 insulated.

$$S_{extra} = \frac{(q_o'' A_{face})}{V_{cv}}$$ → total rate of heat transfer

Volumetric heat generation (Heat source terms)

Case-I is usually treated equivalently as Case-II by dumping all the flux through corresponding heat source term $(S)_{extra}$. Already present source terms un affected. Why do so?

More rapid convergence to solution.

In case I, heat flux has to penetrate through 1 & then through face 2-3. While in case II, with the introduction of $S_{extra}$, flux through 2-3 face is handled in 1 shot.

— Difference observed is marginal in most cases.

Can we specify the B/C at an internal Grid pt & still expect the problem to be well-posed?

$$\frac{d}{dx}\left(k\frac{dT}{dx}\right) = 0 \quad (\text{Take } S=0 \text{ for this eg}).$$
$$k = const.$$

$$\Rightarrow \frac{dT}{dx} = c_1$$

$$T = c_1 x + c_2$$

$x = L/2, \ T = T^* \longrightarrow$
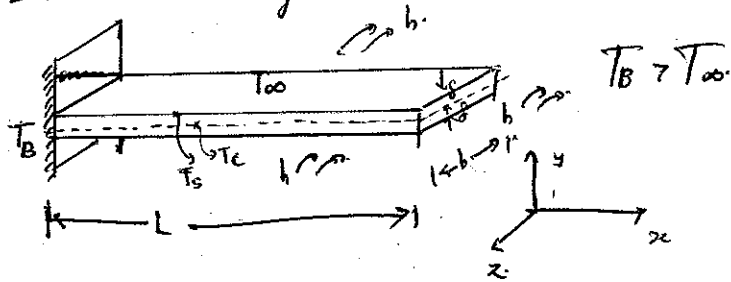$x = L, \ \ T = T_L \longrightarrow$ } Then can obtain both $c_1$ & $c_2$.

↳ perfectly valid BC. Only thing is that it isnt specified at the physical boundary.

Cannot do the same for IVP. Why?

B/C time is a one way coordinate system. B/c whatever is happening now cannot influence what has already happened a while back.

Ex. 1-D steady state heat conduction in a fin.



$T_B > T_\infty$

Thermal resistance in $x$ is most important. Out of $y$ & $z$, since $b >> 2\delta$, thermal resistance along $z$ direction is the next significant one.
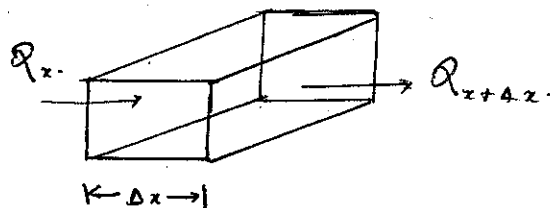
$$\frac{\delta}{kA_1}, \quad \frac{b}{kA_2}$$

Either the consideration $b >> 2\delta$, or $2\delta >> b$ makes this problem 2D.

$$k\frac{T_c - T_s}{\delta} \sim h(T_s - T_\infty)$$

$$\frac{T_c - T_s}{T_s - T_\infty} \sim \boxed{\frac{h\delta}{k}} \longrightarrow Bi_\delta \rightarrow \text{assume this to be small.}$$

If $\dfrac{T_c - T_s}{T_s - T_\infty} << 1 \Rightarrow T_c \sim T_s \Rightarrow$ no necessity to analyse the temp difference b/w centre line & outside surface

$\Rightarrow$ no necessity to analyse heat transfer in $y$ direction

$\Rightarrow$ 1-D problem.

Take a section out of the fin:



$$\overset{0}{Q_x} - A\,Q_{conv} = \overset{0}{Q_{x+\Delta x}}$$

$$\overset{0}{Q_x} + \overset{0}{Q_{x+\Delta x}} \quad Q_x^0 = q_x'' A$$

$$Q_{x+\Delta x}^0 = q_{x+\Delta x}'' A$$

$$Q_x^0 - Q_{x+\Delta x}^0 = (q_x'' - q_{x+\Delta x}'')A$$

$$q_{x+\Delta x}'' = q_x'' + \Delta x \frac{dq_x''}{dx} + \cdots$$

$$\Rightarrow Q_x^0 - Q_{x+\Delta x}^0 = -\left[\frac{dq_x''}{dx}\Delta x + \cdots\right]A.$$

* $T_{op} = 2hb\,\Delta x\,(T-T_\infty)$

$Front = 2h(2\delta)\,\Delta x\,(T-T_\infty)$

$\Delta Q_{conv} = Ph\,(T-T_\infty)\,\Delta x$

$\qquad P = \text{perimeter} = 2(b+2\delta)$.

$Q'_x - Q'_{x+\Delta x} + \Delta Q_{conv}$

$\therefore \quad -\dfrac{dq''_{Lx}}{dx}\,\Delta x + \cdots = \dfrac{Ph\cdot(T-T_\infty)\,\Delta x}{A}$

Take limit as $\Delta x \to 0$.

$\Rightarrow -\dfrac{dq''_{Lx}}{dx}A = Ph\,(T-T_\infty)$

and $q''_{Lx} = -k\dfrac{dT}{dx}$

$+\dfrac{d}{dx}\left(kA\dfrac{dT}{dx}\right) = Ph\,(T-T_\infty)$

$\dfrac{d}{dx}\left(kA\dfrac{dT}{dx}\right) - Ph\,(T-T_\infty) = 0.$

$\quad\hookrightarrow$ Governing DE.

· Discretize using FVM:

The above eq$^n$ is of the form

$\dfrac{d}{dx}\left(k\dfrac{dT}{dx}\right) + S = 0$

with $S = -\dfrac{Ph}{A}(T-T_\infty)$

$\qquad\uparrow$

implicitly linear source terms.

$S = S_c + S_p T$

$S_c = \dfrac{Ph\,T_\infty}{A}$, $S_p = -\dfrac{Ph}{A}$

$\qquad\qquad (-ve)$.

$\Rightarrow$ It's a well posed problem.

(H/W) (1) At $x=0$, $T=T_B$

(2) At $x=L$, $\dfrac{dT}{dx}=0$.

Non-dim entire equations: $\bar{x}=\dfrac{x}{L}$.
Find Temp distribution with FDM, FEM & FVM.

<u>Ex</u> 1D steady state heat conduction in cylindrical coordinates.

Consider Z large & no variation along z.
Axially symmetric ⇒ no variation along θ.

Governing DE:

$\dfrac{1}{r}\dfrac{d}{dr}\left(rk\dfrac{dT}{dr}\right) + S = 0.$

$dA = r\,d\theta\,dr$
$dV = r\,d\theta\,dr\,dz$

$dV = (r\,dr)(d\theta)(dz)$,
form integral eq$^n$ using this in
$\int\dfrac{1}{r}\dfrac{d}{dr}\left(rk\dfrac{dT}{dr}\right)r\,dr + \int Sr\,dr = 0$
$\quad\hookrightarrow$ Multiplying with $r\,dr$
remove singularity induced by $\dfrac{1}{r}$.
$\quad\hookrightarrow$ length of the CV.

$rk\dfrac{dT}{dr}\Big]_{r+\Delta r} - rk\dfrac{dT}{dr}\Big]_r + \int_w^e(S_c+S_pT)r\,dr = 0$

(Assume, $S = S_c + S_pT$
Constant temp over CV)

Make precewise linear profile ass assumptns b/w the grid pts.



$$r_e k_e \frac{T_E - T_P}{\delta r_e} - r_w k_w \frac{T_P - T_w}{\delta r_w}$$

$$+ \left( S_c + S_p T_P \right) \frac{r_e^2 - r_w^2}{2} = 0$$

Organise the eqⁿ in the form:

$$a_P T_P = a_E T_E + a_w T_w + b.$$

where

$$a_P = \frac{r_e k_e}{\delta r_e} + \frac{r_w k_w}{\delta r_w}$$

$$a_E = \frac{r_e k_e}{\delta r_e} \qquad a_w = \frac{r_w k_w}{\delta r_w}$$

$$a_P = a_E + a_w - S_p \left( \frac{r_e^2 - r_w^2}{2} \right)$$

$$b = S_c \left( \frac{r_e^2 - r_w^2}{2} \right)$$

— If it were spherical coord systems, eqⁿ would have been

$$\frac{1}{r^2} \frac{d}{dr} \left( r^2 \frac{dT}{dr} \right) + S = 0.$$

and while forming the integral eqⁿ, we'd have to multiply with $r^2 dr$ instead of $r dr$ ~~x dx~~ before integrating to remove singularities induced by $\frac{1}{r^2}$ term.

FVM for:
1D unsteady state diffusion problem:

$$\frac{\partial}{\partial t} (\ell C_p T) + \nabla \cdot (\ell C_p T \vec{V}) \overset{0 \text{ (diffusion; so no fluid flow terms)}}{=} \underset{\frac{\partial}{\partial x}(k \frac{\partial T}{\partial x})}{\nabla (k \nabla T)} + S$$

$$\Rightarrow \underbrace{\frac{\partial}{\partial t} (\ell C_p T)}_{\text{unsteady term}} = \underbrace{\frac{\partial}{\partial x} \left( k \frac{\partial T}{\partial x} \right)}_{\text{conduction term}} + S.$$

Our obj: how do take care of the new term:
$$\frac{\partial}{\partial t} (\ell C_p T) ?$$

④ Integrate eqⁿ over the domain.
Here domain consists of both t & x.
So our elemental domain consists of both dt & dx.     { Take S = 0 }

$$\therefore \int_w^e \int_t^{t+\Delta t} \frac{\partial}{\partial t} (\ell C_p T) \, dt \, dx = \int_t^{t+\Delta t} \int_w^e \frac{d}{dx} \left( k \frac{\partial T}{\partial x} \right) \, dx \, dt.$$

Treat both terms separately.



$$\int_w^e \int_{t+\Delta t}^{e} \frac{\partial}{\partial t} (\ell C_p T) \, dt \, dx$$

$$= \int_w^e (\ell C_p T)^{t+\Delta t} - (\ell C_p T)^t \, dx$$

Make profile assumptn for T variatn in space:

Simplest assumptn: piecewise constant profile

(why? Here no need to take piecewise linear forms about grid pts b/c there are no derivative terms involved).

Assume $\ell C_p$ a constant.

$$\therefore \int_{v}^{e} \ell C_p T]^{t+\Delta t} - \ell C_p T]^{t} \, dx$$

$$= \ell C_p \cdot \left(T_p^{t+\Delta t} - T_p^{t}\right) \Delta x$$

Term = 2:

$$\int_{t}^{t+\Delta t}\int_{w}^{e} \frac{\partial}{\partial x}\left(k \frac{\partial T}{\partial x}\right) dx \, dt.$$

$$= \int_{t}^{t+\Delta t}\left(\left[k \frac{dT}{dx}\right]^{e} - k\frac{dT}{dx}\right]^{w}\right) dt$$

Profile assumption for t:

piecewise linear t.

$$\Rightarrow = \int_{t}^{t+\Delta t}\left\{k_e\left(\frac{T_E - T_P}{\delta x_e}\right) - k_w \frac{(T_P-T_w)}{\delta x_w}\right\} dt$$

$$\int_{t}^{t+\Delta t} T dt = ??$$

Make profile assumption:

$$= \left[(1-f)T_p^{t} + f T^{t+\Delta t}\right]\Delta t$$

$$\Rightarrow = \left[\frac{k_e}{\delta x_e}\left\{(1-f)T_E^{t} + f T_E^{t+\Delta t}\right\}\right.$$

$$- \frac{k_e}{\delta x_e}\left\{(1-f)T_p^{t} + f T_p^{t+\Delta t}\right\}$$

$$- \frac{k_w}{\delta x_w}\left\{(1-f)T_p^{t} + f T_p^{t+\Delta t}\right\}$$

$$+ \frac{k_w}{\delta x_e}\left\{(1-f)T_\omega^{t} + f T_\omega^{t+\Delta t}\right\}\right]\Delta t.$$

For notational convenience, write

$$T^{t} = T^{o}.$$
$$T^{t+\Delta t} = T^{1} = T.$$

Term① = Term②.

$$\Rightarrow a_p T_p = a_E T_E + a_w T_\omega + a_p^{o}T_p^{o} + b.$$

new term appears!

where,

$$a_E = \frac{k_e}{\delta x_e}f \quad, \quad a_w = \frac{k_w f}{\delta x_e}.$$

$$a_p^{o} = -\frac{(1-f)}{f}(a_E + a_w) + \ell C_p \frac{\Delta x}{\Delta t}$$

$$b = \frac{k_e}{\delta x_e}(1-f)T_E^{o} + \frac{k_w}{\delta x_w}(1-f)T_\omega^{o}$$

$$a_p = \ell C_p \frac{\Delta x}{\Delta t} + \frac{k_e}{\delta x_e}f + \frac{k_w}{\delta x_w}f$$

$$= a_E + a_w + \ell C_p \frac{\Delta x}{\Delta t}$$

Earlier there were two neighbors E and $\underbrace{\text{(spatial neighbors)}}$
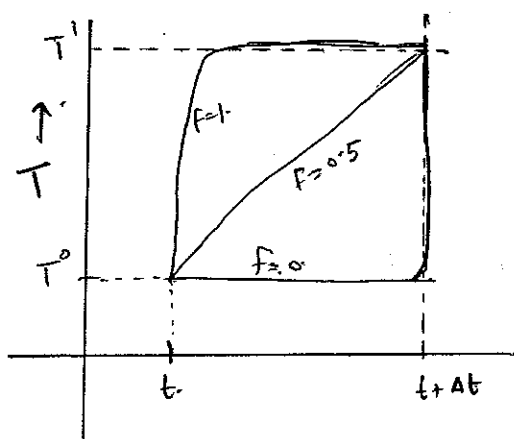cv. Now we have a temporal neighbor $p^{o}$.

We only have one neighbor for time b/c whatever has happened at time $t+\Delta t$ is influenced by the events at time $t$ and not by events at time $t+2\Delta t$.

That's why 2 space neighbors & 1 time neighbor.

## Choices of f



$- \text{Remember:}$
$0 \leq f \leq 1$

$T^t = T^o$

$T^{t+\Delta t} = T^1 = T$

$f = 0 \longrightarrow$ fully explicit scheme

$f = 1 \longrightarrow$ implicit scheme

$f = \frac{1}{2} \longrightarrow$ Crank-Nicholson scheme.

$f = 0:$

$$a_E = 0, \quad a_w = 0, \quad a_p = \rho C_p \frac{\Delta x}{\Delta t},$$

$$a_p^{\,o} = \rho C_p \frac{\Delta x}{\Delta t} - \frac{k_e}{\delta x_e} - \frac{k_w}{\delta x_w}$$

$$b = \frac{k_e}{\delta x_e} T_E^{\,o} + \frac{k_w}{\delta x_w} T_w^{\,o}$$

$$a_p T_p = a_p^{\,o} T_p^{\,o} + b$$

$$T_p = \frac{a_p^{\,o} T_p^{\,o}}{a_p} + \frac{b}{a_p}.$$

$T_p = f(T_E^{\,o}, T_p^{\,o}, T_w^{\,o})$ expressed in an explicit form.

---

$f = 1:$

$$a_E = \frac{k_e}{\delta x_e}, \quad a_w = \frac{k_w}{\delta x_w}, \quad a_p = a_E + a_w + \rho C_p \frac{\Delta x}{\Delta t}$$

$$a_p^{\,o} = \rho C_p \frac{\Delta x}{\Delta t}, \quad b = 0.$$

---

> **Lecture 23: 1-D unsteady state Diffusion.**
> **Problem (contd)**

We should have $a_p^{\,o} \geq 0$.

$\underline{f = 0}.$

$$a_p^{\,o} = \rho C_p \frac{\Delta x}{\Delta t} - \frac{k_e}{\delta x_e} - \frac{k_w}{\delta x_w}.$$

Let $k_e = k_w = k$ (for algebraic simplicity)
$\delta x_e = \delta x_w = \delta x = \Delta x.$



$$a_p^{\,o} = \rho C_p \frac{\Delta x}{\Delta t} - \frac{2k}{\Delta x}.$$

Conditions for $a_p^{\,o} \geq 0$:

$$\longrightarrow \rho C_p \frac{\Delta x}{\Delta t} \geq \frac{2k}{\Delta x}.$$

$$\Rightarrow \alpha \frac{\Delta t}{\Delta x^2} \leq \frac{1}{2} \qquad \alpha - \text{thermal diffusivity.}$$
$$= \frac{k}{\rho C_p}$$

$$\frac{\alpha t_c}{L_c^2} = F_o \longrightarrow \text{Fourier number}$$

$\hookrightarrow$ characteristic time
characteristic length.

Stability criterion for the explicit scheme.

Round off errors can propagate & amplify with calculations. If such a thing happens & it is inherent to the scheme itself, then such a scheme is an unstable scheme.

Key requirement of stable scheme:

→ physically consistent: coeffs are of same sign. If temp is increased at a pt, then that change will cause the temperature at neighbouring pts to increase as well & not decrease.

B/c of the condition $\dfrac{\alpha \Delta t}{\Delta x^2} \le \dfrac{1}{2}$, change in grid spacing also affects the size of the time step to be chosen.

$$\dfrac{\alpha t_c}{L^2}, \dfrac{t_1}{t_c}, \quad t_1 \longrightarrow \Delta t$$

$$t_c \longrightarrow \dfrac{\Delta x^2}{\alpha} \rightarrow \text{characteristic}$$

time over which a Thermal disturbance propagates by thermal diffusion in a medium.

→ conditionally stable schemes: The scheme works as long as $\Delta t, \Delta x$ constraint is followed.

---

$f = 1 \rightarrow$

all coeffs same sign

⇒ unconditionally stable.

$f = 0.5 \longrightarrow$

$$\rho C_p \dfrac{\Delta x}{\Delta t} - \dfrac{2k}{\Delta x} \cdot \dfrac{1}{2} \ge 0$$

$$\Rightarrow \dfrac{\alpha \Delta t}{\Delta x^2} \le 1.$$

---

| Lecture 24: Consequences of Discretization of Unsteady State Problems. |
|---|

Consequences of time-discretization in Finite Difference:

Errors associated with any Taylor series based discretization:

→ Consistency: — characteristic of a numerical scheme.

Consistent if in the limit of grid size & time step size → 0, the algebraic eqⁿs mimic the same behavior as that of its parent diff. eqⁿ.

This happens when error is nullified at such refined scales. ⇒ nullification of truncation error as grid size & time step size finds to... Zero in the limit.

⇒ discretized eqⁿ tends to behave same as g.d.e.

→ Stability: Just like consistency talks about truncation errors, stability talks of round off

errors. $\underset{\wedge}{\text{Errors}}$ Round off. in a numerical scheme

is similar to physical perturbation. How strongly these perturbations propagate / amplify in the presence of numerical calculation determines stability.

Stable ⟹ no amplification of numerical perturbations due to propagation of round-off errors.

→ Convergence: ~~As~~ In the limit of grid size and time step size tends to 0, numerical soln → exact solution.

Lax equivalence theorem: for <u>linear problems</u>,

Consistency + stability ⟹ convergence.

Consistency + stability need not ensure convergence for a non-linear problems. Such problems can have multiple solutions. To test for convergence in non linear problems, the following is done. At a finite grid size & time step size evaluate the soln. Then take finer grid & smaller time step & evaluate soln. As finer & finer grid & time step is used, if the soln is found to be grid-independent & time-step independent, the non-linear problem has convergence.

Finite difference schemes on 1-b unsteady state diffusion problems:

G.de: $\rho C_p \dfrac{\partial T}{\partial t} = k \dfrac{\partial^2 T}{\partial x^2}$

(assume constant properties).

$\dfrac{\partial T}{\partial t} = \alpha \dfrac{\partial^2 T}{\partial x^2} \quad \left(\alpha = \dfrac{k}{\rho C_p}\right).$

Use FTCS.

FT ⟶ Forward Time

CS ⟶ Central Space.

Corresponding Taylor series:

```
   x-Δx    x    x+Δx
  ---|-----|-----|---
   i-1     i    i+1
   (ω)    (P)   (E)
```

Superscript:
n ⟶ current time (t)
n-1 ⟶ previous time (t-Δt)
n+1 ⟶ next time. (t+Δt)

$T_i^{n+1} \overset{\to T_x(t+\Delta t)}{=} T_i^n + \dfrac{\partial T}{\partial t}\bigg]_i \Delta t + \dfrac{\partial^2 T}{\partial t^2}\bigg]_i \dfrac{\Delta t^2}{2} +$

$O(\Delta t^3).$

$\dfrac{\partial T}{\partial t}\bigg]_i = \dfrac{T_i^{n+1} - T_i^n}{\Delta t} - \dfrac{\partial^2 T}{\partial t^2}\bigg]_i \dfrac{\Delta t}{2} + O(\Delta t^2)$

Truncated representation of forward time derivative:

$\dfrac{\partial T}{\partial t}\bigg|_i = \dfrac{T_i^{n+1} - T_i^n}{\Delta t}$

$T_{i+1}^n = T_i^n + \dfrac{\partial T}{\partial x}\bigg|_i^n \Delta x + \dfrac{\partial^2 T}{\partial x^2}\bigg|_i^n \dfrac{\Delta x^2}{2!} +$

$\dfrac{\partial^3 T}{\partial x^3}\bigg|_i^n \dfrac{\Delta x^3}{6} + \dfrac{\partial^4 T}{\partial x^4}\bigg|_i^n \dfrac{\Delta x^4}{24} + O(\Delta x^5)$

$$T_{i-1}^n = T_i^n - \frac{\partial T}{\partial x}\Big[_i \Delta x + \frac{\partial^2 T}{\partial x^2}\Big]_i \frac{\Delta x^2}{2}$$

$$- \frac{\partial^3 T}{\partial x^3}\frac{\Delta x^3}{6} + \frac{\partial^4 T}{\partial x^4}\frac{\Delta x^4}{24} + O(\Delta x^5)$$

$$T_{i+1}^n + T_{i-1}^n = 2 T_i^n + \frac{\partial^2 T}{\partial x^2}\Big|_i \Delta x^2 +$$

$$\frac{\partial^4 T}{\partial x^4}\Big|_i \frac{\Delta x^4}{12} + O(\Delta x^6).$$

$$\frac{T_{i+1}^n + T_{i-1}^n - 2T_i^n}{\Delta x^2} = -\frac{\partial^4 T}{\partial x^4}\frac{\Delta x^2}{12} + O(\Delta x^4)$$

$$= \frac{\partial^2 T}{\partial x_i^2}.$$

Sub into g.d.e.

$$\frac{T_i^{n+1} - T_i^n}{\Delta t} - \overbrace{\left[\frac{\partial^2 T}{\partial t^2}\Big|_i \frac{\Delta t}{2}\right]}^{\text{Term-1}} + O(\Delta t^2)$$

$$= \alpha \frac{T_{i+1}^n + T_{i-1}^n - 2T_i}{\Delta x^2} - \overbrace{\left[\alpha\frac{\partial^4 T}{\partial x^4}\frac{\Delta x^2}{12}\right]}^{\text{Term-2}} + O(\Delta x^4)$$

$$g.d.e \longrightarrow \frac{\partial^2}{\partial x^2} \Rightarrow \frac{\partial^3 T}{\partial t \partial x^2} = \alpha \frac{\partial^4 T}{\partial x^4}$$

$$g.d.e \longrightarrow \frac{\partial}{\partial t} \Rightarrow \frac{\partial^2 T}{\partial t^2} = \alpha\frac{\partial^3 T}{\partial t \partial x^2}$$

$$\frac{\partial^2 T}{\partial t^2} = \alpha\frac{\partial^3 T}{\partial t \partial x^2} = \alpha^2 \frac{\partial^4 T}{\partial x^4}$$

Term ① - term ②.

$$\Rightarrow \alpha\frac{\partial^4 T}{\partial x^4}\left[\frac{\alpha \Delta t}{2} - \frac{\Delta x^2}{12}\right]$$

$$\hookrightarrow O(\Delta t), O(\Delta x^2)$$
First order in time,
second order in space.

---

$$\alpha\frac{\partial^4 T}{\partial x^4}\left[\alpha\frac{\Delta t}{2} - \frac{\Delta x^2}{12}\right] = 0.$$

when $\alpha\frac{\Delta t}{2} - \frac{\Delta x^2}{12} = 0$

$$\Rightarrow \frac{\alpha \Delta t}{\Delta x^2} = \frac{1}{6}.$$

order of
Then $O(\Delta t^2), O(\Delta x^4)$ are errors
in time & space respectively.

Is it consistent? Yes! as $\Delta x \to 0$,
$\Delta t \to 0$

error $\to 0$.

$$\boxed{\text{Lecture 25: FTCS scheme}}$$

FTCS: $\dfrac{\partial T}{\partial t} = \alpha \dfrac{\partial^2 T}{\partial x^2}$.

$$\frac{T_i^{n+1} - T_i^n}{\Delta t} = \alpha\frac{T_{i+1}^n + T_{i-1}^n - 2T_i^n}{\Delta x^2}$$

Checking for stability: (whether numerical perturbations amplify able propagating or not).

$$T_i^{n+1} - T_i^n = r\left[T_{i+1}^n + T_{i-1}^n - 2T_i^n\right]$$

$$\Rightarrow T_i^{n+1} = (1-2r) T_i^n + r T_{i+1}^n + r T_{i-1}^n \quad —①$$

→ analogy with F.V. Discretization.

→ Explicit scheme. $T_i^{n+1}$ described in terms of previous time step.

Say, we get an approximate soln. That soln must also satisfy above eq?
i.e.

$$T_i^{*n+1} = (1-2r)T_i^{*n} + r T_{i+1}^{*n} + r T_{i-1}^{*n}$$
$$—②.$$

① − ③ ⟹

$$\left(T_i - T_i^{*K}\right)^{n+1} = (1-2r)\left(T_i - T_i^*\right)^{\circ} +$$

$$r\left(T_{i+1} - T_{i+1}^*\right)^{\circ} + r\left(T_{i-1} - T_{i-1}^*\right)^{\circ}$$

Say  $\epsilon = T - T^*$

Then  ↳(error)

$$\epsilon_i^{n+1} = (1-2r)\,\epsilon_i^{\circ} + r\,\epsilon_{i+1}^{\circ} + r\,\epsilon_{i-1}^{\circ}$$

errors satisfy the same eq? as discretization
  variable.

$$\epsilon = \epsilon(x,t)$$

Write this in terms of a Fourier series.

$$\epsilon(x,t) = \sum_j A Z_j^{\,m} e^{ikx} \qquad i = \sqrt{-1}$$

What is the particular form of $Z$?

A convenient form is  $e^{at}$, where $a = fn(j)$

why $e^{at}$?  exponential fn helps evaluate
the growth/decay behavior easily.

If  $e^{at+\Delta t} < e^{at} \longrightarrow$ decay;

$e^{at+\Delta t} > e^{at} \longrightarrow$ growth.

to assess whether there is exponential
growth/decay.

curr space, there is a periodicity.
  each
Over length of ∧ CV, repeatability observed.
      to
∴ corresponding ∧ wave number ($k$) (number of
waves over a time period), corresponding
wave length = cell length.

---

$$e^{at} e^{jkx} \qquad j = \sqrt{-1}$$

$$e^{a(t+\Delta t)} \cdot e^{jkx}$$

$$= (1-2r)\,e^{at}e^{jkx} + r\,e^{at}e^{jk(x+\Delta x)}$$

$$+ r\,e^{at}e^{jk(x-\Delta x)}$$

$A$ = Amplification factor = $\dfrac{e^{a(t+\Delta t)}}{e^{at}}$

$$A = (1-2r) + r\left(e^{jk\Delta x} + e^{jk(-\Delta x)}\right)$$

Check for stability.  $|A| < 1$.

If regardless of $r$, $|A| < 1$, then
  unconditionally stable.

If regardless of $r$, $|A| > 1$, then unconditionally
  unstable.

if for some values of $r$, $|A| < 1$, then
  conditionally stable.

$$A = (1-2r) + r\left(e^{j\theta} + e^{-j\theta}\right)$$

$$= (1-2r) + r\left(\cos\theta + j\sin\theta + \cos\theta - j\sin\theta\right)$$

$$= (1-2r) + 2r\cos\theta.$$

$$= 1 - 2r(1-\cos\theta)$$

$$= 1 - 4r\sin^2\!\left(\tfrac{\theta}{2}\right).$$

For stability.  $|A| \le 1$.

$$\Rightarrow \quad -1 \le A \le 1.$$

$$\boxed{-1 \le 1 - 4r\sin^2\!\left(\tfrac{\theta}{2}\right) \le 1}$$

$$\underset{\le 1}{\le 1} \quad \hookrightarrow \quad 4r\sin^2\!\left(\tfrac{\theta}{2}\right) \le 2.$$

$$r \le \frac{1}{2\sin^2(\theta/2)}$$

↳ conservative upper limit ⟹ $r \le \tfrac{1}{2}$

Right hand limit

$$1 - 4\gamma \sin^2\left(\frac{\theta}{2}\right) \leq +1$$

$$\Rightarrow 4\gamma \sin^2\left(\frac{\theta}{2}\right) \geq 0. \longrightarrow \text{always true!}$$

$$\gamma = \frac{\alpha \Delta t}{\Delta x^2} > F_o \leq \frac{1}{2}.$$

---

**Lecture 26:** CTCS scheme (Leap Frog Scheme) and Dufort-Frankel Scheme

---

B/ CTCS Scheme (Leap frog scheme).

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

$$\frac{T_i^{n+1} - T_i^{n-1}}{2\Delta t} = \alpha \frac{T_{i+1}^0 + T_{i-1}^0 - 2T_i^0}{\Delta x^2}$$

$$T_i^{n+1} - T_i^{n-1} = \frac{2\gamma}{\cancel{2}} \left[ T_{i+1}^0 + T_{i-1}^0 - 2T_i^0 \right] - ①$$

Approx sol^n:

$$T_i^{*n+1} - T_i^{*n-1} = 2\gamma \left[ T_{i+1}^{*0} + T_{i-1}^{*0} - 2T_i^{*0} \right]$$

$$① - ① \longrightarrow \qquad \qquad - ②$$

$$\mathcal{E}_i^{n+1} - \mathcal{E}_i^{n-1} = 2\gamma \left[ \mathcal{E}_{i+1}^n + \mathcal{E}_{i-1}^0 - 2\mathcal{E}_i^n \right]$$

Sub the form $\theta e^{at} e^{jkx}$.

$$e^{a(t+\Delta t)} e^{jkx} - e^{a(t-\Delta t)} e^{jkx}$$

$$= 2\gamma \left[ e^{at} e^{jk(x+\Delta x)} + e^{at} e^{jk(x-\Delta x)} \right.$$
$$\left. - 2 e^{at} e^{jkx} \right]$$

---

Divide both sides by $e^{at} e^{jkx}$:

$$\underbrace{e^{a\Delta t}}_{A} - e^{-a\Delta t} = 2\gamma \left[ e^{jk\Delta x} + e^{-jk\Delta x} - 2 \right]$$

$$A - \frac{1}{A} = 2\gamma \left[ \cancel{\cos} e^{j\theta} + e^{-j\theta} - 2 \right]$$

$$A - \frac{1}{A} = 2\gamma \left[ 2\cos\theta - 2 \right]$$

$$A - \frac{1}{A} = 4\gamma (\cos\theta - 1) = -8\gamma \sin^2(\theta/2)$$

$$A^2 + 8\gamma \sin^2(\theta/2) A - 1 = 0.$$

Quadratic eq^n in A.

Product of the roots has a magnitude of 1. But

$$A = \frac{-8\gamma \sin^2(\theta/2) \pm \sqrt{64\gamma^2 \sin^4(\theta/2) + 4}}{2.}$$

$$= -4\gamma \sin^2(\theta/2) \pm \sqrt{16\gamma^2 \sin^4(\theta/2) + 1}.$$

Greater magnitude of $A \longrightarrow$

$$\left| -4\gamma \sin^2(\theta/2) - \sqrt{16\gamma^2 \sin^4(\theta/2) + 1} \right| \longrightarrow \geq 1$$

$\Rightarrow$ unconditionally unstable!

Q/ The greed for higher accuracy $(CT \sim O(\varepsilon^2), FT \sim O(\varepsilon))$ led to us using CTCS over FTCS. But CTCS is unconditionally unstable while FTCS is conditionally stable.

Modification of this scheme will provide accuracy as well as stability. Such a modification → Dufort-Frankel Scheme.

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

$$\frac{T_i^{n+1} - T_i^{n-1}}{2\Delta t} = \alpha \left[ \frac{T_{i+1}^n + T_{i-1}^n - 2T_i^n}{\Delta x^2} \right].$$

Make adhoc change:

$$T_i^n = \frac{T_i^{n-1} + T_i^{n+1}}{2}$$

(temperature at a grid pt at an instant = average temperature at that grid pt at previous time and the next time).

This adhoc change :: CTCS ⟶ Dufort-Frankel Scheme.

$$T_i^{n+1} - T_i^{n-1} = 2r \left[ T_{i+1}^n + T_{i-1}^n - T_i^{n-1} - T_i^{n+1} \right]$$

$$\varepsilon_i^{n+1} - \varepsilon_i^{n-1} = 2r \left[ \varepsilon_{i+1}^n + \varepsilon_{i-1}^n - \varepsilon_i^{n-1} - \varepsilon_i^{n+1} \right]$$

$$e^{a(t+\Delta t)} e^{jkx} - e^{a(t-\Delta t)} e^{jkx}$$

$$= 2r \left[ e^{at} e^{jk(x+\Delta x)} + e^{at} e^{jk(x-\Delta x)} - e^{a(t-\Delta t)} e^{jkx} - e^{a(t+\Delta t)} e^{jkx} \right]$$

Divide by $e^{at} e^{jkx}$,

$$A - A^{-1} = 2r \left[ e^{jk\Delta x} + e^{-jk\Delta x} - A - A^{-1} \right]$$

$$A = \frac{A}{A} a$$

$$A(1+2r) - \frac{1}{A}(1 \mp 2r) = 2r(2\cos\theta).$$

$$(1+2r)A^2 - 4r\cos\theta\, A - (1-2r) = 0.$$

$$A = \frac{4r\cos\theta \pm \sqrt{16r^2\cos^2\theta + 4(1-4r^2)}}{2(1+2r)}.$$

$$= \frac{2 \cdot 4r\cos\theta \pm \sqrt{4r^2\cos^2\theta + 1 - 4r^2}}{2(1+2r)}$$

$$= \frac{2r\cos\theta \pm \sqrt{1 - 4r^2\sin^2\theta}}{(1+2r)}$$

Case 1: $2r < 1$.

Then $4r^2 < 1$.

$$\therefore 4r^2 \sin^2\theta < 1.$$

$$A = \frac{2r\cos\theta \pm \sqrt{1 - 4r^2\sin^2\theta}}{(1+2r)}$$

Take (+) for more conservative estimate,

$$A = \left( 2r\cos\theta + \sqrt{1 - 4r^2\sin^2\theta} \right) / (1+2r)$$

$$\leq \frac{1 + 2r\cos\theta}{1+2r} \leq 1.$$

Case 2:  $2r > 1$

$$A = \frac{2r\cos\theta \pm j\sqrt{4r^2\sin^2\theta - 1}}{1 + 2r}$$

$$|A|^2 = \frac{4r^2\cos^2\theta + 4r^2\sin^2\theta - 1}{(1+2r)^2}$$

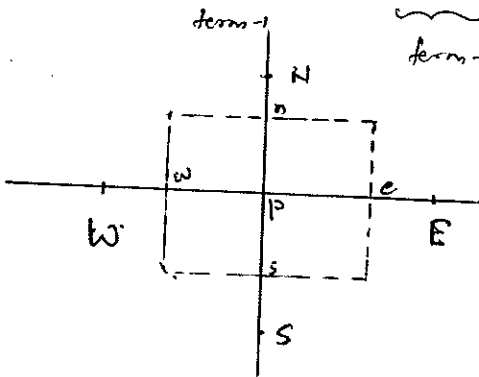$$= \frac{4r^2 - 1}{(1+2r)^2} = \frac{(2r+1)(2r-1)}{(1+2r)(1+2r)}$$

$$\approx \quad < 1.$$

This scheme → Conditionally stable, but not consistent (H/w :prove)

(i.e. truncation error as $\Delta t \to 0$ & $\Delta x \to 0$ doesn't nullify itself.

---

| Lecture 27: FV Discretization of 2D unsteady State Diffusion Type Problems: |
|---|

Eg: heat conduction:

$$\underbrace{\rho C_p \frac{\partial T}{\partial t}}_{\text{term-1}} = \underbrace{\frac{\partial}{\partial x}\left(k\frac{\partial T}{\partial x}\right)}_{\text{term-2}} + \underbrace{\frac{\partial}{\partial y}\left(k\frac{\partial T}{\partial y}\right)}_{\text{term-3}} + \underbrace{S}_{\text{term-4}}$$



Term-1 $\to \displaystyle\int_s^n\int_w^e\int_t^{t+\Delta t} \rho C_p \frac{\partial T}{\partial t}\, dt\, dx\, dy.$

→ can use piece wise const temperature profile (no derivative terms involved).

$\oint \mathcal{E} = \rho C_p (T_P - T_P^0)\Delta x \Delta y$

Then $T = T_P$.

---

Term-2:

$$\int_t^{t+\Delta t}\int_s^n\int_w^e \frac{\partial}{\partial x}\left(k\frac{\partial T}{\partial x}\right) dx\, dy\, dt.$$

$$\int_w^e \frac{\partial}{\partial x}\left(k\frac{\partial T}{\partial x}\right) dx. \quad \text{Take piecewise-linear profile.}$$

$$= k_e \frac{T_E - T_P}{\delta x_e} - k_w \frac{T_P - T_w}{\delta x_w}$$

Final integration form:

= Consider a fully implicit scheme.

$$\Rightarrow \left(k_e \frac{T_E - T_P}{\delta x_e} - k_w \frac{T_P - T_w}{\delta x_w}\right)\Delta t\, \Delta y.$$

Term 3:

$$\int_t^{t+\Delta t}\int_w^e\int_s^n \frac{\partial}{\partial y}\left(k\frac{\partial T}{\partial y}\right) dy\, dx\, dt$$

$$\to \left(k_n \frac{T_N - T_P}{\delta y_n} - k_s \frac{T_P - T_s}{\delta y_s}\right)\Delta t\, \Delta x.$$

Term 4:

$$\int_t^{t+\Delta t}\int_s^n\int_w^e (S_c + S_p T_P)\, dx\, dy\, dt.$$

$$= (S_c + S_p T_P)\, \Delta x\, \Delta y\, \Delta t$$

Assemble terms:

Divide all terms by $\Delta t$:

Then we get eqⁿ of the form:

$$a_P T_P = a_E T_E + a_w T_w + a_s T_s + a_N T_N + a_P^0 T_P^0 + b$$

$$a_E = \frac{k_e}{\delta x_e}\Delta y, \quad a_w = \frac{k_w}{\delta x_w}\Delta y,$$

$$a_s = \frac{k_s}{\delta y_s}\Delta x, \quad a_N = \frac{k_n}{\delta y_n}\Delta x, \quad a_P^0 = \frac{\rho C_p \Delta x \Delta y}{\Delta t}$$

$$a_p = a_E + a_W + a_S + a_N + \rho C_p \frac{\Delta x \Delta y}{\Delta t}$$

$$b = S_c \Delta x \Delta y. \qquad - S_p \Delta x \Delta y.$$

of the form:

$$a_p T_b = \sum a_{nb} T_{nb} + b \qquad (nb \rightarrow neighbour)$$

The coeffs $a_E, a_W, a_N, a_S$ etc physically represent thermal conductance. Why?

Take $a_E = \frac{k_e \Delta y}{\delta x_e}$. Assuming unit length normal to the $x$-$y$ plane, $(\Delta y \times 1) = $ area of face. So $a_E$ is of the form $a_E = \frac{k_e A}{L}$, which is the formula for conductance, $\frac{L}{k_e A} \rightarrow$ resistance.

By setting $\Delta t$ to very large number, this unsteady problem can be converted into a steady state problem as terms containing $(\Delta t)^{-1}$ becomes negligibly small.

(H/w): repeat same exercise for a fully explicit scheme.

---

- Solving system of algebraic eqⁿs:—

① $x + y = 2$

$2x + 3y = 5$

$E_1 \times 2 - E_2$

$y - y = -1$

$\Rightarrow y = 1$

$E_1: \quad x = 2 - y = 2 - 1 = 1 \qquad (x, y) = (1, 1)$

using method of elimination

② $x + y = 2.$
$2x + 2y = 4$

linearly dependent eqⁿs.
No. of independent eqs.
< no. of unknowns.

④    ⑤ $x + y = 0.$

$y = 2 - x.$    $2x + 2y = 0.$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ 2 - x \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} x + \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$

③ $x + y = 2.$

$2x + 2y = 5.$

$E_1 \times 2 \longrightarrow 2x + 2y = 4$

From $E_2 \rightarrow 2x + 2y = 5.$

Equate them both $\Rightarrow 4 = 5$ ✗

$\longrightarrow$ no solⁿ

(parallel straight
lines — no intersection
pts).

④ $x + y = 0.$    } If RHS = 0 $\rightarrow$

$2x + 3y = 0$    homogeneous.
system of
eqⁿs.

$x = 0, y = 0 \rightarrow$
trivial solⁿs.

no non-trivial solⁿs.

⑤ $x + y = 0.$

$2x + 2y = 0.$

trivial solⁿ $x = 0, y = 0$

$\infty$ non-trivial solⁿs.

$y = -x.$

what if number of eqns are large?

Use matrix forms.

$$① \rightarrow \begin{bmatrix} 1 & 1 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

$\underbrace{\qquad}$ Coeff matrix.

$$\begin{bmatrix} 1 & 1 & | & 2 \\ 2 & 3 & | & 5 \end{bmatrix} \rightarrow \text{Augmented matrix.}$$

---

**Lecture 28: Solutions to linear algebraic equations (contd).**

---

Identify rank of coeff matrix & augmented matrix. Why rank? B/c it gives an idea about the linear dependence/independence of equations in the system.

Rank of a matrix is $r$, if

(i) It has atleast one non-zero minor of order $r$.

(ii) all minors of order $> r$ vanishes ($= 0$).

- For eg:1,

~~Take~~ $R_c = 2 \longrightarrow$ coeff matrix

$R_A = 2 \longrightarrow$ aug. matrix.

Observation: $R_c = R_A = 2 = n$ (no. of eqns = no. of unknowns).

- For eg:2,

$$\begin{bmatrix} 1 & 1 & | & 2 \\ 2 & 2 & | & 4 \end{bmatrix}$$

$R_c = 1$

$R_A = 1$

Observation: $R_A = R_c = 1 < n$

- For eg:3,

$$\begin{bmatrix} 1 & 1 & | & 2 \\ 2 & 2 & | & 5 \end{bmatrix}$$

$R_c = 1$

$R_A = 2$

Obs: $R_c \neq R_A$

For a system of non-homogeneous eqns:

$\rightarrow R_c = R_A = n \Rightarrow$ unique soln.

$\rightarrow R_c = R_A < n \Rightarrow$ infinitely large no. of solutions.

$\rightarrow R_c \neq R_A \Rightarrow$ no solution.

For homogeneous eqns, since RHS is all 0, there is no need for augmented matrix.

For system of homogeneous eqns:

$\Delta \neq 0 \longrightarrow$ only trivial soln.

$\Delta = 0 \longrightarrow$ infinitely large no. of solutions

Its important to see the nature of soln for system of algebraic eqns. Why? B/c for (For eg,) well-posed problems, uniqueness is an important criteria. But if the system of non-homogeneous eqns has infinitely many solutions, the problem itself is ill defined & needs modifications.

Solution techniques for systems of linear algebraic equations:

→ Elimination

→ Iteration

→ Gradient search method.

Elimination method:

- $x_1 + x_2 = 2.$ $(E_1)$

  $2x_1 + 3x_2 = 5.$ $(E_2)$

  $2E_1 - E_2 →$ Effort is to eliminate $x_1$

- $x_1 + x_2 + x_3 = 3$ $(E_1)$

  $2x_1 + 2x_2 + 3x_3 = 7$ $(E_2)$

  $3x_1 + x_2 + 2x_3 = 6.$ $(E_3)$

$2 \times E_1 - E_2 ⟹ x_3 = 1.$

$E_3 - 3E_1 ⟹ -2x_2 - x_3 = -3$

$-2x_2 - 1 = -3.$

$x_2 = 1$

$$\begin{bmatrix} 1 & 1 & 1 & | & 3 \\ 2 & 2 & 3 & | & 7 \\ 3 & 1 & 2 & | & 6 \end{bmatrix} → \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

Reorder:

$x_1 + x_2 + x_3 = 3.$

$3x_1 + x_2 + 2x_3 = 6.$

$2x_1 + 2x_2 + 3x_3 = 7.$

$$\begin{bmatrix} 1 & 1 & 1 & | & 3 \\ 0 & -2 & -1 & | & -3 \\ 0 & 0 & 1 & | & 1 \end{bmatrix}$$ → Converted into upper triangular form.

Now, $x_3 = 1$

$$x_2 = \frac{3 - x_3}{2} = 1$$

$⟹ x_1 = 3 - x_2 - x_3 = 1$

We can break this method into two parts:

Part-I: forward elimination.

⟹ Convert $c$ to upper $\Delta$ form.

Part-II: Backwards substitution.

- $x_1 + x_2 + x_3 = 3.$

  $2x_1 + x_2 + 3x_3 = 6.$

  $3x_1 + 4x_2 + 2x_3 = 9.$

$$\begin{bmatrix} 1 & 1 & 1 & | & 3 \\ 2 & 1 & 3 & | & 6 \\ 3 & 4 & 2 & | & 9 \end{bmatrix}$$ 

$E_2 - 2E_1 → -x_2 + x_3 = 0$

$E_3 - 3E_1 → x_2 - x_3 = 0$

linear dependency involved!

$( E_3 = 5E_1 - E_2 ).$

**Ex**  
$$10x_1 + x_2 + x_3 = 12. \quad —\text{(E}_1\text{)}$$
$$x_1 + 10x_2 + x_3 = 12. \quad —\text{(E}_2\text{)}$$
$$x_1 + x_2 + 10x_3 = 12. \quad —\text{(E}_3\text{)}$$

**Step-1** : Row 1 as the pivotal row

$$\text{(E}_2\text{)} \rightarrow \text{(E}_2\text{)} - \frac{1}{10} \times \text{(E}_1\text{)}$$

$$\text{(E}_3\text{)} \rightarrow \text{(E}_3\text{)} + \frac{1}{10}\text{(E}_1\text{)}$$

$$\begin{bmatrix} 10 & 1 & 1 & | & 12 \\ 1 & 10 & 1 & | & 12 \\ 1 & 1 & 10 & | & 12 \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} 10 & 1 & 1 & | & 12 \\ 0 & 9.9 & 0.9 & | & 10.8 \\ 0 & 0.9 & 9.9 & | & 10.8 \end{bmatrix}$$

**Step-2** :

$$\text{(E}_3\text{)} \rightarrow \text{(E}_3\text{)} - \frac{0.9}{9.9}\text{(E}_2\text{)}$$

$$\rightarrow \begin{bmatrix} 10 & 1 & 1 & | & 12 \\ 0 & 9.9 & 0.9 & | & 10.8 \\ 0 & 0 & 9.818 & | & 9.818 \end{bmatrix}$$

2 steps required for 3 equations.

So in general, for n equations, it would require n-1 steps.

→ Upto this forward elimination.

Next: Backward substitution.

From (E₃),
$$x_3 = \frac{9.818}{9.818} = 1.$$

From (E₂),
$$9.9 x_2 = 10.8 -$$

& finally get $x_1$ from (E₁)

---

**Generalization of Gaussian Elimination:**

$$a_{11} x_1 + a_{12} x_2 + a_{13} x_3 + \cdots + a_{1n} x_n = b_1 \quad —E_1$$
$$a_{21} x_1 + a_{22} x_2 + \cdots + a_{2n} x_n = b_2 \quad —E_2$$
$$\vdots$$
$$a_{n1} x_1 + a_{n2} x_2 + \cdots + a_{nn} x_n = b_n \quad —E_n.$$

Row-1 pivotal

$$E_2 \rightarrow E_2 - \frac{a_{21}}{a_{11}} E_1$$

$$E_3 \rightarrow E_3 - \frac{a_{31}}{a_{11}} E_1$$

$$\vdots$$

$$E_i \rightarrow E_i - \frac{a_{i1}}{a_{11}} E_1$$

for step-1:
$$a_{ij} \rightarrow a_{ij} - \frac{a_{i1}}{a_{11}} \times a_{1j}$$

for step 2:
$$a_{ij} \rightarrow a_{ij} - \frac{a_{i2}}{a_{22}} \times a_{2j}$$

Generalize for step no. k. { (k for pivot)

$$a_{ij} \rightarrow a_{ij} - \frac{a_{ik}}{a_{kk}} \times a_{kj}$$

## Lecture 30: Gaussian Elimination and LU Decomposition. methods.

Formalize forward elimination:

for k=1 to n-1

    for j=k+1 to n    ∵ augmented matrix considered.

        for i=k+1 to n

$a_{ik}$ is already ready if 3rd loop comes in 2nd loop's position. So.

for k=1 to n-1

    for i=k+1 to n

        $R = a_{ik}/a_{kk}$;

        for j=k+1 to n+1

            $a_{ij} = a_{ij} - R \times a_{kj}$;

        end

    end

end

Formalize backward substitution:

$$x_n = \frac{b_n}{a_{nn}} \rightarrow a_{nn} x_n = b_n$$

$$a_{n-1,n-1} x_{n-1} + a_{n-1,n} x_n = b_{n-1}.$$

$$x_{n-1} = \frac{b_{n-1} - a_{n-1,n} x_n}{a_{n-1,n-1}}$$

$$a_{n-2,n-2} x_{n-2} + a_{n-2,n-1} x_{n-1} + a_{n-2,n} x_n = b_{n-2}.$$

$$x_{n-2} = \frac{1}{a_{n-2,n-2}}\left[ b_{n-2} - a_{n-2,n-1} x_{n-1} - a_{n-2,n} x_n \right]$$

---

$$x_{n-i} = \frac{1}{a_{n-i,n-i}}\left[ b_{n-i} - \left\{ \sum_{j=0}^{j=i-1} a_{n-i,n-j} x_{n-j} \right\} \right].$$

· Formal algo for backward substitution:

$$x_n = \frac{b_{nn}}{a_{nn}};$$

for i=1 to n-1

    Sum=0

    for y=0 to i-1

        Sum = Sum + $a_{n-i,n-j} x_{n-j}$

    end.

    $x_{n-i} = (b_{n-i} - Sum)/a_{n-i,n-i}$;

end

Assessment of number of computations:-

Forward elimination → $O(n^3)$ → (n for each loop).

Backward substitution → $O(n^2)$

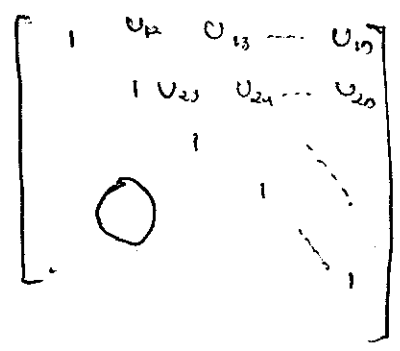Rate determining step: forward elimination

∴ computational complexity of the algo → $O(n^3)$.

LU decomposition evolved to reduce the complexity.
$$\sim O(n^2).$$

L-U decomposition technique:

Factorize $A = LU$. → Upper triangular matrix
                    ↳ Lower triangular matrix.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 & 0 & 0 \\ l_{21} & l_{22} & 0 & 0 \\ l_{31} & l_{32} & l_{33} & 0 \\ \vdots & & & 0 \\ l_{n1} & l_{n2} & l_{n3} & \cdots & l_{nn} \end{bmatrix} \times$$

$$\begin{bmatrix} 1 & U_{12} & U_{13} & \cdots & U_{1n} \\ & 1 & U_{23} & U_{24} \cdots & U_{2n} \\ & & 1 & & \\ & & & 1 & \ddots \\ & & & & \ddots \\ & & & & 1 \end{bmatrix}$$

## Crout's method:

$$a_{11} = l_{11} \qquad l_{11} = a_{11}$$
$$a_{21} = l_{21} \qquad l_{21} = a_{21}$$
$$\qquad\qquad\qquad l_{31} = a_{31}$$
$$\qquad\qquad\qquad \vdots$$
$$a_{n1} = l_{n1} \qquad l_{n1} = a_{n1}$$

$$\begin{array}{l|l}
A x = b & z = L^{-1} b \to O(n^2) \\
L[Ux] = b & Ux = z \\
L z = b & x = U^{-1} z \to O(n^2) \\
\;\;\;\hookrightarrow \text{triangular} & \left(\begin{array}{l}\text{Back} \\ \text{sub.}\end{array}\right) \\
\;\;\;\;\;\text{matrix}
\end{array}$$

(Forward elim.)

Even though complexity may seem $O(n^2)$, calculations required to factorize is has complexity $O(n^3)$. Thus it is no better than Gaussian elimination.

Then why use L-U decomposition method at all?

---

Lecture 31: Illustrative example of elimination method.

---

$$Ax = b$$
$$\downarrow \qquad \downarrow$$
Fixed   Variable.

→ Special case: when $A$ is symmetric and +ve definite ⇒ $U = L^T$ (Cholesky's L-U factorization).

Here number of calculations become half (But doesn't become $O(n^3) \to O(n^2)$; remains $O(n^3)$ itself)

Symmetric : $A = A^T$

+ve definite : $V^T A V > 0$.

Ex. $\varepsilon x_1 + x_2 = 1 \quad (E_1)$
$x_1 + x_2 = 2 \quad (E_2)$   $\varepsilon \to$ a small no.

Solve by Gaussian elimination

Forward elimination:

$$(E_2) - (E_1) \times \frac{1}{\varepsilon} \;\Rightarrow\; \left(1 - \frac{1}{\varepsilon}\right) x_2 = \left(2 - \frac{1}{\varepsilon}\right)$$

$$x_2 = \frac{2 - \frac{1}{\varepsilon}}{1 - \frac{1}{\varepsilon}} \approx 1$$

$$x_1 = \frac{1 - x_2}{\varepsilon} = 0 . \;\Big\}\; \times$$

∵ these sol^n doesn't satisfy the equations, they aren't the correct solutions.

Test whether reordering eq^n help or not:

Reorder the eq^ns:

$$x_1 + x_2 = 2 \;-\; (E_1)$$
$$\varepsilon x_1 + x_2 = 1 \;-\; (E_2).$$

Forward elimination:

$$(E_2) - \varepsilon(E_1) . \;\Rightarrow\; (1 - \varepsilon) x_2 = 1 - 2\varepsilon$$

Backward Substitution:
$$x_2 = \frac{1 - 2\varepsilon}{1 - \varepsilon} \approx 1 \to x_1 = 2 - x_2 = 1 \Big\}\checkmark$$

Reordering works!

Origin of the problem: Pivotal coefficient/diagonal entry being small

Reordering the equation makes diagonal entries not small.

What is the issue with diagonals being small?

- during forward elimination, division by diagonal element is required. Division by small number $\varepsilon$ makes the resulting number very large. This blowing off can over weigh any other coeff- number in the equation making them insignificant & making it difficult to extract out the difference b/w coeff numbers. $\longrightarrow$ leading to errors.

Reorder equations to reassign pivotal rows $\longrightarrow$ pivotization.

- Gaussian elimination need not work well in cases where there are issues in diagonal dominance.

---

- has complexity: $O(n)$
- also known as Thomas algorithm.



$$a_P \phi_P = a_E \phi_E + a_w \phi_w + b.$$
$$= \sum a_{nb} \phi_{nb} + b.$$
$$a_i \phi_i = b_i \phi_{i+1} + c_i \phi_{i-1} + d_i \quad --- \circledast.$$



Tridiagonal matrix.

Instead of using 2 indices for storage of one element, use 1 index for storing entries related to three diagonals together. $\longrightarrow$ 2D array storage system $\frac{}{to}$ linear storage system.

$a_1 \phi_1 = b_1 \phi_2 + d_1 \qquad (c_1 = 0)$

$\phi_1 = \dfrac{b_1 \phi_2 + d_1}{a_1} \longrightarrow f_1(\phi_2)$

$a_2 \phi_2 = b_2 \phi_3 + c_2 \overset{f(\phi_1)}{\cancel{\phi_1}} + d_2.$

$\phi_2 \overset{\downarrow}{=} f_2(\phi_3).$

$\vdots$

$\phi_3 = f_n (\cancel{\phi_{n+1}}^{\nearrow \text{const}})$  (∵ there is no $\phi_{n+1}$, its just a const)

In general,

$\phi_i = P_i \phi_{i+1} + Q_i \longrightarrow$ is the form of a  ③⊛  linear function.

$\mathcal{I}$mmediate step before:

$\phi_{i-1} = P_{i-1} \phi_i + Q_{i-1}. \quad \boxed{\text{⊛⊛}}$

Sub this recursive formula in ⊛ :  ④⊛

$a_i \phi_i = b_i \phi_{i+1} + c_i \left[ P_{i-1} \phi_i + Q_{i-1} \right] + d_i$

$\Rightarrow (a_i - c_i P_{i-1}) \phi_i = b_i \phi_{i+1} + d_i + c_i Q_{i-1}.$

$\phi_i = \dfrac{b_i}{a_i - c_i P_{i-1}} \phi_{i+1} + \dfrac{d_i}{a_i - c_i P_{i-1}} + \dfrac{c_i}{a_i - c_i P_{i-1}} Q_{i-1}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{④⊛}$

$\mathcal{C}$ompare ③⊛ and ④⊛ :-

$\gg P_i = \dfrac{b_i}{a_i - c_i P_{i-1}}$

$Q_i = \dfrac{d_i + c_i Q_{i-1}}{a_i - c_i P_{i-1}}.$

---

where $c_1 = 0$. (∵ there is not $\phi_0$) &
$b_N = 0$ (∵ there is no $\phi_{n+1}$)

$P_1 = \dfrac{b_1}{a_1}$ , & $Q_1 = \dfrac{d_1}{a_1}.$

$\phi_n = Q_N.$  (∵ $\phi_n = P_n \phi_{n+1} + Q_n$

$\qquad\qquad\qquad$ ∵ there is no $\phi_{n+1}$,

Above ↑ stiff → forward   $\overset{\text{like}}{\qquad}$  ∴ $\phi_n = Q_n$).
$\qquad\qquad$ elimination.

$\phi_i = P_i Q_{i+1} + Q_i \longrightarrow \text{w}^r$ to backward sub.

$\mathcal{S}$ummary of TDMA:

· Input  $a_i, b_i, c_i, d_i$

$P_1 = \dfrac{b_1}{a_1}$ , $Q_1 = \dfrac{d_1}{a_1}$

forward elimination:
for  $i = 2, N$
$P_i = \dfrac{b_i}{a_i - c_i P_{i-1}}$ ;  $Q_i = \dfrac{d_i + c_i Q_{i-1}}{a_i - c_i P_{i-1}}$
end

$\cancel{\phi_n} \phi_n = Q_N.$

Backward substitution:
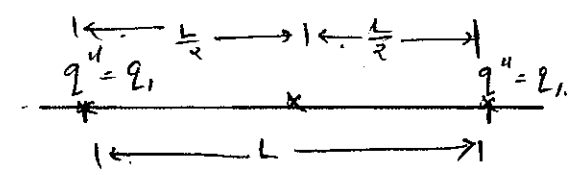for  $i = n-1, i, i--$
$\phi_i = P_i \phi_{i+1} + Q_i$
end

$O(N)$  complexity!

Will TDMA always work?

Take an example:-  $S = 0.$



$q'' = q_1 \qquad\qquad\qquad\qquad q'' = 2 q_1$

Find steady state T distribution.

Steps:

→ FV discretization / FD discretization.
→ sol$^n$ of eq$^n$ by TDMA

Governing equation:

$$\frac{d}{dx}\left(k\frac{dT}{dx}\right) = 0.$$

If $k$ const $\longrightarrow$ $\frac{d^2T}{dx^2} = 0.$

Applied to grid-point (2)

Fb discretization:

$$\frac{T_3 + T_1 - 2T_2}{(L/2)^2} = 0$$

$$\Rightarrow T_2 = \frac{T_1 + T_3}{2} \longrightarrow 2T_2 = T_1 + T_3.$$

For b/c at ①,

$$q'' = -k\frac{dT}{dx}\Big|_1 = k\left(\frac{T_1 - T_2}{(L/2)}\right)$$

$$T_1 = T_2 + \frac{q''L}{2k}$$

$$\longrightarrow T_1 = T_2 + \frac{1}{2}\alpha$$

BC at ③,

$$T_3 = T_2 - \frac{1}{2}\alpha \quad \text{is } a \text{ similar way.}$$

$\left\{\begin{array}{l}\text{Physically obvious as heat is transferred} \\ \text{from higher to lower temperature.} \\ \hspace{4cm} T_3 < T_2\end{array}\right\}$

$$a_1 T_1 = b_1 T_2 + d_1$$

$$a_1 = 1, \quad b_1 = 1, \quad d_1 = \alpha.$$

$$a_2 T_2 = b_2 T_3 + c_2 T_1 + d_2.$$

$$a_2 = 2, \quad b_2 = 1, c_2 = 1, d_2 = 0.$$

$$a_3 T_3 = b_3 c_3 T_1 + d_3.$$

$$a_3 = 1, \quad c_3 = 1, \quad d_3 = -\alpha$$

$$P_1 = \frac{b_1}{a_1} = 1, \quad Q_1 = \frac{d_1}{a_1} = \alpha.$$

$$P_2 = \frac{b_2}{a_2 - c_2 P_1} \qquad\qquad Q_2 = \frac{d_2 + c_2 Q_1}{a_2 - c_2 P_1}$$

$$= \frac{1}{2 - 1\times1} = \frac{1}{1} \qquad\qquad = \frac{0 + 1\times\alpha}{1}$$

$$P_3 = \frac{b_3}{a_3 - c_3 P_2}. \qquad\qquad = \alpha$$

$$= \frac{0}{1 - 1\times1} = \frac{0}{0} \longrightarrow \text{indeterminate.}$$

TDMA breaks down is this case.

why? Coeff matrix is singular
$(det = 0)$.

In $Ax = b$

$x = A^{-1}b$,

inverse of $A$ $A^{-1}$

$$= \frac{adj(A)}{det(A)}.$$

If $det(A) = 0 \longrightarrow A^{-1}$ has a problem.

✓

ill posed BVP.

Why? 2 flux conditions doesn't give additional information regarding the system.

~~ill pose~~ If ill posedness & physical unreality of the problem unt detected, the mathematics of the problem will naturally reveal it.

## Lecture 33: Elimination Methods: Error Analysis

**Prob** Given for L-U factorization by Croout's method, ($A = LU$) following steps are to be executed.

$$l_{i1} = a_{i1} \quad \text{for } i \in [1, n]$$

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \quad \text{where}$$

$$u_{1j} = \frac{a_{1j}}{l_{11}} \quad \text{for } j \in [2, n]$$

$$u_{ij} = \frac{1}{l_{ii}} \left[ a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj} \right] \quad \text{for } j \in [i+1, n]$$

Estimate the operational count if $A_{N \times N}$.

Numerical error → combined effect of errors intrinsic to method and errors due to the machine on which the algo is performed.

— Errors intrinsic to elimination method:—
{Terminology}

* norm of a vector:

Say, we have a vector $x$,

$$x = \{x_1, x_2, \dots, x_n\}$$

element of $x \to x_i$

$$\|x\|_p = \left[ \sum_i |x_i|^p \right]^{1/p}$$

**Ex** $\quad x = \{1, -2, 3, -4\}$

$$\|x\|_1 = |1| + |-2| + |3| + |-4|$$
$$\underset{=}{= 10}$$

$$\|x\|_2 = (1^2 + 2^2 + 3^2 + 4^2)^{1/2} = \sqrt{30}$$

↗ length of a vector.

---

$$\|x\|_\infty = \max |x_i|$$

$$\underset{=}{= 4}$$

**Norm of a matrix:**

$$\|A\| ?$$

Introduce a vector & find $\|Ax\|$.

$$\|A\| \longrightarrow \frac{\|Ax\|}{\|x\|} \leq \quad \|x\| = 1$$

**Prob:** $A = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$. Find $\|A\|_2$.

$$Ax = b ; \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \leq \|x\|_2 = 1$$

$$Ax = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 + x_2 \\ x_1 \end{bmatrix}.$$

$$\|Ax\|_2 = \sqrt{(x_1 + x_2)^2 + x_1^2}$$

subject to the constraint:

$$\|x\|_2 = 1 \Rightarrow \sqrt{x_1^2 + x_2^2} = 1. \quad \textcircled{1}$$

$$\|Ax\|_2 = \sqrt{2x_1^2 + x_2^2 + 2 x_1 x_2}$$

Using ①, $\|Ax\|_2$ can be written in terms of either $x_1$ or $x_2$ → unambiguously.

B/c of that we say, $\|Ax\|_2 = \max \dfrac{\|Ax\|_2}{\|x\|_2}$

Effectively find out

$$\max \left( \sqrt{2 x_1^2 + x_2^2 + 2 x_1 x_2} \right), \text{ given}$$

$$\underset{=}{x_2 = \sqrt{1 - x_1^2}}$$

$$y = 2x_1^2 + 1 - x_1^2 + 2x_1\sqrt{1-x_1^2}.$$

$$y = x_1^2 + 1 + 2x_1\sqrt{1-x_1^2}.$$

for max $y$, $\dfrac{dy}{dx} = 0$.

$\longrightarrow \quad 2x_1 + 2\left[\sqrt{1-x_1^2} + \dfrac{x_1(-2x_1)}{2\sqrt{1-x_1^2}}\right] = 0.$

$\Longrightarrow \quad x_1\sqrt{1-x_1^2} + \left[1-x_1^2 - x_1^2\right] = 0.$

$\Longrightarrow \quad x_1\sqrt{1-x_1^2} = 2x_1^2 - 1$

$x_1^2(1-x_1^2) = 4x_1^4 - 4x_1^2 + 1.$

$\Longrightarrow \quad 5x_1^4 - 5x_1^2 + 1 = 0.$

$x_1^2 = \dfrac{5 \pm \sqrt{5}}{2\times 5} = \dfrac{1}{2} \pm \dfrac{1}{2\sqrt{5}}$

$x_1^2 = \dfrac{1}{2} \pm \dfrac{\sqrt{5}}{10}$

Find condition for maxima, out of the two possible roots.

↳ That gives $x_1^2 \longrightarrow x_1 \longrightarrow x_2 \longrightarrow$ elements of $Ax$.

---

$$y = 2x_1^2 + 1 - x_1^2 + 2x_1\sqrt{1-x_1^2}.$$

$$y = x_1^2 + 1 + 2x_1\sqrt{1-x_1^2}.$$

For max $y$, $\dfrac{dy}{dx_1} = 0$

$\Rightarrow 2x_1 + 2\sqrt{1-x_1^2} + 2x_1 \times \dfrac{1(-2x_1)}{2\sqrt{1-x_1^2}} = 0$

$\Rightarrow x_1\sqrt{1-x_1^2} + 1 - x_1^2 - x_1^2 = 0$

$\Rightarrow x_1\sqrt{1-x_1^2} = 2x_1^2 - 1.$

$x_1^2(1-x_1^2)^2 = 4x_1^2$

Norms with special meaning for matrices: 1-norm & ∞-norm.

- 1-norm $\quad \|A\|_1 \longrightarrow$ Column sum norm.

$$\max_j \left\{\sum_i |A_{ij}|\right\}$$

- ∞-norm $\quad \|A\|_\infty \longrightarrow$ maximum row sum norm.

$$\max_i \sum_j |A_{ij}|.$$

$\begin{bmatrix} 1 & 4 & 6 \\ -5 & -2 & 1 \\ -8 & -1 & 3 \end{bmatrix} \quad \|A\|_1 \longrightarrow \max(14, 7, 10) = 14$

$\qquad\qquad\qquad\quad \|A\|_\infty \longrightarrow \max(11, 8, 12) = 12$

$A$

Some important properties of matrix norms:

1. $\|kA\| = |k|\|A\|$

2. $\|A+B\| \leq \|A\| + \|B\|$

3. $\|AB\| \leq \|A\|\cdot\|B\| \longrightarrow \max \|ABy\| = \dfrac{\max \|ABx\|}{\|x\|}$

$= \max \dfrac{\|ABx\|}{} \cdot \max \dfrac{\|Bx\|}{\|x\|}$

$= \max \dfrac{\|Ay\|}{\|y\|} \cdot \dfrac{\|Bx\|}{} \cdot \max \|Bx\|/\|x\|$

Etc.

$$\|AB\| = \max \frac{\|Ay\|}{\|y\|}, \max \frac{\|Bx\|}{\|x\|}$$

$$\leq \|A\| \cdot \|B\|$$

Error analysis of elimination methods:
Consider
$$Ax = b.$$

Let $x_{approx}$ be the approximate numerical sol?

$$A\underbrace{(x - x_{approx})}_{e\,(error)} = \underbrace{b - Ax_{approx}}_{\mathcal{R}\,(residue)}.$$

In error analysis, we look for an upper bound of error without actually knowing the exact solution. It involves estimating, not exact quantification.

Say we have $e_1 \sim 10^{-10}$ and $e_2 \sim 10^{-5}$. In which case error is more/less? We cannot really tell as it depends on the $x$ itself. What we require is relative error, not absolute error.

Relative error indicator: $\frac{\|e\|}{\|x\|}$

$$\|Ax\| \leq \|A\| \cdot \|x\|.$$
$$\|b\| \leq \|A\| \cdot \|x\|$$
$$\|x\| \geq \frac{\|b\|}{\|A\|}$$
$$x = A^{-1}b$$
$$\|x\| = \|A^{-1}b\| \leq \|A^{-1}\| \cdot \|b\|$$

$\|x\|$ If you have $Ae = \mathcal{R}$,

$$\|e\| \geq \frac{\|\mathcal{R}\|}{\|A\|} \quad and$$
$$\|e\| \leq \|A^{-1}\| \cdot \|\mathcal{R}\|.$$

Bounds:

$$\frac{\|e\|_{min}}{\|x\|_{max}} \leq \frac{\|e\|}{\|x\|} \leq \frac{\|e\|_{max}}{\|x\|_{min}}$$

$$\frac{\|\mathcal{R}\|}{\|A\| \cdot \|A^{-1}\| \cdot \|b\|} \leq \frac{\|e\|}{\|x\|} \leq \frac{\|A^{-1}\| \|\mathcal{R}\| \|A\|}{\|b\|}$$

Focus on upper bound (for conservative approach):

$$\frac{\|e\|}{\|x\|} \leq \boxed{\|A^{-1}\| \|A\|} \cdot \frac{\|\mathcal{R}\|}{\|b\|}$$

Conclusion: Even with small residual, the relative error may be large if $\|A^{-1}\| \cdot \|A\|$ is large. Therefore the largeness of $\|A^{-1}\| \|A\|$ determines the condition for accuracy of the system of equations one is solving.

Condition number: $C(A) = \|A\| \cdot \|A^{-1}\|$

Large $C(A) \Rightarrow$ even a small $\frac{\|\mathcal{R}\|}{\|b\|}$ can lead to large $\frac{\|e\|}{\|x\|}$.

$$\|AA^{-1}\| \leq \|A\| \cdot \|A^{-1}\|.$$
$$1 \leq \|A\| \cdot \|A^{-1}\|.$$
$$\Rightarrow C(A) \geq 1$$

Closer to 1, is better.

$C(A) \rightarrow$ very critical parameter for estimating error.

Example problem:

Given, $A = \begin{bmatrix} 2 & 1 \\ 2 & 1.01 \end{bmatrix}$.

Find $C(A)$.

$C(A) = \|A\| \cdot \|A^{-1}\|$

$A^{-1} = \begin{bmatrix} 1.01 & -1 \\ -2 & 2 \end{bmatrix} \times \dfrac{1}{\underbrace{\phantom{xxx}}_{(2.02-2)}}$

$\quad = \dfrac{1}{0.02} \begin{bmatrix} 1.01 & -1 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 50.05 & -50 \\ -100 & 100 \end{bmatrix}$

$\|A\|_\infty = 3.01$

$\|A^{-1}\|_\infty = 200$

$C(A) = \|A\|_\infty \|A^{-1}\|_\infty$

$\quad = 602 \longrightarrow$ quite large

$\quad\quad \Rightarrow$ ill conditioned system.

Source of this largeness: large smallness of the determinant

Cond$^t$ Inferences: Smaller the determinant, greater the chances of a larger Condition number & ill-condition of the system.

---

Iteration methods:

Basic philosophy: start with an initial guess for solution & iterate on it till you get a final solution that converges.

Say, we have equations:

$5x_1 + x_2 = 6 \quad\quad (E_1)$

$x_1 + 5x_2 = 6 \quad\quad (E_2)$

$x_1 = \dfrac{6 - x_2}{5} \quad\quad x_2 = \dfrac{6 - x_1}{5}$

Make an initial guess:

$x_1^{(0)} = 0 \quad , \quad x_2^{(0)} = 0.$

Now try to update on this initial guess:
Make an iterative formula out of the given system:

$x_1^{(k+1)} = \dfrac{6 - x_2^{(k)}}{5} \quad , \quad x_2^{(k+1)} = \dfrac{6 - x_1^{(k)}}{5}$

∵ If we write iteration formula in this manner, it is called Jacobii's method / Jacobii's iteration scheme.

S-1
$x_1^{(1)} = \dfrac{6-0}{5} = \dfrac{6}{5} \quad x_2^{(1)} = \dfrac{6-0}{5} = \dfrac{6}{5}$

S-2
$x_1^{(2)} = \dfrac{6 - x_2^{(1)}}{5} = \dfrac{6 - (6/5)}{5} = \dfrac{24}{25}$

$x_2^{(2)} = \dfrac{6 - x_1^{(1)}}{5} = \dfrac{6 - (6/5)}{5} = \dfrac{24}{25}$

S-3
S-4
S-5
⋮

$\longrightarrow x_1 \longrightarrow 1 \;, \; x_2 \longrightarrow 1 \quad$ convergence

Convergence $\Rightarrow$ result b/w the current & previous steps doesn't differ substantially. (within some tolerance) ←

To update iterations with a faster rate, use

$$x_1^{(k+1)} = \frac{6 - x_2^k}{5}$$

$$x_2^{(k+1)} = \frac{6 - x_1^{(k+1)}}{5}$$

(instead of $x_1^k$, we use more updated version).

If you do that, then this becomes Gauss-Siedel method.

$$x_1^{(1)} = \frac{6 - x_2^{(0)}}{5} = 6/5.$$

$$x_2^{(1)} = \frac{6 - x_1^{(1)}}{5} = \frac{6 - 6/5}{5} = \frac{24}{25}$$

$$x_1^{(2)} = \frac{6 - x_2^{(1)}}{5} = \frac{6 - 24/25}{5} = \frac{126}{125}$$

$$x_2^{(2)} = \frac{6 - x_1^{(2)}}{5} = \frac{6 - 126/125}{5} = \frac{624}{625}$$

Here within two steps, they sol$^n$ is converged very fast.

• where is the guarantee that the scheme will converge or not ?

---

Generalized analysis of the iterative methods: ④

$$a_{11} x_1 + a_{12} x_2 + a_{13} x_3 + \cdots + a_{1n} x_n = b_1$$
$$a_{21} x_1 + a_{22} x_2 + a_{23} x_3 + \cdots + a_{2n} x_n = b_2$$
$$\vdots \qquad \vdots \qquad \vdots \qquad \ddots \qquad \vdots \qquad \vdots$$
$$a_{m1} x_1 + a_{n2} x_2 + a_{n3} x_3 + \cdots + a_{nn} x_n = b_n$$

$$[A][x] = [b]$$

$$x_1^{(k+1)} = \frac{b_1 - \left( a_{12} x_2^{(k)} + a_{13} x_3^{(k)} + \cdots + a_{1n} x_n^{(k)} \right)}{a_{11}}$$

Jacobi's method:

$$x_2^{(k+1)} = \frac{b_2 - \left( a_{21} x_1^{(k)} + a_{23} x_3^{(k)} + \cdots + a_{2n} x_n^{(k)} \right)}{a_{12}}$$

Gauss-Siedel:

$$x_2^{(k+1)} = \frac{b_2 - \left( a_{21} x_1^{(k+1)} + a_{23} x_3^{(k)} + \cdots + a_{2n} x_n^{(k)} \right)}{a_{22}}$$

$$[A][x] = [b]$$
↓
$$[L] + [D] + [U].$$

{ Don't confuse with LU decomposition }

lower diagonal upper
diag. diag.

$$L = \begin{bmatrix} a_{21} & & & O \\ a_{31} & a_{32} & & \\ a_{41} & a_{42} & a_{43} & \\ & \ddots & & \\ a_{n1} & a_{n2} & \cdots & a_{nn-1} \end{bmatrix}$$

$$D = \begin{bmatrix} a_{11} & & & O \\ & a_{22} & & \\ & & a_{33} & \\ O & & & \ddots \\ & & & & a_{nn} \end{bmatrix}$$

$$U = \begin{bmatrix} a_{12} & a_{13} & \cdots & a_{1n} \\ & a_{23} & a_{24} & \cdots & a_{2n} \\ & & a_{34} & a_{35} & a_{3n} \\ O. & & & \\ & & & a_{n-1,n} \end{bmatrix}$$

Jacobi: $x_2^{(k+1)} = \dfrac{b_2 - (\overbrace{a_{21} x_1^{(k)}}^{L} + \overbrace{a_{23} x_3^{(k)}}^{U} + \cdots + a_{2n} x_n^{(k)})}{a_{22}}$

Gauss-Siedel:

$x_2^{(k+1)} = \dfrac{b_2 - (\overbrace{a_{21} x_1^{(k+1)}}^{L} + \overbrace{a_{23} x_2^{(k)}}^{U} + \cdots + a_{2n} x_n^{(k)})}{a_{22}}$

In place of $Ax = b$,

$(L + D + U) x = b$.

Jacobi: $D x^{k+1} + (L+U) x^k = b$.

$x^{(k+1)} = -D^{-1}(L+U) x^{(k)} + D^{-1} b$.

$x^{(k+1)} = M x^{(k)} + C$,

where $M = -D^{-1}(L+U)$

$C = D^{-1} b$

Gauss-Siedel:

$D X^{(k+1)} + L X^{(k+1)} + U x^{(k)} = b$.

$(L + D) x^{(k+1)} = -U x^{(k)} + b$

$x^{(k+1)} = -(L+D)^{-1} U x^{(k)} + (L+D)^{-1} b$.

$X^{(k+1)} = M x^{(k)} + C$

$M = -(L+D)^{-1} U$

$C = (L+D)^{-1} b$

---

Say, $x^{(k+1)} = M x^{(k)} + C$

$x^*$ is actual sol$^n$.

$x^* = M x^* + C$.

$x^{(k+1)} - x^* = M(x^k - x^*)$

error in the $(k+1)^{th}$ step.

$e^{(k+1)}$

$e^{(k+1)} = M e^{(k)}$

$e^{(1)} = M e^{(0)}$

$e^{(2)} = M e^{(1)}$

$\qquad = M^2 e^{(0)}$

$\vdots$

$e^{(k)} = M^k e^{(0)}$

$\dfrac{\|e^k\|}{\|e^0\|} < 1 \longrightarrow$ requirement for convergence.

$\Rightarrow \|M^k\| < 1$

Not easy for raw computation of $M^k$.
Use eigenvalues & eigenvectors of $M$ for ease of representation.

• Let $d_i$ & $v_i$ be correspondingly the eigenvalues & eigenvectors of $M$.

$e^0 = a_1 v_1 + a_2 v_2 + a_3 v_3 + \cdots + a_n v_n$.

(eigenvalues are arranged in a way that

$|\lambda_1| > |\lambda_2| > |\lambda_3| \cdots > |\lambda_n|)$

$M e^0 = a_1 M v_1 + a_2 M v_2 + a_3 M v_3 + \cdots$
$\qquad\qquad + a_n M v_n$

$M v_1 = d_1 v_1 ; \quad M v_2 = d_2 v_2 \cdots ; \quad M v_n = \lambda_n v_n$.

$$M^2 e^0 = a_1 \lambda_1 M v_1 + a_2 \lambda_2 M v_2 + \cdots + a_n \lambda_n M v_n$$

$$= a_1 \lambda_1^2 v_1 + a_2 \lambda_2^2 v_2 + \cdots + a_n \lambda_n^2 v_n$$

$$e^k : \qquad M^k e^0 = a_1 \lambda_1^k v_1 + a_2 \lambda_2^k v_2 + \cdots + a_n \lambda_n^k v_n.$$

Compare leading order term of $M^k e^0$
with that of $e^0$.

$$\Rightarrow \quad |\lambda_1|^k < 1 \longrightarrow \frac{\|e^k\|}{\|e^0\|}$$

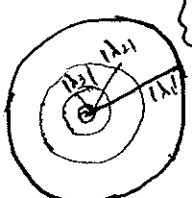$$= \frac{\| a_1 \lambda^k v_1 + LOT \|}{\| a_1 \lambda v_1 + LOT \|}$$

$$\sim |\lambda|^k.$$

$|\lambda_{max}| \rightarrow$ spectral radius of convergence.

{ In certain physical problems, eigenvalues represents the natural frequencies of the system. Something to do with frequency $\rightarrow$ "spectral" radius }

For convergence

{ Given
$$|\lambda_1| > |\lambda_2| > |\lambda_3| \cdots$$
$$|\lambda_1| = |\lambda_{max}|$$ }

$\ast$ radius of the 'biggest circle should be less than 1.

• Sufficient conditions for convergence :
(If that condition is satisfied, you'll definitely satisfy convergence. But you may also have convergence without satisfying that condition).

---

• Rate of convergence :

Requirement : no. of iterations to converge.

Say we require 'm' decimal accuracy.
Then,
$$\frac{\|e^k\|}{\|e^0\|} < 10^{-m}$$

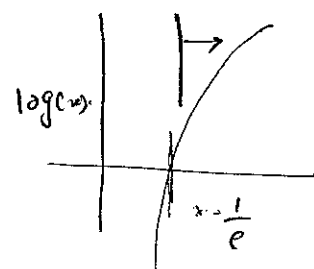$$e^k < 10^{-m}$$

$$k \log_{10} P < -m$$

$$m < k \log_{10} \left(\frac{1}{e}\right)$$

$$k > \boxed{\frac{m}{\log_{10}\left(\frac{1}{e}\right)}} \rightarrow R = \text{rate of convergence.}$$

Lecture 37: Further discussion on Iterative Methods.

• Smaller spectral radius $\Rightarrow$ better rate of convergence.

$$Mv = \lambda v$$

$$\|Mv\| \leq \|M\| \cdot \|v\|$$
$$\downarrow$$
$$|\lambda| \|v\| \leq \|M\| \cdot \|v\|$$

$$\Rightarrow |\lambda| \leq \|M\| \Rightarrow P \leq \|M\|.$$

Spectral radius upper bound.

• Which norm to choose?
An estimate of $P \rightarrow \max [\|M\|_1, \|M\|_\infty]$

- Sufficient condition for convergence:

$$\max\left(\|M\|_1, \|M\|_\infty\right) \leq 1$$

Jacobi's method:

$$M = -D^{-1}(L+U)$$

$$\begin{bmatrix} a_{11} & & & \bigcirc \\ & a_{22} & & \\ & & a_{33} & \\ & & & \ddots \\ \bigcirc & & & & a_{nn} \end{bmatrix} , \begin{bmatrix} a_{21} & & & \bigcirc \\ a_{31} & a_{32} & & \\ a_{41} & a_{42} \cdots & & \\ a_{n1} & a_{n2} \cdots & & a_{nn-1} \end{bmatrix}$$

$$\underbrace{\qquad}_{D} \qquad \underbrace{\qquad}_{L}$$

$$\begin{bmatrix} & a_{12} & a_{13} \cdots & a_{1n} \\ & & a_{23} & a_{24} \cdots a_{2n} \\ & & & a_{34} \cdots a_{3n} \\ \bigcirc & & & \ddots \\ & & & & a_{n-1n} \end{bmatrix}$$

$$\underbrace{\qquad}_{U}$$

$$\|M\|_R = \|M\|_\infty \rightarrow \frac{\sum\limits_{i \neq j} |a_{ij}|}{|a_{ii}|} \longrightarrow$$

Should be $\leq 1$ to satisfy sufficient condition for convergence.

$$\downarrow$$

$$\frac{\sum |a_{nb}|}{|a_p|} \leq 1.$$

for 1D systems $\longrightarrow$ TDM.

2D " $\longrightarrow$ Penta diagonal system.

3D " $\longrightarrow$ 7-diag system.

But matrix is generally sparse (sparse matrices).

---

Scarborough Criteria for of sufficient condition for convergence in Gauss-Siedel method.

$$\sum \frac{|a_{nb}|}{|a_p|} \leq 1 \quad \text{for all eq}^{os}$$

$$< 1 \quad \text{for at least one eq}^o.$$

Lecture 38: Illustrative Examples of Iterative methods

Ex 1-D steady state heat conduction in a rod with uniform $k$, $S=0$



$q''_1$ (given) ... $q''$ (given)

gde: $\dfrac{d^2 T}{dx^2} = 0$

CD: $\dfrac{T_3 + T_1 - 2T_2}{\Delta x^2} = 0$

$$\Rightarrow T_2 = \frac{T_1 + T_3}{2}$$

For $\$1$, $q'' = -k\left(\dfrac{T_2 - T_1}{\Delta x}\right)$

$$T_1 = T_2 + \left(\frac{q'' \Delta x}{k}\right)^c$$

For $\$3$, $q'' = -k\dfrac{(T_3 - T_2)}{\Delta x}$

$$T_3 = T_2 - \left(\frac{q'' \Delta x}{k}\right)^c$$

Eq 1 $\longrightarrow$ $T_1 = T_2 + \mathcal{T}_c$ $\longrightarrow$ $\sum \dfrac{|a_{bb}|}{|a_p|} = 1$

Eq 2 $\longrightarrow$ $2T_2 = T_1 + T_3$ $\longrightarrow$ $\sum \dfrac{|a_{nb}|}{|a_p|} = \dfrac{2}{2} = 1$

Eq 3 $\longrightarrow$ $T_3 = T_2 - T_c$ $\longrightarrow$ $\sum \dfrac{|a_{nb}|}{|a_p|} = \dfrac{1}{1} = 1$

Ill-posed problem ∴ we need atleast Temp specified at a boundary.

In all cases $\sum \dfrac{|a_{ob}|}{|a_p|} \leq 1$. In none of

The cases is it $< 1$. Can't satisfy

Scarborough's criteria.

So specify at gp1 $T_{given} = T_1^*$ instead of $2''$

Then for gp 2,

$2T_2 = T_1 + T_3$.

$\Rightarrow 2T_2 = T_1^* + T_3$.

$\left\{ \begin{array}{l} E_2'' \text{ no more valid} \\ \text{as } 2'' \text{ at gp1} \\ \text{not known} \end{array} \right\}$

∴ $T_1^*$ is already known, mathematically it is no longer a neighbour to $T_2$. only $T_3$ is the neighbour

∴ for modified eq$^n$ ②,

$\sum \dfrac{|a_{nb}|}{|a_p|} = \dfrac{1}{2} \neq 1 < 1$

With re-definition of the problem, Scarborough criterion is satisfied.

---

Ex Consider the system,

$2x_1 + 3x_2 + 10x_3 = 10$

$5x_1 - 2x_2 + 2x_3 = 5$

$x_1 + 10x_2 + 5x_3 = 6$.

Q: Is it possible to follow an iterative method (say Jacobi iteration) for the above system with guaranteed convergence?
$\rightarrow$ If yes, what is the estimated no. of iterations to achieve 4 decimal accuracy?

. $\sum\limits_{i \neq j} \dfrac{|a_{ij}|}{|a_{ii}|} \leq 1 \longrightarrow$ representative of the diagonal dominance.

Diagonal terms are dominating over the sum of the off diagonal terms.

$R_1 \longrightarrow \dfrac{|3|+|10|}{|2|} \nleq 1$ doesn't satisfy sufficient condition.

How to ≥ somehow satisfy the suff cond?
Find an eq$^n$ from the set of eq$^n$s where the first term's coeff is largest. Here eq-2 has met that criteria. So swap eq-1 with eq-2 $[eq-1^* = eq-2]$.

Now $R_1^* \longrightarrow \dfrac{|-2|+|+2|}{|5|} = \dfrac{4}{5} \leq 1$.

sufficient condition is satisfied!

$ll^{rly}$ make eq-3 $\longrightarrow$ eq-2$^*$ & eq-1 $\rightarrow$ eq-3$^*$

Thus reordering equations can help satisfy the suff. cond.

eq-1* $\dfrac{|2|+|2|}{|5|} = \dfrac{4}{5}$

eq-2* $\dfrac{|5|+|1|}{|10|} = \dfrac{6}{16}$

eq-3* $\dfrac{|2|+|3|}{|10|} = \dfrac{5}{10}$

Row sum norm = 0.8 as max = 0.8

$l = \max\left(\|M\|_R, \|M\|_c\right)$

Column sum norm $\|M\|_c$

$$\begin{array}{ccc} 5 & -2 & 2 \\ 1 & 10 & 5 \\ 2 & 8 & 10 \end{array}$$

col-1: 3/5
col-2: 5/10    max: 0.7
col-3: 7/10    $\|M\|_c = 0.7$

$l = \max(0.8, 0.7) = 0.8$

$\therefore h = \log_{10}\left(\dfrac{1}{\rho}\right) = \log_{10}(1.25)$

Ex. For a linear system of size N, the $i^{th}$ eigenvalue of the Jacobi iteration matrix $[M = -D^{-1}(L+U)]$ is given by $\lambda_i = \cos\left(\dfrac{i\pi}{N+1}\right)$ where $i=1,2,\dots,N$. It is also known that for sufficiently large values of N, $\cos\dfrac{\pi}{N+1} \approx \dfrac{1-\pi^2}{2(N+1)^2} \approx \exp\left(\dfrac{-\pi^2}{2N^2}\right)$.

If the size of the coefficient matrix changes from size $10^3 \times 10^3$ to $10^2 \times 10^2$, to what proportion would the total no. of iterations expected to achieve a desired total level of accuracy will decrease?

---

for $\|M\|$ at, also i=1,

$\cos\left(\dfrac{\pi}{N+1}\right) \approx \exp\left(\dfrac{-\pi^2}{2N^2}\right)$

$k \longrightarrow \dfrac{M}{R}$

$\dfrac{k_2}{k_1} \approx \dfrac{R_1}{R_2}$ ( keeping same level of accuracy, but diff coeff matrix x)

$\rightarrow \dfrac{k_2}{k_1} = \dfrac{-\log_{10}(|\lambda|_{max}) \to N=10^3}{-\log_{10}(|\lambda_i|_{max}) \to N=10^2}$

$\approx 10^{-2}$

$k_2 = 10^{-2} k_1.$

Lecture 38: Gradient Search Based Methods

f: function

Grad f → represents maximum rate of change of f.

Say, $f = \dfrac{1}{2}x^T A x$    Obj: Solve Ax=b.

$- b^T x + c,$

C arbitrary const.

Two restrictions:
- A symmetric
- A +ve definite.

Next obj, find cond. for min f=?

min f ⟹ $\nabla f = 0$.

$f = \dfrac{1}{2}[x_1 \ x_2 \cdots x_n] \begin{bmatrix} a_{11} & a_{12} & a_{13} \cdots a_{1n} \\ a_{21} & a_{22} & a_{23} \cdots a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & a_{n3} \cdots a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$

Left column:

$$- [b_1 \; b_2 \; b_3 \cdots b_n] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} + c$$

$$\nabla f = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\[2mm] \dfrac{\partial f}{\partial x_2} \\ \vdots \\ \dfrac{\partial f}{\partial x_n} \end{bmatrix}$$

$$f = \frac{1}{2} [ x_1 \; x_2 \cdots x_n] \begin{bmatrix} a_{11} x_1 + a_{12} x_1 + \cdots + a_{1n} x_n \\ a_{21} x_1 + a_{22} x_2 + \cdots + a_{2n} x_n \\ a_{n1} x_1 + a_{n2} x_2 + \cdots + a_{nn} x_n \end{bmatrix}$$

$$- [ b_1 \; b_1 \cdots b_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + c$$

$$f = \frac{1}{2} x_1 ( a_{11} x_1 + \cdots + a_{1n} x_n) + \frac{1}{2} x_2 ( a_{21} x_1 + \cdots + a_{2n} x_n)$$

$$+ \cdots + \frac{1}{2} x_n ( a_{n1} x_1 + \cdots + a_{nn} x_n)$$

$$- ( b_1 x_1 + b_2 x_2 + \cdots + b_n x_n) + c.$$

$$\frac{\partial f}{\partial x_1} = \frac{1}{2} ( a_{11} x_1 + \cdots + a_{1n} x_n) + \frac{1}{2} x_1 ( a_{11}) +$$

$$\frac{1}{2} a_{21} x_2 + \frac{1}{2} a_{31} x_3 + \cdots + \frac{1}{2} a_{n1} x_n$$

$$- b_1$$

$$= \frac{1}{2} ( a_{11} x_1 + \cdots + a_{1n} x_n) +$$

$$\frac{1}{2} ( a_{11} x_1 + \cdots + a_{n1} x_n) - b_1 + \cdots$$

$$\because A \text{ is symmetric}, \quad a_{ij} = a_{ji}$$

$$\Rightarrow \frac{\partial f}{\partial x_1} = ( a_{11} x_1 + a_{12} x_2 + \cdots + a_{1n} x_n) - b_1$$

$$= [ a_{11} \; a_{12} \cdots a_{1n}] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} - b_1.$$

Right column:

$$\frac{\partial f}{\partial x_2} = [ a_{21} \; a_{22} \cdots a_{2n}] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} - b_2.$$

$$\Rightarrow \nabla f = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\[2mm] \dfrac{\partial f}{\partial x_2} \\ \vdots \\ \dfrac{\partial f}{\partial x_n} \end{bmatrix}$$

$$= \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} - \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

$$= Ax - b$$

$$\nabla f = 0$$

$$\Rightarrow Ax - b = 0$$

$$\Rightarrow Ax = b.$$

Getting a sol$^n$ $Ax = b$ is as good as extremizing $f = \frac{1}{2} x^T A x - b^T x + c.$

we're actually doing minimization of $f$ & not maximization

How to show it? Why minimization?

$$f_{(x)} = \frac{1}{2} x^T A x - b^T x + c$$

$$f_{(x+e)}$$

if $f_{(x)} < f_{(x+e)} \; \forall e$, then $f_{(x)}$ is minimum.

$$f_{(x+e)} = \frac{1}{2} (x+e)^T A (x+e) - b^T (x+e)$$

$$+ c$$

$$= \frac{1}{2} ( x^T A x + e^T A x + x^T A e + e^T A e)$$

$$- b^T x - b^T e + c.$$

$$= \left[ \frac{\nu}{2}(x^T A x) - b^T x + c \right] + e^T A x$$

• $x^T A e = (x^T A e)^T$

  (transpose of scalar is scalar itself)

  $= (A e)^T (x^T)^T$

  $= e^T A^T x$    $\Big($ A symmetric;

  $= e^T A x$       $\quad A^T = A \Big).$

• $b^T e = (b^T e)^T$    (again, scalar).

  $= e^T (b^T)^T$

  $= e^T b.$

• $f(x+e) = f(x) + e^T A x - e^T b + \frac{1}{2} e^T A e$

  $= f(x) + e^T(Ax - b) + \frac{1}{2} e^T A e$

B/c $Ax = b$, $e^T(Ax-b) = 0$.

  $= f(x) + e^T(Ax-b) + \boxed{\frac{1}{2} e^T A e}$

  $\geq 0$ for arbitrary $e$

  $\Big($ ∵ $\overset{A \text{ is}}{C}$ positive definite $\Big)$.

$\Rightarrow f(x+e) \geq f(x).$

  $\Rightarrow f(x)$ is a minimum.

— Many gradient search methods!

1. Steepest descent method.

Ex   Solve $Ax = b$, where $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$, $b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

---

And $f = \frac{1}{2} x^T A x - b^T x + c$

$f = \frac{1}{2}[x_1 \ x_2]\begin{bmatrix}1 & 0 \\ 0 & 2\end{bmatrix}\begin{bmatrix}x_1 \\ x_2\end{bmatrix} - [1 \ 1]\begin{bmatrix}x_1 \\ x_2\end{bmatrix} + c$

$= \frac{1}{2}[x_1 \ x_2]\begin{bmatrix}x_1 \\ 2x_2\end{bmatrix} - [x_1 + x_2] + c$

$= \frac{1}{2}x_1^2 + x_2^2 - x_1 - x_2 + c$

what does $f = 0$ represent?

$\frac{1}{2}(x_1^2 - 2x_1 + \underset{1^2}{1}) + (x_2^2 - x_2 + (\frac{1}{2})^2)$

$\qquad\qquad\qquad +c = 0$

$\qquad\qquad\qquad -\frac{3}{4}$

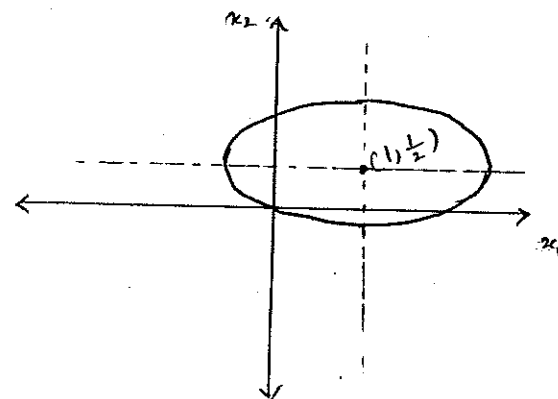$= \frac{(x_1-1)^2}{2} + \frac{(x_2 - \frac{1}{2})^2}{1} = \frac{3}{4} - c$

If $f$ satisfies $(0,0)$,

then $\quad \frac{1}{2} + \frac{1}{4} = \frac{3}{4} - c$

$\qquad\qquad \Rightarrow c = 0$

{ ∵ choice of $c$ is arbitrary, we can choose our curve to pass through a pt s $c=0$. }

$\frac{(x_1-1)^2}{2} + \frac{(x_2-\frac{1}{2})^2}{1} = \frac{3}{4}$

$\frac{(x_1-1)^2}{(\sqrt{3/2})^2} + \frac{(x_2-\frac{1}{2})^2}{(\sqrt{3}/2)^2} = 1.$

Say we start with $(x_1, x_2) = (0,0)$ & try to reach the actual solution. We move to a direction along the gradient $(1, \frac{1}{2})$ direction.
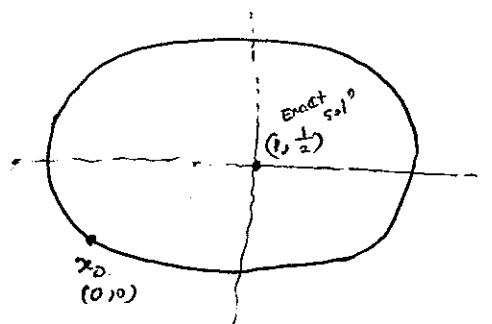
$$\nabla f = Ax - b.$$
$$= -r_1.$$

After moving some distance. We again stop & move in the gradient direction. The objective then becomes to find how much to travel in each segment.

---

## Lecture 40: Steepest descent method
### contd's

$f(x) =$ Steepest descent method:

— move along direction of maximum rate of change.



$A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ $b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$\nabla f = Ax - b$
$$= -r_1.$$

$x^1 = x^0 + \alpha_0 r_0.$

$r_0 = \nabla f |_{x^0}$

$\alpha_0 \rightarrow$ how much to move in that direction.

$r_0 \rightarrow$ tells direction in which you're moving

$$f(x^1) = \frac{1}{2}(x^0 + \alpha_0 r_0)^T A (x^0 + \alpha_0 r_0)$$
$$- b^T (x^0 + \alpha_0 r_0) + c$$

For $f$ to be min, $\frac{\partial f}{\partial \alpha_0} = 0.$

$$\frac{1}{2} r_0^T A (x^0 + \alpha_0 r_0) + \frac{1}{2}(x^0 + \alpha_0 r_0)^T A r_0$$
$$- b^T r_0 = 0.$$

---

$$\alpha_0 r_0^T A r_0 + \frac{1}{2} r_0^T A x^0 + \frac{1}{2} x^{0T} A r_0$$
$$- b^T r_0 = 0.$$

$x^{0T} A r_0 = (x^{0T} A r_0)^T$
$$= (A r_0)^T (x^0)^T$$
$$= r_0^T A^T x^0$$
$$= r_0^T A x^0 \qquad (A^T = A)$$

$b^T r_0 = (b^T r_0)^T = r_0^T b.$

$$\alpha_0 r_0^T A r_0 + r_0^T \underbrace{(A x^0 - b)}_{-r_0} = 0$$

$$\Rightarrow \alpha_0 = \frac{r_0^T r_0}{r_0^T A r_0} \qquad —— ①$$

$x^1 = x^0 + \alpha_0 r_0.$

$\alpha_1 = \dfrac{r_1^T r_1}{r_1^T A r_1.}$

$x^2 = x^1 + \alpha_1 r_1$

Relation b/w directions of $r_0$ & $r_1$:

$r_0^T r_1 = r_0^T [b - A x_1]$
$$= r_0^T [b - A(x^0 + \alpha_0 r_0)].$$
$$= r_0^T [\underbrace{(b - A x^0)}_{+r_0^0} - \alpha_0 A r_0]$$

$$= r_0^T r_0 \quad - \alpha_0 r_0^T A r_0.$$

$$= 0 \longrightarrow \{ \text{Using } ① \}$$

$\rightarrow$ $r_0$ and $r_1$ are orthogonal to each other.

$\Rightarrow$ we'll be moving mutually perpendicular directions till we reach the solution.

$r_0 = b - Ax^0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$r_0^T r_0 = \begin{bmatrix} 1 & 1 \end{bmatrix}\begin{bmatrix} 1 \\ 1 \end{bmatrix} = 2$

$r_0^T A r_0 = \begin{bmatrix} 1 & 1 \end{bmatrix}\begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}\begin{bmatrix} 1 \\ 1 \end{bmatrix} = 3$.

$\alpha_0 = \dfrac{r_0^T r_0}{r_0^T A r_0} = \dfrac{2}{3}$

$x^1 = x^0 + \alpha_0 r_0$

$\qquad = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \dfrac{2}{3}\begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2/3 \\ 2/3 \end{bmatrix}$

$r_1 = b - Ax^1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}\begin{bmatrix} 2/3 \\ 2/3 \end{bmatrix}$

$\downarrow$

$\alpha_1 \to \ x_2 \to \cdots$ until $r_k = 0$.

Substantial no. of steps required for the naïve method.

Improvement $\longrightarrow$ Conjugate Gradient method.

---

| Lecture 41: Conjugate Gradient Method |
| --- |

$P_0 = r_0$

$\qquad \qquad x_c \left(\frac{b_1}{2}\right)$

$P_1 = k(x_c - x')$

$\rightarrow$ Initial direction same as steepest descent method

Suppose we want to reach $x_c$ from $x'$ is 1 step. Then we'll have to move according to.

$AP_1 = k(Ax_c - Ax')$

$\qquad \qquad \underbrace{\qquad}_{b}$

$(\because x_c$ is correct sol$^n)$.

$\qquad = k(r_1)$

---

$r_0^T r_1 = 0$

Multiply LHS & RHS by $r_0^T$

$\rightarrow \quad r_0^T AP_1 = k r_0^T r_1 = 0$

$\Rightarrow P_0^T A P_1 = 0$

$P_0$ is 'A orthogonal' to $P_1$.

Form,

$P_1 = r_1 - \beta_1 P_0$

$\quad \hookrightarrow$ make new direction from old directions (Grame-Schmidt Conjecture).

$P_0^T A (r_1 - \beta_1 P_0) = 0$.

$\Rightarrow \beta_1 = \dfrac{P_0^T A r_1}{P_0^T A P_0}$.

We get direction $\overline{P_1}$ from this.
Next to find: how much to go?

$x^2 = x' + \alpha_1 P_1$.

In steepest descent, we used
$\qquad \qquad x^2 = x' + \alpha_1 r_1$

(we moved along $r_1$. Now we move along different direction $P_1$ with the 'hope of' reaching the target is 1 shot).

$\alpha_1$ should be such that $f$ should be a min.

$\Rightarrow f(x^2) = \dfrac{1}{2}(x' + \alpha_1 P_1)^T A (x' + \alpha_1 P_1)$

$\qquad \qquad - b^T(x' + \alpha_1 P_1) + c$

For $f$ to be minimum

$\dfrac{\partial f}{\partial \alpha_1} = 0 \longrightarrow$