

## System of Linear Equations

①

In the matrix notation, the system of  $n$  linear algebraic equations in  $n$ -unknowns:

$$\left. \begin{array}{l} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{array} \right\} \quad (1)$$

Can be written as:

$$AX = b \quad (2)$$

Where

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}_{n \times n}, \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1}, \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}_{n \times 1}$$

- The matrix  $[A|b]$  of order  $n \times (n+1)$  is called the augmented matrix.
- If all  $b_i$ 's are zero, then the system (1) or (2) is said to be homogeneous, and if at least one  $b_i$  is non-zero then it is in-homogeneous or non-homogeneous.
- System (2) has unique solution if and only if determinant of  $A$  (denoted by  $|A|$ ) is non-zero. And the solution is:

$$X = A^{-1}b$$

(2)

- The homogeneous system ( $b_i = 0, i = 1, 2, \dots, n$ ) possesses only a trivial solution, that is,  $x_1 = x_2 = \dots = x_n = 0$  (when  $|A| \neq 0$ )
- The methods for solving system of linear equations may broadly be classified into two types:
  - i) Direct Methods: These methods produce the exact solution after a finite number of steps.
  - ii) Iterative Methods: These methods give a sequence of approximate solutions, which converges to the exact solution when the number of steps tend to infinity.

### Direct Methods

- i) Cramer rule
  - ii) Gauss-elimination method
  - iii) Gauss-Jordan elimination method
- } — You are familiar with these from Mathematics-II course.

### Iterative Method

A general linear iterative method for the solution of the system (2) may be defined in the following form:

$$X^{(k+1)} = H X^{(k)} + C, \quad k = 0, 1, 2, \dots \quad (3)$$

where  $X^{(k+1)}$  and  $X^{(k)}$  are the approximations for  $X$  (exact sol.) at the  $(k+1)^{th}$  and  $k^{th}$  iterations, respectively.  $H$  is called iteration matrix depending on  $A$ ,  $C$  is a column vector.

When  $k \rightarrow \infty$ ,  $x^{(k)} \xrightarrow{(3)} x$  (exact sol.), which is given by

$$x = A^{-1}b \quad \text{--- (4)}$$

Therefore, from (3) :

$$A^{-1}b = HA^{-1}b + c \quad \text{--- (5)}$$

$\therefore$  the column vector  $c$  is given by :

$$c = (I - H)A^{-1}b \quad \text{--- (6)}$$

where  $I$  denotes identity matrix of order  $n \times n$ .

- We now determine the iteration matrix  $H$  and the column vector  $C$  for a few well known iteration methods:

### Jacobi Iteration Method :

Assume that the quantities  $a_{ii}$  ( $i=1, 2, \dots, n$ ) in (2) are pivot elements. The equation (2) may be re-written as:

$$\left. \begin{aligned} a_{11}x_1 &= -(a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n) + b_1 \\ a_{22}x_2 &= -(a_{21}x_1 + a_{23}x_3 + \dots + a_{2n}x_n) + b_2 \\ &\vdots \\ a_{nn}x_n &= -(a_{n1}x_1 + a_{n2}x_2 + \dots + a_{n,n-1}x_{n-1}) + b_n \end{aligned} \right\} \quad \text{--- (7)}$$

Note: pivot must always be non-zero.  
If not, can be made by row exchange (eq. exchange)

The Jacobi (also known as Gauss-Jacobi) iteration method may now be defined as :

(4)

$$\left. \begin{aligned} x_1^{(k+1)} &= -\frac{1}{a_{11}} (a_{12}x_2^{(k)} + a_{13}x_3^{(k)} + \dots + a_{1n}x_n^{(k)} - b_1) \\ x_2^{(k+1)} &= -\frac{1}{a_{22}} (a_{21}x_1^{(k)} + a_{23}x_3^{(k)} + \dots + a_{2n}x_n^{(k)} - b_2) \\ &\vdots \\ x_n^{(k+1)} &= -\frac{1}{a_{nn}} (a_{n1}x_1^{(k)} + a_{n2}x_2^{(k)} + \dots + a_{n,n-1}x_{n-1}^{(k)} - b_n) \end{aligned} \right\} \quad \text{--- (4)}$$

~~keep track of k~~  
k = 0, 1, 2, ...

- Since we replace the complete vector  $x^{(k)}$  in the right hand side of (4) at the end of each iteration. This method is also called as Method of Simultaneous Displacement.

- In the matrix form, Jacobi method can be represented as

$$\begin{aligned} x^{(k+1)} &= -D^{-1}(L+U)x^{(k)} + D^{-1}b \\ &= Hx^{(k)} + C \end{aligned}$$

where  $H = -D^{-1}(L+U)$  and  $C = D^{-1}b$

Here,  $L$  and  $U$  are respectively lower and upper triangular matrices (with zero diagonal). And  $D$  is the diagonal matrix such that

$$A = L + D + U$$

Note: Any <sup>square</sup> matrix can be decomposed as  $A = L + D + U$ :

$$\boxed{\begin{bmatrix} A \\ \end{bmatrix} = \begin{bmatrix} L \\ \end{bmatrix} + \begin{bmatrix} D \\ \end{bmatrix} + \begin{bmatrix} U \\ \end{bmatrix}}$$

Note: Eq. (4) is well-defined as each  $a_{ii}$  ( $i=1, 2, \dots, n$ ) is a pivot, hence non-zero.

— Here  $k=0, 1, 2, \dots$  represent the iteration index, and each value of  $k$  will generate one approximation of the solution  $x$

(5)

### Example 1: (Jacobi method)

Solve the system of equations:

$$\begin{bmatrix} 4x_1 + x_2 + x_3 = 2 \\ x_1 + 5x_2 + 2x_3 = -6 \\ x_1 + 2x_2 + 3x_3 = -4 \end{bmatrix}$$

Matrix notation:

$$\begin{bmatrix} 4 & 1 & 1 \\ 1 & 5 & 2 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ -6 \\ -4 \end{bmatrix}$$

A                    X                    b

$$L = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 2 & 0 \end{bmatrix}, \quad D = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 3 \end{bmatrix}, \quad U = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{bmatrix}$$

Iteration matrix:  $H = -D^{-1}(L+U)$

$$\begin{aligned} &= -\begin{bmatrix} 4 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 2 \\ 1 & 2 & 0 \end{bmatrix} = -\begin{bmatrix} y_4 & 0 & 0 \\ 0 & y_5 & 0 \\ 0 & 0 & y_3 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 2 \\ 1 & 2 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 & -\frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{5} & 0 & -\frac{2}{5} \\ -\frac{1}{3} & -\frac{2}{3} & 0 \end{bmatrix} \end{aligned}$$

$$\text{Column vector } c = D^{-1}b = \begin{bmatrix} y_4 & 0 & 0 \\ 0 & y_5 & 0 \\ 0 & 0 & y_3 \end{bmatrix} \begin{bmatrix} 2 \\ -6 \\ -4 \end{bmatrix} = \begin{bmatrix} y_2 \\ -6y_5 \\ -4y_3 \end{bmatrix}$$

$\therefore$  Jacobi iterative method becomes:

$$x^{(k+1)} = \begin{bmatrix} 0 & -\frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{5} & 0 & -\frac{2}{5} \\ -\frac{1}{3} & -\frac{2}{3} & 0 \end{bmatrix} x^{(k)} + \begin{bmatrix} y_2 \\ -6y_5 \\ -4y_3 \end{bmatrix}, \quad k = 0, 1, 2, \dots$$

————— (\*)

(6)

Note that to start the iterative process given by (4) we need  
 $x^{(0)}$  — known as initial guess.

$$x^{(0)} = \begin{bmatrix} x_1^{(0)} \\ x_2^{(0)} \\ x_3^{(0)} \end{bmatrix}$$

Choose  $x^{(0)} = \begin{bmatrix} 0.5 \\ -0.5 \\ -0.5 \end{bmatrix}$

Then using (4):

$$x^{(1)} = H x^{(0)} + C$$

$$\begin{bmatrix} 0 & -\frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{5} & 0 & -\frac{2}{5} \\ -\frac{1}{3} & -\frac{2}{3} & 0 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.5 \\ -0.5 \end{bmatrix} + \begin{bmatrix} \frac{1}{2} \\ -\frac{6}{5} \\ -\frac{4}{3} \end{bmatrix} = \begin{bmatrix} 0.7500 \\ -1.1000 \\ -1.1667 \end{bmatrix}$$

Using  $x^{(1)}$ , we can get  $x^{(2)}$  from (4):  $x^{(2)} = \begin{bmatrix} 1.0667 \\ -0.8033 \\ -0.8500 \end{bmatrix}$

Similarly, using  $x^{(2)}$ , we can get  $x^{(3)}$  from (4):  $x^{(3)} = \begin{bmatrix} 0.9300 \\ -1.0733 \\ -1.1000 \end{bmatrix}$

Likewise, we can continue to compute 4<sup>th</sup>, 5<sup>th</sup>, ... approximations to the solution.

Note: Exact sol<sup>n</sup> for this system is:  $x_1 = 1, x_2 = -1, x_3 = -1$

□

## Gauss-Seidel Iterative method

(1)

Gauss-Seidel method can be obtained from Jacobi method as:

- Use only on the right hand side of (8), all the available values from the present iteration. We write the Gauss-Seidel method as:

$$\left. \begin{aligned} x_1^{(k+1)} &= -\frac{1}{a_{11}} (a_{12}x_2^{(k)} + a_{13}x_3^{(k)} + \dots + a_{1n}x_n^{(k)}) + \frac{b_1}{a_{11}} \\ x_2^{(k+1)} &= -\frac{1}{a_{22}} (a_{21}x_1^{(k+1)} + a_{23}x_3^{(k)} + \dots + a_{2n}x_n^{(k)}) + \frac{b_2}{a_{22}} \\ &\vdots \\ x_n^{(k+1)} &= -\frac{1}{a_{nn}} (a_{n1}x_1^{(k+1)} + a_{n2}x_2^{(k+1)} + \dots + a_{n,n-1}x_{n-1}^{(k+1)}) + \frac{b_n}{a_{nn}} \end{aligned} \right\} \quad (10)$$

Which may be re-arranged in the form:

$$\left. \begin{aligned} a_{11}x_1^{(k+1)} &= -\sum_{i=2}^n a_{1i}x_i^{(k)} + b_1 \\ a_{21}x_1^{(k+1)} + a_{22}x_2^{(k+1)} &= -\sum_{i=3}^n a_{2i}x_i^{(k)} + b_2 \\ &\vdots \\ a_{n1}x_1^{(k+1)} + \dots + a_{nn}x_n^{(k+1)} &= b_n \end{aligned} \right\} \quad (11)$$

- Since we replace the vector  $x^{(k)}$  in the right hand side of (8) element by element, thus, this method is also called the method of successive displacement.

- In matrix notation, (11) becomes:

$$(D+L)x^{(k+1)} = -UX^{(k)} + b$$

(2)

$$\text{OR} \quad X^{(k+1)} = -(D+L)^{-1}U X^{(k)} + (D+L)^{-1}b \\ = H X^{(k)} + c, \quad k=0, 1, 2, \dots \quad \text{--- (12)}$$

where,  $H = -(D+L)^{-1}U$  and  $c = (D+L)^{-1}b$

Example 2: (Gauss-Seidel method) :

Solve the system of equations:

$$\begin{aligned} 2x_1 - x_2 + 0 \cdot x_3 &= 7 \\ -x_1 + 2x_2 - x_3 &= 1 \\ 0x_1 - x_2 + 2x_3 &= 1 \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \quad \begin{array}{l} \text{Use initial guess:} \\ X^{(0)} = [0, 0, 0]^T \end{array}$$

$$A = \begin{bmatrix} 2 & 1 & 0 \\ -1 & 2 & -1 \\ 0 & 1 & 2 \end{bmatrix}, \quad b = \begin{bmatrix} 7 \\ 1 \\ 1 \end{bmatrix}$$

$$(D+L) = \begin{bmatrix} 2 & 0 & 0 \\ -1 & 2 & 0 \\ 0 & 1 & 2 \end{bmatrix}, \quad U = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix}$$

Gauss-Seidel method is

$$X^{(k+1)} = H X^{(k)} + c, \quad \text{where } H = -(D+L)^{-1}U, \quad c = (D+L)^{-1}b \\ k=0, 1, 2, \dots$$

$$(D+L)^{-1} = \begin{bmatrix} 1/2 & 0 & 0 \\ 1/4 & 1/2 & 0 \\ 1/8 & 1/4 & 1/2 \end{bmatrix}, \quad (D+L)^{-1}U = \begin{bmatrix} 1/2 & 0 & 0 \\ 1/4 & 1/2 & 0 \\ 1/8 & 1/4 & 1/2 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & -1/2 & 0 \\ 0 & -1/4 & -1/2 \\ 0 & -1/8 & -1/4 \end{bmatrix}$$

$$(D+L)^{-1}b = \begin{bmatrix} 1/2 & 0 & 0 \\ 1/4 & 1/2 & 0 \\ 1/8 & 1/4 & 1/2 \end{bmatrix} \begin{bmatrix} 7 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 7/2 \\ 9/4 \\ 13/8 \end{bmatrix}$$

$$\therefore X^{(k+1)} = \begin{bmatrix} 0 & 1/2 & 0 \\ 0 & 1/4 & 1/2 \\ 0 & 1/8 & 1/2 \end{bmatrix} X^{(k)} + \begin{bmatrix} 7/2 \\ 9/4 \\ 13/8 \end{bmatrix}, \quad X^{(0)} = [0, 0, 0]^T \\ k=0, 1, 2, \dots$$

(3)

Hence, we can calculate approximations as :

$$x^{(1)} = \begin{bmatrix} 0 & 1/2 & 0 \\ 0 & 1/4 & 1/2 \\ 0 & 1/8 & 1/4 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 7/2 \\ 9/4 \\ 13/8 \end{bmatrix} = \begin{bmatrix} 3.5000 \\ 3.6250 \\ 2.3125 \end{bmatrix}$$

$$x^{(2)} = \dots = \begin{bmatrix} 4.6250 \\ 3.6250 \\ 2.3125 \end{bmatrix}$$

$$x^{(3)} = \dots = \begin{bmatrix} 5.3125 \\ 4.3125 \\ 2.6563 \end{bmatrix}$$

and so on.

Note: Exact sol. is  $x = [6, 5, 3]^T$

□

(1)

## System of Linear Equations (Contd.)

So far we have discussed about two numerical methods to solve system of linear equations of the form  $Ax=b$ . The basic idea of iterative method is to construct a sequence of vectors  $X^{(k)}$  that enjoy the property of convergence, that is:

$$X = \lim_{k \rightarrow \infty} X^{(k)} \quad \text{--- (1)}$$

where  $X$  is the exact solution to  $Ax=b$  and  $X^{(k)}$  is the  $k^{\text{th}}$  approximation of the exact solution  $X$  using an iterative method.

- The iterative process is stopped at the minimum value of  $n$  such that  $|X^{(n)} - X| < \epsilon$ .  $\text{--- (2)}$   
where  $\epsilon$  is a fixed tolerance (given by the user)
- However, in many engineering problems, exact sol. is obviously not available, it is necessary to introduce suitable stopping criteria to monitor the convergence of the iterative process. (This will be discussed later).
- Next, we will establish the convergence criteria for Jacobi and Gauss-Seidel methods.

Note: In case if two (or more) iterative methods converge, the betterness of an iterative method will be decided by the Rate of Convergence (that is how fast approximations approach to the exact sol. as in Eq. (1)).

(2)

Take the linear iterative process as:

$$X^{(k+1)} = HX^{(k)} + C, \quad k=0, 1, 2, \dots \quad \text{--- (3)}$$

for  $Ax=b$ . Here  $X$  is the exact solution of  $Ax=b$ .

$$\therefore X \text{ satisfy (3)} \Rightarrow X = HX + C \quad \text{--- (4)}$$

Subtracting (4) from (3) :

$$\Rightarrow X^{(k+1)} - X = H(X^{(k)} - X) \quad \text{--- (5)}$$

Substitute  $X^{(k)} - X = \epsilon^{(k)}$ , where  $\epsilon^{(k)}$  denotes error at  $k^{\text{th}}$  approximation. Equation (5) becomes:

$$\epsilon^{(k+1)} = H\epsilon^{(k)}, \quad k=0, 1, 2, \dots \quad \text{--- (6) (Recurrence relation)}$$

Note that:

$$\begin{aligned} \epsilon^{(k)} &= H\epsilon^{(k-1)} = H(H\epsilon^{(k-2)}) = H^2\epsilon^{(k-2)} = H^2(H\epsilon^{(k-3)}) \\ &= H^3\epsilon^{(k-3)} = \dots = H^k\epsilon^{(0)} \end{aligned}$$

Where  $\epsilon^{(0)} = X^{(0)} - X$ , that  $\epsilon^{(0)}$  denotes error in initial guess

$$\epsilon^{(k)} = H^k\epsilon^{(0)}, \quad k=0, 1, 2, \dots \quad \text{--- (7)}$$

Note: Here we assumed that the iteration matrix  $H$  remains constant for each iteration — STATIONARY method

From ⑦, it can be noted that:

(3)

$$\epsilon^{(k)} \rightarrow 0 \quad \text{if and only if } H^k \rightarrow 0 \quad \text{as } k \rightarrow \infty$$

- That is error in  $k^{\text{th}}$  approximation approaches to zero as  $k \rightarrow \infty$  if and only if  $H^k \rightarrow 0$  as  $k \rightarrow \infty$ .

Note:  $\epsilon^{(0)}$  is finite, as this is error in initial guess

- Since,  $H$  is a matrix and  $k$  is iteration index (+ve integer), therefore, it is not easy to compute  $H^k$  (Think of case when  $k$  is large). Hence, we will derive some other criteria to establish the convergence of linear iterative methods.

Before the convergence results, let's have few definitions:

### Spectral Radius

The largest eigenvalue (in magnitude) is called the spectral radius of a matrix  $A$ . Spectral radius of a matrix  $A$  is denoted by  $\rho(A)$ .

For example: If  $A_{3 \times 3}$  have eigenvalues as:  $\lambda_1 = 1, \lambda_2 = -5, \lambda_3 = 3$

then

$$\rho(A) = 5$$

### Diagonal Dominance:

If  $A = [a_{ij}]_{n \times n}$ , then  $A$  is diagonally dominant if  $|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, i=1, 2, \dots, n$

That is absolute value of the diagonal element in each row is greater than or equal to absolute sum of non-diagonal elements in that row

(4)

and if

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i=1, \dots, n$$

then A is said to be strictly diagonally dominant.

For example:  $A = \begin{bmatrix} 5 & -1 & 2 \\ 3 & -6 & -1 \\ 2 & -5 & 8 \end{bmatrix}$  is strictly diagonally dominant.

### Result 1 (Necessary and Sufficient Condition for Convergence)

A necessary and sufficient condition for convergence of an iterative method of the form:

$$X^{(k+1)} = H X^{(k)} + C, \quad k=0, 1, 2, \dots$$

is that the eigenvalues of the iteration matrix H satisfy

$$\rho(H) < 1$$

### Result 2 (Sufficient condition for the convergence of Jacobi method)

If A is strictly diagonally dominant, then the Jacobi iteration scheme converges for any choice of initial guess.

### Result 3 (Sufficient condition for the convergence of Gauss-Seidel method)

If A is strictly diagonally dominant, then the Gauss-Seidel iteration method converges for any choice of initial guess ( $X^{(0)}$ ).

NOTE: Each result have a very nice mathematical proof. (Not included here).

Rate of Convergence: The rate of convergence of an iterative (5) method is given by:

$$\gamma = -\log_{10}(\rho(H)) \text{ OR } \gamma = -\ln(\rho(H)) \quad \text{--- (8)}$$

where,  $\rho(H)$  is the spectral radius of  $H$ .

Note: In (8) we are using the final result. To derive (8) you need a mathematical proof!

### Remarks:

spectral radius

a) We have a relation between the iteration matrices of Jacobi and Gauss-Seidel methods as

$$\rho(H_{GS}) = [\rho(H_J)]^2 \rightarrow \boxed{\text{Note: I will prove this relation in the next week}}$$

Where

$H_{GS}$  - iteration matrix for Gauss-Seidel

$H_J$  - iteration matrix for Jacobi

Now if (say)  $\rho(H_J) = 0.5 \Rightarrow \rho(H_{GS}) = 0.25$  — Now relate this with Result 1 (or Eq. 8) you can conclude here that if both the methods converges, then rate of convergence of G-S method is more rapid than Jacobi method.

b) Next assume if  $\rho(H_J) = 1.5 \Rightarrow \rho(H_{GS}) = 2.25$  — Now using Result 1 you can conclude that if both the methods diverges (Result 1), then G-S method diverges more rapidly as compared to Jacobi method (Eq. 8).

□

### Successive Over-Relaxation Method

This method is a generalization of Gauss-Seidel method.

Gauss-Seidel method for  $Ax = b$  is given as:

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right] \quad \text{--- (1)}$$

$i = 1, 2, \dots, n$

Successive relaxation method is obtained from (1) by introducing a relaxation parameter  $\omega$  as:

$$x_i^{(k+1)} = \frac{\omega}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right] + (1-\omega) x_i^{(k)} \quad \text{--- (2)}$$

$i = 1, 2, \dots, n$

Equation (2) can be written in vector/matrix form as:

$$\cancel{x_i^{(k+1)}} \quad X^{(k+1)} = \omega \left[ D^{-1} b - D^{-1} L X^{(k+1)} - D^{-1} U X^{(k)} \right] + (1-\omega) I X^{(k)}$$

$$\text{OR, } (I + \omega D^{-1} L) X^{(k+1)} = [(1-\omega) I - \omega D^{-1} U] X^{(k)} + \omega D^{-1} b$$

$$\text{OR, } \boxed{X^{(k+1)} = (D + \omega L)^{-1} [(1-\omega) D - \omega U] X^{(k)} + \omega (D + \omega L)^{-1} b} \quad \text{--- (3)}$$

(2)

Now, comparing ③ with general iterative form that is

$$X^{(k+1)} = H X^{(k)} + C, \quad k=0,1,2,\dots$$

For Successive Relaxation method we have:

$$H_{SR} = (D + \omega L)^{-1} [(1-\omega)D - \omega U]$$

$$C = \omega(D + \omega L)^{-1} b$$

This method is consistent for any  $\omega \neq 0$  and for  $\omega=1$  it coincides with Gauss-Seidel method. In particular, if  $0 < \omega < 1$  the method is called under-relaxation, while if  $\omega > 1$  it is called over-relaxation.

Q: How to determine the <sup>optimum</sup> value  $\omega_{opt}$  for which the convergence rate is the highest?

To answer this question, consider the iteration matrices of Jacobi and Seidel methods Successive relaxation methods:

$$H_J = -D^{-1}(L+U)$$

$$H_{SR} = (D + \omega L)^{-1} [(1-\omega)D - \omega U]$$

Let  $\mu$  be an eigenvalue of  $H_J$  and  $\lambda$  be an eigenvalue of  $H_{SR}$ .

(3)

The relation between  $\lambda$  and  $\mu$  is :

$$\mu = \frac{\lambda + \omega - 1}{\lambda^{\frac{1}{2}} \omega} \quad \text{--- (4)}$$

$$\Rightarrow \lambda - \mu \omega \lambda^{\frac{1}{2}} + (\omega - 1) = 0 \quad (\text{which is quadratic in } \lambda^{\frac{1}{2}})$$

$$\therefore \lambda^{\frac{1}{2}} = \frac{1}{2} [\mu \omega \pm \sqrt{\mu^2 \omega^2 - 4(\omega - 1)}]$$

$$= \frac{1}{2} [\mu \omega] \pm \frac{\mu}{2} \sqrt{(\omega - \omega_1)(\omega - \omega_2)}$$

$$\text{where } \omega_{1,2} = \frac{2}{\mu^2} \left[ 1 \mp \sqrt{1 - \mu^2} \right]$$

When  $\omega = \omega_1$ ,  $\rho(H_{SR})$  is smallest (implying fastest convergence), hence the optimal relaxation parameter is given as:

$$\omega_{opt} = \frac{2}{\mu^2} \left[ 1 - \sqrt{1 - \mu^2} \right] = \frac{2}{1 + \sqrt{1 - \mu^2}} \quad \text{--- (5)}$$

Also, when  $\omega = \omega_1$ , we get

$$\lambda = \frac{1}{4} \mu^2 \omega^2 = \omega - 1$$

For convergence, we require

$$|\lambda| < 1 \Rightarrow |\omega - 1| < 1$$

$$\text{OR} \quad 0 < \omega < 2$$

For  $0 < \omega < 1$  — Under Relaxation.

And  $1 < \omega < 2$  — Over Relaxation.

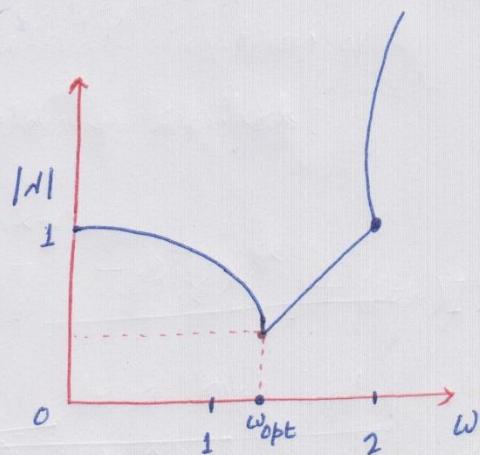


Fig:  $|λ|$  Vs.  $\omega$

(4)

The rate of convergence for SR method is :

$$\gamma_{SR} = -\ln(\omega_{opt} - 1) \quad \text{--- (6)}$$

Note that if  $\omega = 1$  (where SR becomes GS) from (4), we have :

$$\mu = \lambda^{1/2} \quad \text{OR} \quad \lambda = \mu^2$$

Where  $\mu$  is the arbitrary eigenvalue of  $H_J$  and  $\lambda$  is the eigenvalue of  $H_{GS}$ .

$$\Rightarrow \rho(H_{GS}) = [\rho(H_J)]^2 \quad \xrightarrow{\text{This is the same equation as given in Remarks of Lecture-2}}$$

--- (7)

∴ We can derive (7) if we can derive (4).

Also note from (5) that  $\omega_{opt}$  is real if  $|\mu| < 1$ . Therefore, the necessary condition for the SR method to converge is that the corresponding Jacobi method is convergent.

Example 1: (To find  $\omega_{opt}$ )

$$\begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \quad \text{--- (*)}$$

Q. Find  $\omega_{opt}$ .

Ans: Use  $\omega_{opt} = \frac{2}{1 + \sqrt{1 - \mu^2}}$ , where  $\mu$  is the eigenvalue/spectral radius of Jacobi method for system (\*)

Example 2: (SR method)

(5)

$$\begin{bmatrix} -2 & -1 & 0 \\ 1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ 1 \\ 1 \end{bmatrix} \quad \text{--- (**)}$$

Set up the SR method. Find optimal relaxation parameter ( $w_{opt}$ ) and the rate of convergence. Perform three iterations of the SR method. Take initial guess as  $x^{(0)} = [0, 0, 0]^T$

Ans: To find  $w_{opt}$ : First find  $H_J$  and then get spectral radius of  $H_J$ , say  $\mu$ .

$$\text{Then } w_{opt} = \frac{2}{\mu^2} (1 - \sqrt{1 - \mu^2})$$

To find rate of convergence: when  $w = w_{opt}$

Eigenvalues of  
 $H_J: \lambda_J = 0, \pm \frac{1}{\sqrt{2}}$   
 Take  $\mu = +\frac{1}{\sqrt{2}}$   
 ↓  
 Largest (in magnitude)

$$\rho(H_{SR}) = w_{opt}^{-1} = 0.1715 \quad (\text{check it!})$$

$$\text{Rate of convergence } \gamma_{SR} = -\ln(0.1715) = \dots$$

To start the iterations: Put  $w = w_{opt} = 0.1715$  in (3) (SR iterations)

and get

$$\begin{cases} x^{(k+1)} = H_{SR} x^{(k)} + c, \quad k=0, 1, 2, \dots \\ H_{SR} = (D + wL)^{-1} [(1-w)D - wU] \\ c = w(D + wL)^{-1} b \end{cases}$$

Initial guess (for  $k=0$ ) is given.

- Compute  $x^{(1)}, x^{(2)}$  and  $x^{(3)}$  (Three iterations)
- Complete it. Exercise. □

## System of Non-Linear Equations

①

So far, we have discussed about the numerical solution of system of linear equations. Now, we will discuss about the numerical solution to system of non-linear equations.

- As solving non-linear equations is much more challenging as compared to the linear equations.
- Therefore, first we will consider a single non-linear equation, then slowly will move towards system of non-linear equations (that is two or more number of equations)
- Examples of non-linear equations (Transcendental and Polynomials)
  - i)  $y = f(x) = \cos x - xe^{-x}$
  - ii)  $y = f(x) = x^3 - 5x + 1 \square$
- We will consider the real valued functions, that is,  
 $f : \mathbb{R} \rightarrow \mathbb{R}$ , where  $\mathbb{R} = (-\infty, \infty)$  — set of real numbers.
- Solving a non-linear equation means we have to find a scalar  $x$  such that  $f(x) = 0$  (That is  $x$  satisfies  $f(x)$ ). Then  $x$  is called solution or zero or root of  $f(x)$ .
- Methods for the numerical approximation of a zero of  $f(x)$  are usually iterative. Thus the aim is to generate a

sequence of approximations (similar to linear case)  $\{x^{(k)}\}_{k=0}^{\infty}$   
 such that

$$\boxed{\lim_{k \rightarrow \infty} x^{(k)} = \alpha}, \quad \text{where } \alpha \text{ is the exact sol. of } f(x).$$

— ①

### Geometric Approach to Rootfinding (of a non-linear function)

Here, we introduce following methods for finding roots:

- i) The bisection method
- ii) The chord method
- iii) The secant method
- iv) The false position (or Regula Falsi) method
- v) The Newton's method.

Growing complexity of the algorithms.

#### The Bisection method:

In the bisection method, the only information that is being used is the sign of the function  $f(x)$  at the end of any bisection (sub-interval).

The bisection method is based on the following property:

#### Property 1 (IVT Theorem)

Note: IVT - Intermediate

Value Theorem

- You know it from

Given a continuous function

$f: [a, b] \rightarrow \mathbb{R}$ , such that  $f(a)f(b) < 0$   
 then there exists  $\alpha \in (a, b)$  (That is  $a < \alpha < b$ )  
 such that  $f(\alpha) = 0$ .

(3)

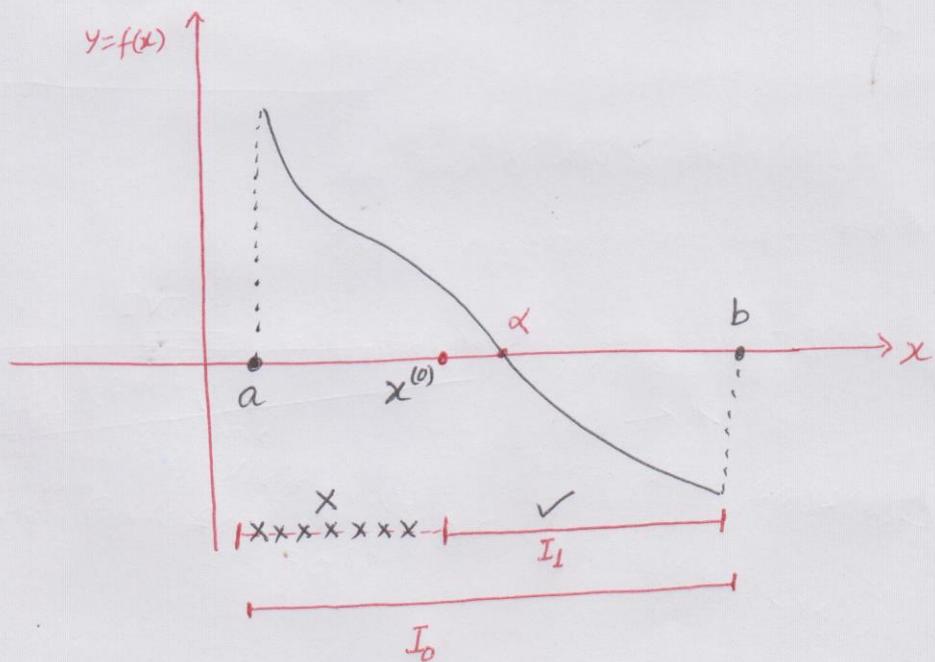
Starting from  $I_0 = [a, b]$ , the bisection method generates a sequence of subintervals  $I_k = [a^{(k)}, b^{(k)}]$ ,  $k \geq 0$  with

$$I_k \subset I_{k-1}, k \geq 1$$

and enjoys the property that  $f(a^{(k)})f(b^{(k)}) < 0$ .

Precisely, we set  $a^{(0)} = a$ ,  $b^{(0)} = b$  and  $x^{(0)} = (a^{(0)} + b^{(0)})/2$ , then, for  $k \geq 0$ :

$$\left. \begin{array}{l} \text{set } a^{(k+1)} = a^{(k)}, b^{(k+1)} = x^{(k)}, \text{ if } f(x^{(k)})f(a^{(k)}) < 0 \\ \text{set } a^{(k+1)} = x^{(k)}, b^{(k+1)} = b^{(k)}, \text{ if } f(x^{(k)})f(b^{(k)}) < 0 \\ \text{finally, set } x^{(k+1)} = (a^{(k+1)} + b^{(k+1)})/2 \end{array} \right\} -②$$



- The bisection iteration terminates at the  $m^{\text{th}}$  step  
for which

$$|x^{(m)} - \alpha| \leq |I_m| \leq \epsilon$$

where  $\epsilon$  is a fixed tolerance and  $|I_m|$  is the length of  $I_m$ . As for the speed of convergence of the bisection method, note that  $|I_0| = b - a$ , while

$$|I_k| = I_0 / 2^k = \frac{(b-a)}{2^k}, \quad k \geq 0 \quad \text{--- (3)}$$

Denoting  $e^{(k)} = |x^{(k)} - \alpha|$  — Absolute error at  $k^{\text{th}}$  iteration, from (3) it follows that

$$|e^{(k)}| \leq \frac{(b-a)}{2^k}, \quad k \geq 0 \quad \text{which implies}$$

$$\lim_{k \rightarrow \infty} |e^{(k)}| = 0 \iff x^{(k)} \rightarrow \alpha \text{ as } k \rightarrow \infty.$$

∴ The bisection method is convergent. To get the value of  $m$  (for which  $|x^{(m)} - \alpha| \leq \epsilon$ , we must take

$$m \geq \log_2(b-a) - \log_2(\epsilon) = \frac{\log((b-a)/\epsilon)}{\log(2)} \underset{\text{--- (4)}}{\approx} \frac{\log((b-a)/\epsilon)}{0.6931}$$

Since  $m$  is an (+ve) integer, we take  $m$  as the next nearest integer.

(5)

- This singles out the bisection method. The (only) drawback of this method is its slow convergence.
- The Methods of Chord, Secant and Regula Falsi and Newton's Method:
- To devise algorithms with better convergence properties than the bisection method, it is necessary to include information from the values attained by  $f$  and, possibly, also by its derivative  $f'$  (if  $f$  is differentiable) or by a suitable approximation.
- Let us expand  $f$  in a Taylor Series around  $\alpha$  and truncate the expansion at the first order, as:

$$f(\alpha) = 0 = f(x) + (\alpha - x)f'(\xi) \quad \text{--- (5)}$$

where  $\xi$  lies between  $\alpha$  and  $x$ .

- Equation (5) prompts the following iterative method:

For any  $k \geq 0$ , given  $x^{(k)}$ , determine  $x^{(k+1)}$  by solving equation

$$\boxed{f(x^{(k)}) + (x^{(k+1)} - x^{(k)}) q_k = 0}$$

where  $q_k$  is a suitable approximation of  $f'(x^{(k)})$ .

- ⑥
- The methods described here amounts to find the intersection between  $x$ -axis and the straight line of slope  $q_k$  passing through the point  $(x^{(k)}, f(x^{(k)}))$ .
  - Therefore, we have:

$$x^{(k+1)} = x^{(k)} - q_k^{-1} f(x^{(k)}), \quad k=0, 1, 2, \dots$$

Q: What are the possibilities for  $q_k$ ?

- We consider four particular choices of  $q_k$ .

#### A) The Chord Method:

Here, we let

$$q = q_k = \frac{f(b) - f(a)}{b - a}, \quad \text{for all } k \geq 0$$

we have

$$x^{(k+1)} = x^{(k)} - \left( \frac{b-a}{f(b)-f(a)} \right) f(x^{(k)}), \quad k=0, 1, 2, \dots$$

— ⑥

Note that to start ⑥, one initial guess  $x^{(0)}$  is required.

#### B) The Secant Method

In this case, take

$$q_k = \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}, \quad k=0, 1, 2, \dots \quad \text{--- ⑦}$$

(7)

∴ we have

$$x^{(k+1)} = x^{(k)} - \left( \frac{x^{(k)} - x^{(k-1)}}{f(x^{(k)}) - f(x^{(k-1)})} \right) f(x^{(k)}), \quad k=0, 1, 2, \dots$$

— (8)

Note that to start iteration with (8), we need two initial guesses.

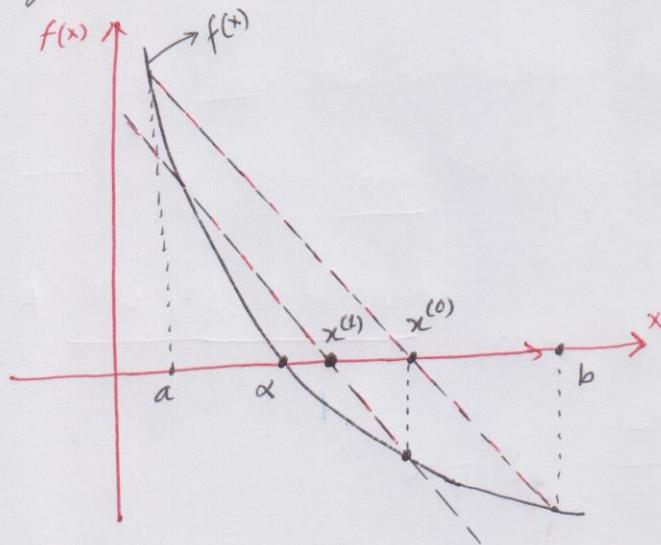


Fig1: Chord Method

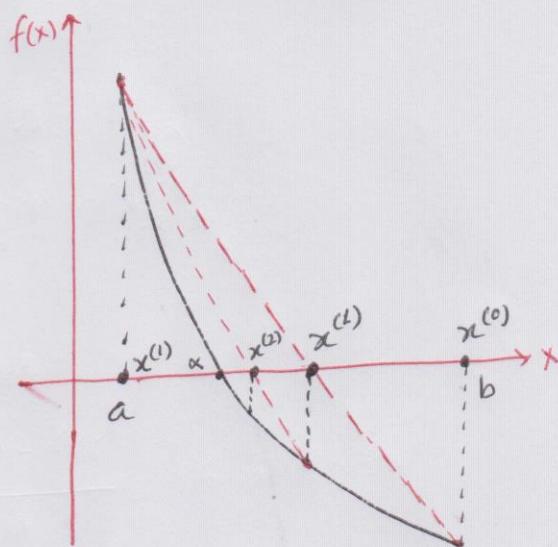


Fig2: Secant Method

(Here we need two initial guesses. Take guesses as :

$$x^{(0)} = b \quad x^{(1)} = a$$

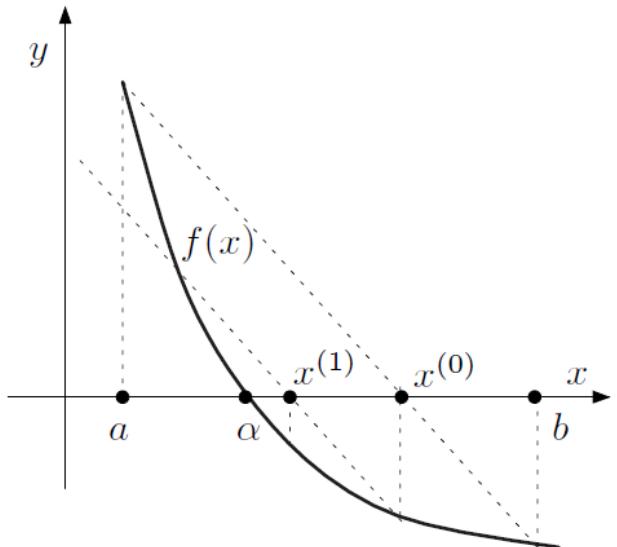


Fig. 1: Chord Method

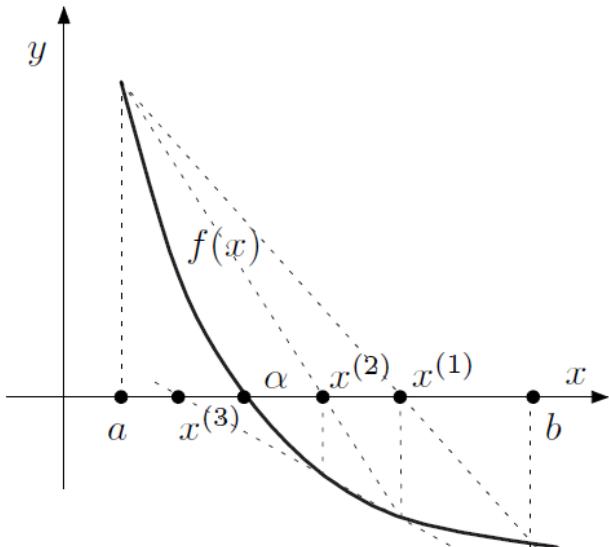


Fig. 2: Secant Methods

## System of Non-Linear Equations

Lec. Notes-05

①

### The Regula Falsi (or false position) method

- Variant of the Secant method.
- Here, instead of taking the secant line through the points  $(x^{(k)}, f(x^{(k)}))$  and  $(x^{(k-1)}, f(x^{(k-1)}))$ , we take the one through  $(x^{(k)}, f(x^{(k)}))$  and  $(x^{(k')}, f(x^{(k')}))$ .

Where  $k'$  denotes the maximum index less than  $k$  such that

$$f(x^{(k)}) \cdot f(x^{(k')}) < 0.$$

Therefore, Once two initial guesses  $x^{(0)}$  and  $x^{(1)}$  have been found such that

$$f(x^{(0)}) \cdot f(x^{(1)}) < 0,$$

we let

$$x^{(k+1)} = x^{(k)} - \left[ \frac{x^{(k)} - x^{(k')}}{f(x^{(k)}) - f(x^{(k')})} \right] f(x^{(k)}), \quad k \geq 0$$

- If error tolerance ( $\epsilon$ ) is given, then the iterative process ① terminates at the  $m^{\text{th}}$ -step such that

$$|f(x^{(m)})| < \epsilon$$

(2)

- Regula Falsi method has the same complexity as the Secant method
- However, unlike the Secant method, the iterates generated by Regula falsi are all contained within the starting interval  $[x^{(0)}, x^{(1)}]$ .

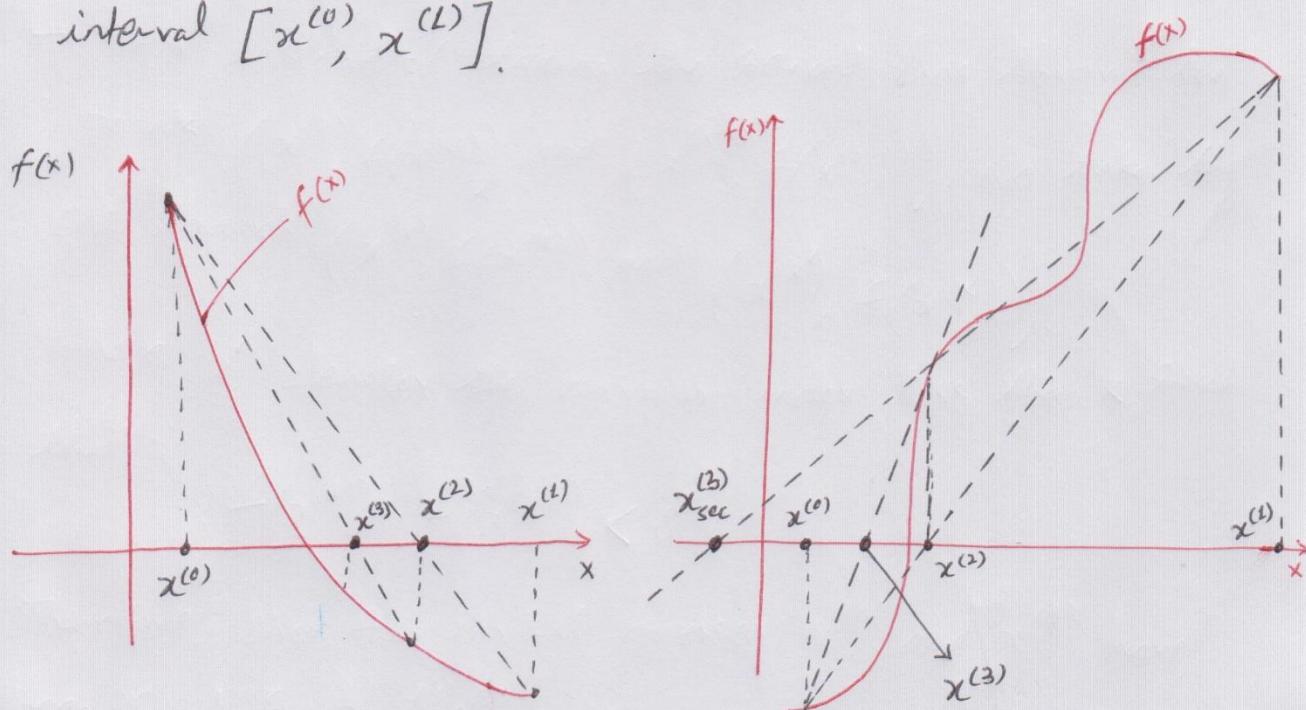


Fig 1: First two step of the Regula Falsi method for two different functions

- On the right side frame, we can notice that initial guesses are chosen as:  $x^{(0)}$  and  $x^{(1)}$ .
- Note that  $x^{(2)}$  computed by the Secant method coincides with Regula Falsi method, however,  $x^{(3)}$  computed by secant (that is  $x_{\text{sec}}^{(3)}$ ) falls outside the searching interval  $[x^{(0)}, x^{(1)}]$ .

(3)

Note:

Regula Falsi and bisection methods can be regarded as globally convergent method.

### Newton's (Newton-Raphson) method :

Assuming that  $f$  is continuously differentiable on a interval  $I$  (that is  $f \in C^1(I)$ ) and  $f'(x) \neq 0$  (that is  $\alpha$  is a simple root of  $f$ ).

If we take

$$g_k = f'(x^{(k)}), \quad \text{for all } k \geq 0$$

and if initial guess  $x^{(0)}$  is given. Then the Newton's  
(or N-R) method is :

$$\boxed{x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k \geq 0} \quad (2)$$

- At  $k^{\text{th}}$  iteration, Newton's method requires two functional evaluations :  $f(x^{(k)})$  and  $f'(x^{(k)})$ .
- (You will see later) that increasing computational cost will be compensated for by a higher order of convergence.

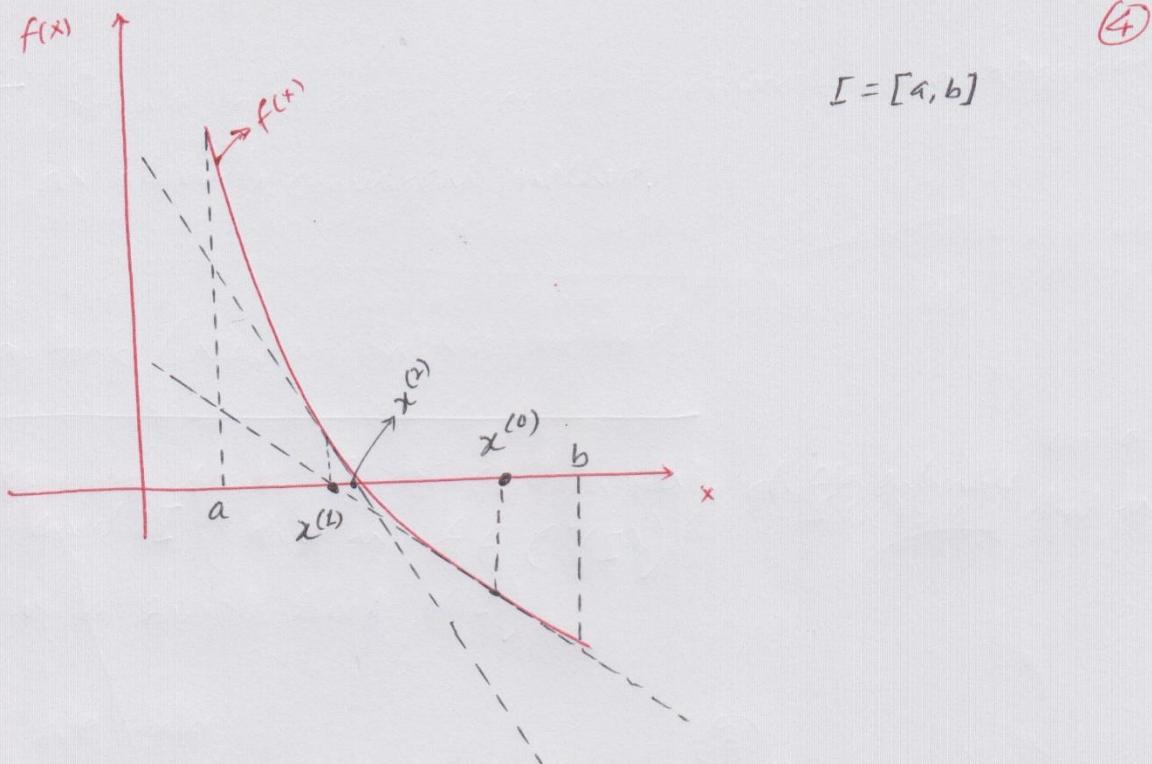


Fig 2: Newton's method.

Rate of Convergence (of iterative methods to find the root of a non-linear equation):

— An iterative method is said to be of order  $p$  (or has the rate of convergence  $p$ ) if  $p$  is the largest positive real number for which there exists a finite constant  $C \neq 0$  such that

$$|\varepsilon_{k+1}| \leq C |\varepsilon_k|^p$$

where  $\varepsilon_k = x_k - \alpha$  is the error in the  $k^{\text{th}}$  approximation,

Here  $\alpha$  is the exact root of  $f(x) = 0$ .

## Rate of convergence of Newton's method:

(5)

Newton's Method

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k=0, 1, 2, 3, \dots \quad \text{--- (3)}$$

On substituting  $x^{(k)} = \alpha + \varepsilon_k$ , we get from (3):

$$\varepsilon_{k+1} = \varepsilon_k - \frac{f(\alpha + \varepsilon_k)}{f'(\alpha + \varepsilon_k)}$$

Using Taylor series expansions for  $f$  and  $f'$ :

$$= \varepsilon_k - \frac{(f(\alpha) + \overset{\nearrow^0}{\varepsilon_k} f'(\alpha) + \frac{\varepsilon_k^2}{2!} f''(\alpha) + \dots)}{(f'(\alpha) + \varepsilon_k f''(\alpha) + \dots)}$$

$$= \varepsilon_k - \left[ \varepsilon_k + \frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)} \varepsilon_k^2 + \dots \right] \left[ 1 + \frac{f''(\alpha)}{f'(\alpha)} \varepsilon_k + \dots \right]^{-1}$$

$\downarrow$

[Use Binomial expansion  
 $(1+x)^n = \sum_{k=0}^n \binom{n}{k} x^k$ ]

OR

$$\varepsilon_{k+1} = \frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)} \varepsilon_k^2 + O(\varepsilon_k^3)$$

On neglecting  $\epsilon_k^3$  and higher powers of  $\epsilon_k$ , we get

(6)

$$\epsilon_{k+1} = C \epsilon_k^2 \quad \dots \quad (4)$$

where  $C = \frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)}$

Thus (using (4)) we conclude that Newton's method has quadratic convergence (or second order convergence)  $\square$

### Exercise

1. Show that Regula Falsi method has linear rate of convergence (that is of order 1)
2. Show that Secant Method converges to  $\alpha$  with order  $p = (1 + \sqrt{5})/2 \approx 1.63$

### Note: (Summary of order of convergence)

- Chord and Regula - Falsi have linear convergence ( $p=1$ )
- Secant has convergence of order  $p=1.63$  (larger than Chord/RF and smaller than N-R)
- Newton's Method has quadratic convergence ( $p=2$ ).  $\square$

## Multiple Roots

Lec. Notes-6a

(1)

Let's consider a (non-linear) function  $f(x)$

- To find the root of  $f(x)$ , write  $f(x)=0$  — (1)
- If we can write (1) as

$$f(x) = (x-\alpha)^m g(x) = 0 \quad (g(x) \text{ bounded})$$

such that  $g(\alpha) \neq 0$ , then  $\alpha$  is called a multiple root of multiplicity  $m$ . In this case:

$$f\Big|_{x=\alpha} = f'\Big|_{x=\alpha} = \dots = f^{(m-1)}\Big|_{x=\alpha} = 0$$

If  $m=1$ , the number  $\alpha$  is said to be simple root.

Example:  $f(x) = x^3 - 5x^2 + 7x - 3 = (x-3)(x-1)^2$

So, here  $x=1$  is the double root of  $f(x)$ .

and  $x=3$  is a simple root of  $f(x)$ .

- Multiple roots pose some difficulties for numerical method (many)
  - i) The function does not change sign at multiple roots!
    - Can not use root bracketing methods (like bisection and Regula Falsi)
  - ii) Another possible problem is related to the fact that not only  $f(x)$  but also  $f'(x)$  goes to zero at the root

- This poses problems for both the Newton's and Secant methods (which both contain the derivative or its estimate in the denominator of their respective formulas)

iii) A very simple way to circumvent these problems is based on the fact that it can be demonstrated theoretically that  $f(x)$  will always reach zero before  $f'(x)$ . Therefore, if a zero check for  $f(x)$  is incorporated, then the computation can be terminated before  $f'(x)$  reaches to zero.

### How to alleviate these problems?

- Following is the modification proposed for Newton's method:

$$x^{(k+1)} = x^{(k)} - m \frac{f(x^{(k)})}{f'(x^{(k)})}, k \geq 0 \quad \text{--- (2)}$$

Where  $m$  is the multiplicity of the root.

- The drawback for alternative (2) is that it hinges on forknowledge of the multiplicity of the root!

- Another alternative, is to define a new function  $u(x)$  such that-

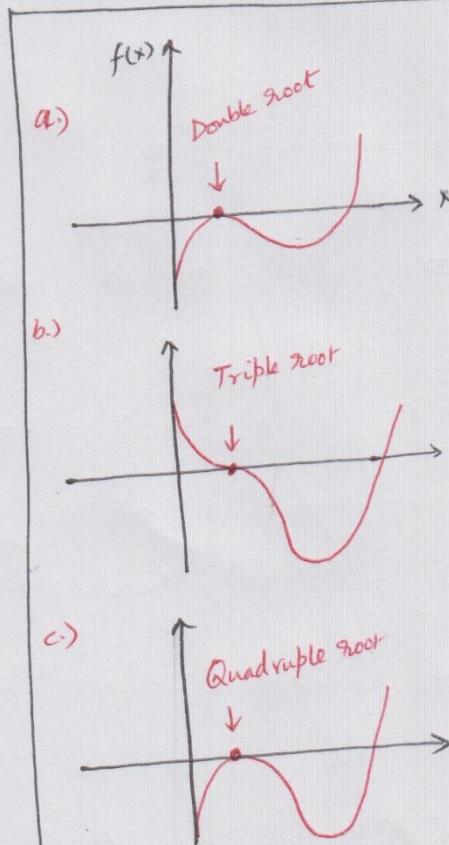


Fig: Double, Triple and Quadruple roots.

(3)

$$u(x) = \frac{f(x)}{f'(x)} \quad \text{--- (3)}$$

— It can be shown that  $u(x)$  has roots at all the same locations as the  $f(x)$ .

Substitute (3) into Newton's iterative formula for simple roots  
 (which was  $x^{(k+1)} = x^{(k)} - f(x^{(k)})/f'(x^{(k)})$ )

$$\Rightarrow x^{(k+1)} = x^{(k)} - \frac{u(x^{(k)})}{u'(x^{(k)})} \quad \text{--- (4)}$$

Differentiate (3)  $\Rightarrow u'(x) = \frac{f'(x)f'(x) - f(x)f''(x)}{[f'(x)]^2} \quad \text{--- (5)}$

Substitute (5) and (3) into (4): we have

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)}) f'(x^{(k)})}{[f'(x^{(k)})]^2 - f(x^{(k)}) f''(x^{(k)})}, \quad k \geq 0$$

(6)

↳ Modified Newton's Method for Multiple Roots.

### Exercise:

Use modified Newton's method to evaluate the multiple root of  $f(x) = x^3 - 5x^2 + 7x + 3$ , with an initial guess  $x^{(0)} = 1$ .

(Perform THREE iterations using (6))



# Accuracy and Precision

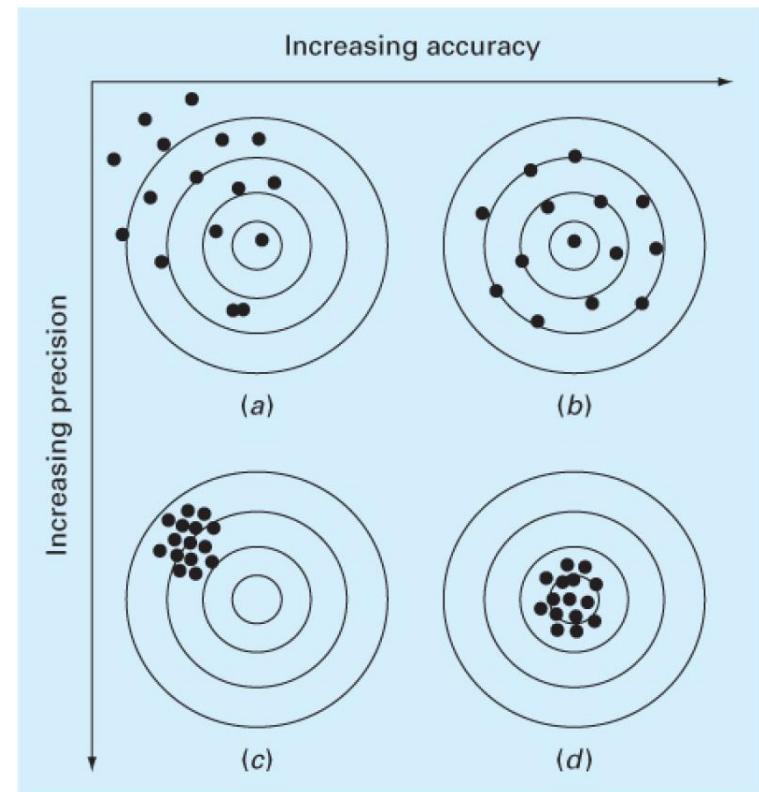
## Accuracy:

Accuracy refers to how closely a computed or measured value agrees with the true value.

## Precision:

Precision refers to how closely individual computed or measured values agree with each other.

- a) inaccurate and imprecise
- b) accurate and imprecise
- c) inaccurate and precise
- d) accurate and precise



## Error Definitions (1/2)

Numerical errors arise due to use of approximations to represent exact mathematical operations and quantitates. In general, these include:

### Truncation errors and Round-off errors.

- **Truncation errors:** arise when approximations are used to represent exact mathematical procedure.
- **Round-off errors:** arise when numbers having limited significant figures are used to represent exact numbers.

For both the cases, we have

$$\text{True value} = \text{approximation} + \text{error}$$

**True error ( $E_t$ ):** the difference between the true value and the approximation

$$E_t = \text{true value} - \text{approximation} \quad \text{----- (1)}$$

## Error Definitions (1/2)

**Absolute error** ( $|E_t|$ ): the absolute difference between the true value and the approximation.

**Note:** the drawback of  $(E_t)$  is that it takes no account of the order of magnitude of the value under examination.

**True fractional relative error:** the true error divided by the true value

$$\text{True fractional relative error} = \frac{\text{true value} - \text{approximation}}{\text{true value}}$$

**Relative error** ( $\varepsilon_t$ ): the true fractional relative error expressed as a percentage

$$\varepsilon_t = \frac{\text{true value} - \text{approximation}}{\text{true value}} \times 100 \% \quad ----- (2)$$

## Error Definitions (1/2)

**Example-1:** Suppose that you have the task of measuring the lengths of a bridge and a rivet and come up with 9999 and 9 cm, respectively. If the true values are 10,000 and 10 cm, respectively, compute (i) the true error and (ii) the true percent relative error for each case.

**Solution:**  $E_t = 10,000 - 9999 = 1 \text{ cm}$  (for bridge)

$$E_t = 10 - 9 = 1 \text{ cm}$$
 (for rivet)

(i) The percent relative error for bridge is:

$$\varepsilon_t = \frac{1}{10000} 100\% = 0.01\%$$

(ii) The percent relative error for rivet is:

$$\varepsilon_t = \frac{1}{10} 100\% = 10\%$$

## Error Definitions (2/2)

The previous definitions of error relied on knowing a true value. If that is not the case, approximations can be made to the error as:

The **approximate percent relative error** can be given as the approximate error divided by the approximation, expressed as a percentage - though this presents the challenge of finding the approximate error!

$$\varepsilon_a = \frac{\text{approximate error}}{\text{approximation}} \times 100 \%$$

For **iterative processes**, the error can be approximated as the difference in values between successive iterations

$$\varepsilon_a = \frac{\text{present approximate} - \text{previous approximation}}{\text{present approximation}} \times 100 \% \quad \text{--- (4)}$$

## Error Tolerance and Error Estimates

Often, when performing calculations, we may not be concerned with the sign of the error but are interested in whether **the absolute value of the percent relative error** is lower than a **pre-specified tolerance ( $\varepsilon_s$ )**.

For such cases, the computation is repeated until  $|\varepsilon_a| < \varepsilon_s$

This relationship is referred to as a **stopping criterion**

Note that for the remainder of our discussions, we almost always employ absolute values when using relative errors

We say that an approximation is correct to **at least  $n$  significant figures (significant digits)** if its  $|\varepsilon_a|$  is smaller than  $\varepsilon_s$  that has a value

$$\varepsilon_s = (0.5 \times 10^{2-n})\% \quad \text{----- (5)}$$

**Example-2:** It is known that the exponential function can be computed using Maclaurin series, as:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots + \frac{x^n}{n!}$$

Starting with the simplest version,  $e^x = 1$ , add terms one at a time to estimate  $e^{0.5}$ . Note that the true value is  $e^{0.5} = 1.648721\dots$ . Add terms until the absolute value of the  $\varepsilon_a$  falls below a prescribed  $\varepsilon_s$  conforming to three significant figures.

**Solution:**

Using Eq. (5) determine the error criterion that ensures a result is correct to at least three significant figures:

$$\varepsilon_s = (0.5 \times 10^{-3}) \% = 0.05\%$$

Thus, we will add terms to the series until  $\varepsilon_a$  falls below this level.

<b>Terms</b>	<b>Result</b>	$\epsilon_t, \%$	$\epsilon_a, \%$
1	1	39.3	
2	1.5	9.02	33.3
3	1.625	1.44	7.69
4	1.645833333	0.175	1.27
5	1.648437500	0.0172	0.158
6	<u>1.648697917</u>	0.00142	<u>0.0158</u>

Thus, after six terms are included, the approximate error falls below  $\epsilon_s = 0.05\%$ , and computations are terminated.

## Round-off Errors

### Roundoff Errors:

Roundoff errors arise because digital computers cannot represent some quantities exactly. There are two major facets of roundoff errors involved in numerical calculations:

- i. Digital computers have size and precision limits on their ability to represent numbers.
- ii. Certain numerical manipulations are highly sensitive to roundoff errors.

## Computer Representation of Numbers (1/2)

**Number System:** A convention used for representing quantities.

We are familiar with **decimal**, or **base-10**, number system. The base-10 system uses the 10 digits:

0, 1, 2, 3, 4, 5, 6, 7, 8, 9

to represent a number.

Numbers on the computer are represented with a **binary**, or **base-2**, system (having just 0 and 1)

**Floating-Point Representation:** Fractional quantities are typically represented in computers using floating point form. Here, the number is expressed as a fraction part, called a **mantissa** or **significand**, and an integer part, called **exponent** or **characteristic**, as an

## Computer Representation of Numbers (2/2)

$$m \cdot b^e$$

Where  $m$  = the mantissa,  $b$  = the base of the number system being used, and  $e$  = the exponent. For example, 156.78 can be represented as  $0.15678 \times 10^3$  in a floating-point base-10 system.

Note that mantissa is usually normalized if it has leading zero digits. For example, suppose the quantity  $1/34 = 0.029411765\dots$  was stored in a floating-point base-10 system that allowed only four decimal places to be stored as

$$0.0294 \times 10^0$$

Which can be normalized as,  $0.2941 \times 10^{-1}$ , thus we retain an additional significant figure when the number is stored.

## Extended Precision

The most commonly used precision for Computers are:

- i. Single precision
  - ii. Double precision
- 

### Single Precision:

Computers that use IEEE format allow 24 bits to be used for the mantissa, which translates into about seven significant base-10 digits of precision with a range of about  $10^{-38}$  to  $10^{39}$ .

### Double Precision:

Here, about fifteen to sixteen significant base-10 digits of precision with a range of approximately  $10^{-308}$  to  $10^{308}$ .

## Round-off Errors with Arithmetic Manipulations (1/3)

Roundoff error can happen in several circumstances other than just storing numbers. For example:

- i. **Large computations:** if a process performs a large number of computations, roundoff errors may build up to become significant

```
-  
function sout = sumdemo()  
s = 0;  
for i = 1:10000  
    s = s + 0.0001;  
% notice that 0.0001 cannot be expressed exactly in base - 2.  
end  
sout = s;
```

```
-  
function sout = sumdemo()  
s = 0;  
for i = 1:10000  
    s = s + 0.0001;  
>> format long  
>> sumdemo  
ans =  
0.999999999999991
```

- ii. **Adding a Large and a Small Number:** Since the small number's mantissa is shifted to the right to be the same scale as the large number, digits are lost

## Round-off Errors with Arithmetic Manipulations (2/3)

What if  $0.0010 + 4000$  is represented with 4 - digit mantissa and 1 - digit exponent?

$$\begin{array}{r} 4.000 \quad \times 10^3 \\ 0.000001 \quad \times 10^3 \\ \hline 4.000\,001 \quad \times 10^3 \end{array} \text{ (chopped to } 4.000 \times 10^3\text{)}$$

iii. **Smearing:** Smearing occurs whenever the individual terms in a summation are larger than the summation itself Round-off errors:

$(x + 10^{-20}) - x = 10^{-20}$  mathematically, but  
 $x = 1$ ;  $(x + 10^{-20}) - x$  gives a 0 in MATLAB!

## Round-off Errors with Arithmetic Manipulations (3/3)

iv. **Subtractive Cancellation:** This refers to the round-off error induced when subtracting two nearly equal floating-point numbers.

For example, one common instance where this can occur involves finding the roots of a quadratic equation (or parabola) with the quadratic formula:

$$\frac{x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}}{x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}}$$

For cases where  $b^2 \gg 4ac$ , the difference in the numerator can be very small. In such cases double precision can mitigate the problem. To avoid this, an alternative formula can be used to minimize subtractive cancellation, as:

$$\frac{x_1 = \frac{-2c}{b + \sqrt{b^2 - 4ac}}}{x_2 = \frac{-2c}{b - \sqrt{b^2 - 4ac}}}$$

## System of Non-linear Equations

(1)

- Numerical Solution of System of non-linear equations
- Let's generalize to the  $n$ -dimensional case to find the zero/root of a function as:

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad n \text{ is any +ve integer.}$$

- Our problem is to find  $x^*$  in  $\mathbb{R}^n$  such that  $F(x^*) = 0$ .

— (1)

For example:

Take the non-linear system

$$\begin{bmatrix} e^{x_1^2+x_2^2} - 1 = 0 \\ e^{x_1^2-x_2^2} - 1 = 0 \end{bmatrix} \quad \left[ \begin{array}{l} \text{Here } n=2 \text{ and} \\ F(x) = \begin{bmatrix} e^{x_1^2+x_2^2} - 1, e^{x_1^2-x_2^2} - 1 \end{bmatrix}^T \end{array} \right]$$

OR

$$F(x) = \begin{bmatrix} F_1(x) \\ F_2(x) \end{bmatrix} \text{ where}$$

$$F_1(x) = e^{x_1^2+x_2^2} - 1$$

$$F_2(x) = e^{x_1^2-x_2^2} - 1$$

$$\text{and } x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- Q: How to extend the iterative methods discussed for 1-dimension to several dimensions?

(2)

- Here, we shall always assume that  $F \in C^1(D)$  that is  
 $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a continuously differentiable on  $D$ ,  
 where  $D \subseteq \mathbb{R}^n$ .

- We denote also by  $J_F(x)$  the Jacobian matrix associated with  $F$  and evaluated at the point  $x = (x_1, x_2, \dots, x_n)^T$  of  $\mathbb{R}^n$ , defined as

$$(J_F(x))_{ij} = \left( \frac{\partial F_i}{\partial x_j} \right)(x), \quad i, j = 1, 2, \dots, n.$$

— (2)

OR

$$J_F(x) = \begin{bmatrix} \frac{\partial F_1(x)}{\partial x_1} & \frac{\partial F_1(x)}{\partial x_2} & \dots & \frac{\partial F_1(x)}{\partial x_n} \\ \frac{\partial F_2(x)}{\partial x_1} & \frac{\partial F_2(x)}{\partial x_2} & \dots & \frac{\partial F_2(x)}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial F_n(x)}{\partial x_1} & \frac{\partial F_n(x)}{\partial x_2} & \dots & \frac{\partial F_n(x)}{\partial x_n} \end{bmatrix}$$

### Newton's Method :

Here, we consider the extension of  
 Newton's method (for simple roots) to  
 the Vector case.

Given  $x^{(0)} \in \mathbb{R}^n$ , for  $k=0, 1, \dots$ , until convergence

Solve

$$\left. \begin{aligned} J_F(x^{(k)}) \delta x^{(k)} &= -F(x^{(k)}) \\ \text{Set } x^{(k+1)} &= x^{(k)} + \delta x^{(k)} \end{aligned} \right\} \quad (3)$$

Therefore, at each iterative-step  $k$  the solution of a linear system with matrix  $J_F(x^{(k)})$  is required.

Example 1:

$$\left. \begin{aligned} e^{x_1^2 + x_2^2} - 1 &= 0 \\ e^{x_1^2 - x_2^2} - 1 &= 0 \end{aligned} \right| \quad \text{exact sol: } x^* = 0$$

- Take  $x^{(0)} = (0.1, 0.1)^T$

Then after 15 iterations you will get

$$x^{(15)} = (0.61 \times 10^{-5}, 0.61 \times 10^{-5})$$

- Again, if  $x^{(0)} = (10, 10)^T$ , then 220 iterations are needed to obtain a solution comparable to previous case.

- And, if  $x^{(0)} = (20, 20)^T$ , then Newton's method fails to converge.

Remarks:

- Thus, this example points out the high sensitivity<sup>(4)</sup> of Newton's method on the choice of the initial guess  $x^{(0)}$ .
- Theoretically, it can be established that Newton's method is quadratically convergent only if  $x^{(0)}$  is sufficiently close to the solution  $x^*$  and if the Jacobian matrix is non-singular.
- Also note that computation complexity (effort) can be excessively high as  $n$  gets large.

Derivation of Newton's method (for  $n=2$  case):

- Let's derive Eq. (3) for  $n=2$  case. Therefore, here

$$F = \left( F_1(x_1, x_2), F_2(x_1, x_2) \right)^T \quad \text{--- (4)}$$

Let  $(x_1^{(k)}, x_2^{(k)})$  be a suitable approximation to the root  $x^* = (x_1^*, x_2^*)$  of the system (4). Let  $\Delta x_1$  be an increment in  $x_1^{(k)}$  and  $\Delta x_2$  be an increment in  $x_2^{(k)}$  such that  $(x_1^{(k)} + \Delta x_1, x_2^{(k)} + \Delta x_2)$  is the exact solution of (1), that is:

Means we have system of two non-linear equations

$$F_1(x_1, x_2) = 0$$

$$F_2(x_1, x_2) = 0$$

(5)

$$\left. \begin{array}{l} F_1(x_1^{(k)} + \Delta x_1, x_2^{(k)} + \Delta x_2) = 0 \\ F_2(x_1^{(k)} + \Delta x_1, x_2^{(k)} + \Delta x_2) = 0 \end{array} \right]$$

Using Taylor's series expansion about point  $(x_1^{(k)}, x_2^{(k)})$ , we get

$$\begin{aligned} & F_1(x_1^{(k)}, x_2^{(k)}) + \left[ \Delta x_1 \frac{\partial}{\partial x_1} + \Delta x_2 \frac{\partial}{\partial x_2} \right] F_1(x_1^{(k)}, x_2^{(k)}) + \frac{1}{2!} \left[ \Delta x_1 \frac{\partial}{\partial x_1} + \Delta x_2 \frac{\partial}{\partial x_2} \right]^2 \\ & F_1(x_1^{(k)}, x_2^{(k)}) + \dots = 0 \end{aligned}$$

$$\begin{aligned} & F_2(x_1^{(k)}, x_2^{(k)}) + \left[ \Delta x_1 \frac{\partial}{\partial x_1} + \Delta x_2 \frac{\partial}{\partial x_2} \right] F_2(x_1^{(k)}, x_2^{(k)}) + \frac{1}{2!} \left[ \Delta x_1 \frac{\partial}{\partial x_1} + \Delta x_2 \frac{\partial}{\partial x_2} \right]^2 \\ & F_2(x_1^{(k)}, x_2^{(k)}) + \dots = 0 \end{aligned}$$

Neglecting second and higher powers of  $\Delta x_1$  and  $\Delta x_2$ , we get:

$$\left. \begin{array}{l} F_1(x_1^{(k)}, x_2^{(k)}) + \Delta x_1 \frac{\partial F_1(x_1^{(k)}, x_2^{(k)})}{\partial x_1} + \Delta x_2 \frac{\partial F_1(x_1^{(k)}, x_2^{(k)})}{\partial x_2} = 0 \\ F_2(x_1^{(k)}, x_2^{(k)}) + \Delta x_1 \frac{\partial F_2(x_1^{(k)}, x_2^{(k)})}{\partial x_1} + \Delta x_2 \frac{\partial F_2(x_1^{(k)}, x_2^{(k)})}{\partial x_2} = 0 \end{array} \right] \quad (5)$$

Solving (5) for  $\Delta x_1$  and  $\Delta x_2$ , we get:

$$\left[ \begin{array}{cc} \frac{\partial F_1(x_1^{(k)}, x_2^{(k)})}{\partial x_1} & \frac{\partial F_1(x_1^{(k)}, x_2^{(k)})}{\partial x_2} \\ \frac{\partial F_2(x_1^{(k)}, x_2^{(k)})}{\partial x_1} & \frac{\partial F_2(x_1^{(k)}, x_2^{(k)})}{\partial x_2} \end{array} \right] \left[ \begin{array}{c} \Delta x_1 \\ \Delta x_2 \end{array} \right] = - \left[ \begin{array}{c} F_1(x_1^{(k)}, x_2^{(k)}) \\ F_2(x_1^{(k)}, x_2^{(k)}) \end{array} \right] \quad (6)$$

(6)

OR

$$\mathcal{J}_F(x^{(k)}) \Delta x^{(k)} = -F(x^{(k)})$$

which is nothing  
but Eq. (3).

The solution of the system (6) is :

$$\Delta x^{(k)} = -\left(\mathcal{J}_F(x^{(k)})\right)^{-1} F(x^{(k)})$$

Where  $\left(\mathcal{J}_F(x^{(k)})\right)^{-1}$  is the inverse of Jacobian matrix

$\mathcal{J}_F(x^{(k)})$ . Then set  $x^{(k+1)} = x^{(k)} + \Delta x^{(k)}$

□

Example 2:

Perform three iterations of the Newton's method to solve the system of equations.

$$x_1^2 + x_1 x_2 + x_2^2 = 7$$

$$x_1^3 + x_2^3 = 9$$

Here exact sol.

$$x^* = (x_1^*, x_2^*) = (2, 1)$$

with  $x^{(0)} = (x_1^{(0)}, x_2^{(0)}) = (1.5, 0.5)$

Sol.<sup>n</sup> Take  $F_1 = x_1^2 + x_1 x_2 + x_2^2 - 7$   
 $F_2 = x_1^3 + x_2^3 - 9$  ] and  $F = (F_1, F_2)$

$$\therefore \mathcal{J}_F(x^{(k)}) = \begin{bmatrix} \frac{\partial F_1(x^{(k)})}{\partial x_1} & \frac{\partial F_1(x^{(k)})}{\partial x_2} \\ \frac{\partial F_2(x^{(k)})}{\partial x_1} & \frac{\partial F_2(x^{(k)})}{\partial x_2} \end{bmatrix}$$

(7)

$$\text{Solve: } \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \end{bmatrix} = -\left(\mathcal{J}_F(x^{(k)})\right)^{-1} \begin{bmatrix} F_1(x^{(k)}) \\ F_2(x^{(k)}) \end{bmatrix}$$

OR

$$\begin{bmatrix} x_1^{(k+1)} \\ x_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \end{bmatrix} + \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \end{bmatrix} = \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \end{bmatrix} - \left(\mathcal{J}_F(x^{(k)})\right)^{-1} \begin{bmatrix} F_1(x^{(k)}) \\ F_2(x^{(k)}) \end{bmatrix}$$

OR (in short) Solve:

$$x^{(k+1)} = x^{(k)} - \left(\mathcal{J}_F(x^{(k)})\right)^{-1} F(x^{(k)})$$

where  $k = 0, 1, 2, \dots$ 

□

## Roots for a function of complex variable

— The root of an equation  $f(z) = 0$ ,  $z$  is a complex variable, can be obtained by using the method (Newton's) for the system of non-linear equations.

$$f(z) = 0, z = x + iy$$

$$\Rightarrow f(z) = u(x, y) + iv(x, y) = 0 \quad \text{--- (1)}$$

OR

$$\begin{cases} u(x, y) = 0 \\ v(x, y) = 0 \end{cases} \quad \text{--- (2)}$$

Thus, the problem of finding a complex root of (1) reduces to solving a system of two non-linear equations (2).

Example:

$$f(z) = z^3 + 1 = 0$$

with initial guess  $(x_0, y_0) = (0.25, 0.25)$

$$f(x+iy) = u(x, y) + iv(x, y) = (x^3 - 3xy^2 + 1) + i(3x^2y - y^3) = 0$$

$$u(x, y) = x^3 - 3xy^2 + 1 = 0$$

$$v(x, y) = 3x^2y - y^3 = 0$$

$$\text{Jacobian matrix } J = \begin{bmatrix} u_x & u_y \\ v_x & v_y \end{bmatrix} = \begin{bmatrix} 3x^2 - 3y^2 & -6xy \\ 6xy & 3x^2 - 3y^2 \end{bmatrix}$$

$$\text{and } |J| = 9(x^2 + y^2)^2$$

$$\therefore J^{-1} = \frac{1}{|J|} \begin{bmatrix} 3(x^2 - y^2) & 6xy \\ -6xy & 3(x^2 - y^2) \end{bmatrix}$$

l2

Newton's method:

$$\begin{bmatrix} x_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ y_k \end{bmatrix} - \frac{1}{|J_k|} \begin{bmatrix} 3(x_k^2 - y_k^2) & 6x_k y_k \\ -6x_k y_k & 3(x_k^2 - y_k^2) \end{bmatrix} \begin{bmatrix} x_k^3 - 2x_k y_k^2 + 1 \\ 3x_k^2 y_k - y_k^3 \end{bmatrix}$$

with  $\begin{bmatrix} x_0 \\ y_0 \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.25 \end{bmatrix}$   $k=0, 1, 2, \dots$

Start iteration to get:  $x^{(1)} = \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}, x^{(2)} = \begin{bmatrix} x_2 \\ y_2 \end{bmatrix}, \dots$

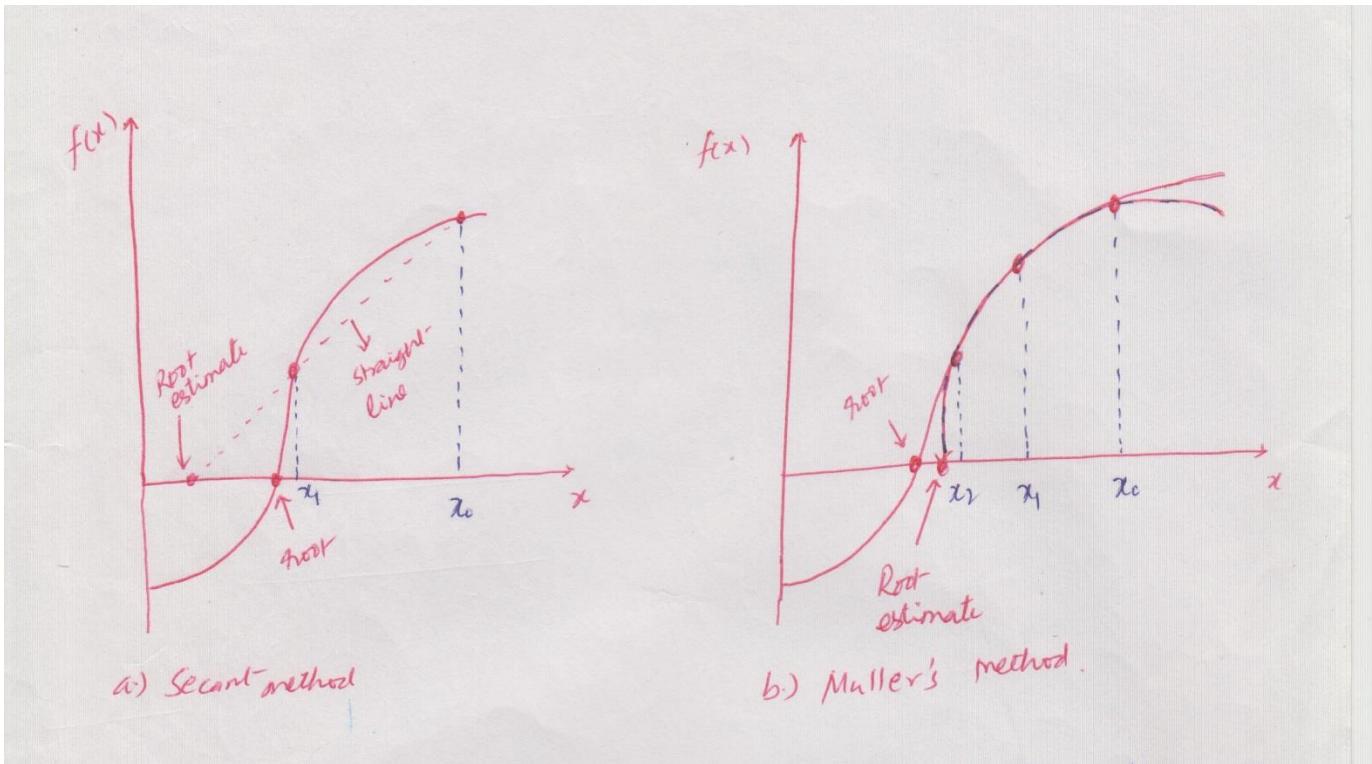
□

To find the real and complex roots of polynomial

- Müller's method ✓
- Bairstow method

Müller's method:

- As Secant method finds a root estimate by projecting a line to the x-axis through two function values. Müller's method follows a similar approach, but projects a parabola through three points.
- The method consists of deriving the coefficients of the parabola passing through three points.



3

These coefficients can then be substituted into the quadratic formula to obtain the point where the parabola intercepts the  $x$ -axis — that is, the root estimate.

$$f_2(x) = a(x-x_2)^2 + b(x-x_2) + c \quad \text{--- (3)}$$

We want this parabola to intersect the three points:

$$(x_0, f(x_0)), (x_1, f(x_1)) \text{ and } (x_2, f(x_2))$$

Coefficient of (3) can be obtained as:

$$\left. \begin{array}{l} f(x_0) = a(x_0-x_2)^2 + b(x_0-x_2) + c \\ f(x_1) = a(x_1-x_2)^2 + b(x_1-x_2) + c \\ f(x_2) = a(x_2-x_2)^2 + b(x_2-x_2) + c \end{array} \right\} \quad \text{--- (4)}$$

$c = f(x_2)$

And,

$$\left. \begin{array}{l} f(x_0) - f(x_2) = a(x_0-x_2)^2 + b(x_0-x_2) \\ f(x_1) - f(x_2) = a(x_1-x_2)^2 + b(x_1-x_2) \end{array} \right\} \quad \text{--- (5)}$$

How to get  $a$  and  $b$ ?

Notations:  $h_0 = x_1 - x_0$ ,  $h_1 = x_2 - x_1$

$$\delta_0 = \frac{f(x_1) - f(x_0)}{x_1 - x_0}, \quad \delta_1 = \frac{f(x_2) - f(x_1)}{x_2 - x_1}$$

- ∵ From (5) :

$$(h_0 + h_1)b - (h_0 + h_1)^2 a = h_0 \delta_0 + h_1 \delta_1$$

$$h_1 b - h_1^2 a = h_1 \delta_1$$

Solve for  $a$  and  $b$ :

$$a = \frac{\delta_1 - \delta_0}{h_1 + h_0}, \quad b = ah_1 + \delta_1 \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \text{--- (6)}$$

and  $c = f(x_2)$

- To find the root, we apply the quadratic formula to Eqn. (3), however, because of potential round-off error, we use the alternative formulation to yield:

$$x_3 - x_2 = \frac{-2c}{b \pm \sqrt{b^2 - 4ac}} \quad \text{OR} \quad x_3 = x_2 + \left[ \frac{-2c}{b \pm \sqrt{b^2 - 4ac}} \right] \quad \text{--- (7)}$$

- Note that the use of quadratic formula means that both real and complex roots can be located. This is the major benefit of the formula method.

15

- Note that ⑦ yields two roots corresponding to  $\pm$  term in the denominator. In Müller's method, the sign is chosen to agree with the sign of  $b$ . This choice will result in the largest denominator, and hence will give the root estimate that is closest to  $x_2$ .
- Once  $x_3$  is determined, the process is repeated. This brings up the issue of which point is discarded. Two general strategies are typically used:

- 1.) If only real roots are being located, we choose the two original points that are nearest the new root estimate  $x_3$ .
- 2.) If both real and complex roots are being evaluated, a sequential approach is employed. That is, just like the Secant method  $x_1, x_2$  and  $x_3$  take the place of  $x_0, x_1$  and  $x_2$ .

Example :

Perform three iterations using Müller's method to find the smallest +ve root of the equation:

$$f(x) = x^3 - 5x + 1 = 0$$

Sol: Note that the smallest +ve root of  $f(x)$  lies in  $(0, 1)$ .

Take :  $x_0 = 0, x_1 = 0.5, x_2 = 1.0$  (initial guesses)

$$\Rightarrow f(x_0) = 1.0, \quad f(x_1) = -1.375, \quad f(x_2) = -3.0$$

6

$$\therefore C = f(x_2) = -3.0, \quad a = \frac{s_1 - s_0}{h_1 + h_0} = 1.5$$

$$b = ah_1 + s_1 = -2.5$$

Note that, since  $b < 0$ , we get:

$$x_3 = x_2 - \frac{2c}{b - \sqrt{b^2 - 4ac}} = 0.191857$$

For new iterations, choose:  $x_0 = 0.5, x_1 = 1.0, x_2 = 0.1918$

Then

$$C = 0.0477, \quad a = 1.6991, \quad b = -5.1305$$

again, as  $b < 0$ , we get

$$x_4 \Rightarrow \text{as} \quad x_3 = x_2 - \frac{2c}{b - \sqrt{b^2 - 4ac}} = 0.2011$$

again, choose:  $x_0 = 1.0, x_1 = 0.1918, x_2 = 0.2011$

1  
1  
1

□

## Interpolation

11

- Consider the problem of approximating a given function by a class of simpler functions (polynomials).
- Main uses of interpolation:
  - a) In reconstructing the function  $f(x)$  when it is not given explicitly, and only the values and/or its certain order derivatives are given at a set of points.
  - b) To replace the  $f(x)$  by a interpolating polynomial  $p(x)$  so that many common operations such as determination of roots, differentiation and integration etc. which are intended for  $f(x)$  may be performed using  $p(x)$ .

Definition: A polynomial  $p(x)$  is called an interpolating polynomial if the values of  $p(x)$  and/or its certain order derivatives coincide with those of  $f(x)$  and/or its same order derivatives at one or more points.

- In general, if there are  $(n+1)$  distinct points:  
$$a \leq x_0 \leq x_1 \leq \dots \leq x_n \leq b$$
Then the problem of interpolation is to obtain  $p(x)$  satisfying

the condition:

(2)

i)  $p(x_i) = f(x_i), i = 0, 1, 2, \dots, n$  — (1)

OR

ii)  $\begin{cases} p(x_i) = f(x_i) \\ p'(x_i) = f'(x_i), i = 0, 1, \dots, n \end{cases}$  — (2)

— Condition in Eq. (1) give rise to Lagrange interpolation, and (2) gives Hermite interpolation.

Note: The derivative condition in Eq. (2) may be replaced with more general derivative condition involving higher-order derivatives.

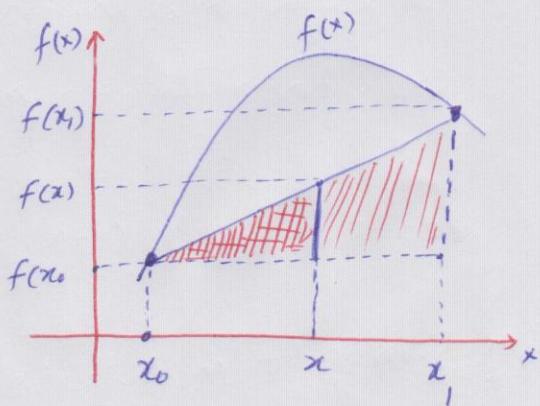
Next, we discuss the interpolations of various degrees.

Linear interpolation (Lagrange):

— Using property of similar triangles:

$$\frac{f(x) - f(x_0)}{x - x_0} = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$$

OR



3

$$\begin{aligned}f(x) &= \left(\frac{x-x_1}{x_0-x_1}\right)f(x_0) + \left(\frac{x-x_0}{x_1-x_0}\right)f(x_1) \\&= l_0(x)f(x_0) + l_1(x)f(x_1) \quad \text{--- (3)}$$

Where,  $l_0(x) = \frac{x-x_1}{x_0-x_1}$ ,  $l_1(x) = \frac{x-x_0}{x_1-x_0}$

— functions  $l_0(x)$  and  $l_1(x)$  are called Lagrange fundamental polynomials. It can be verified that:

$$l_0(x) + l_1(x) = 1$$

$$\begin{cases} l_0(x_0) = 1, & l_0(x_1) = 0 \\ l_1(x_0) = 0, & l_1(x_1) = 1 \end{cases}$$

OR

$$l_i(x_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

↓  
Kronecker delta

— Also, Eqn.(3) is a weighted average of  $f(x_0)$  and  $f(x_1)$ . Here, weights ( $l_i(x)$ ) are inversely related to the distance from end points to the unknown point; the closer point has more influence than the farther point.

Thus the weights  $\left(\frac{x-x_1}{x_0-x_1}\right)$  and  $\left(\frac{x-x_0}{x_1-x_0}\right)$  which are normalized distances between the unknown point and each end point.

### Truncation error(bounds):

14

In linear interpolation (Lagrange) case, polynomial  $f_1(x)$  coincides with the function  $f(x)$  at  $x_0$  and  $x_1$ , and it deviates at all other points in interval  $(x_0, x_1)$  (as shown in Fig. 1). This deviation will introduce truncation error and denoted as :

$$E_1 = f(x) - f_1(x) \quad \text{--- (4)}$$

### Expression for $E_1(f, x)$ :

— Note that at  $x=x_0$  or  $x=x_1 \Rightarrow E_1(f, x) = 0$   
 $\therefore$  If  $x \in (x_0, x_1)$ , then define a function  $g(t)$  as:

$$g(t) = f(t) - f_1(t) - (f(x) - f_1(x)) \frac{(t-x_0)(t-x_1)}{(x-x_0)(x-x_1)} \quad \text{--- (5)}$$

Check that:  $g(t) = 0$  if  $t = x_0, t = x_1$ , or  $t = x$

$\Rightarrow g(t)$  satisfies Rolle's theorem.

Applying Rolle's theorem on the intervals  $(x_0, t)$  and  $(t, x_1)$  separately, we get:

$$g'(\xi_1) = 0, \quad x_0 < \xi_1 < t \quad \text{and} \quad g'(\xi_2) = 0, \quad t < \xi_2 < x_1$$

15

—  $\Rightarrow g'(t)$  also satisfies the conditions of Rolle's theorem.

$\therefore$  Applying Rolle's theorem for  $g'(t)$  on interval  $(\bar{x}_1, \bar{x}_2)$

we obtain  $g''(\bar{z}) = 0, \quad \bar{x}_1 < \bar{z} < \bar{x}_2$  or  $x_0 < \bar{z} < x_1$

From ⑤:

$$g''(t) = f''(t) - \frac{2(f(x) - f_1(x))}{(x-x_0)(x-x_1)} \quad \text{--- (6)}$$

$\therefore g''(\bar{z}) = 0$ , from ⑥, we have

$$f(x) = f_1(x) + \frac{1}{2}(x-x_0)(x-x_1)f''(\bar{z})$$

$$\Rightarrow E_1(f, x) = f(x) - f_1(x) = \frac{1}{2}(x-x_0)(x-x_1)f''(\bar{z})$$

— If we can determine a bound for  $f''(x)$  in  $[x_0, x_1]$ ,  
that is,

$$|f''(x)| \leq M_2, \quad \text{for all } x \in [x_0, x_1]$$

then

$$|E_1(f, x)| = \left| \frac{1}{2}(x-x_0)(x-x_1)f''(\bar{z}) \right|$$

OR

$$\boxed{|E_1(f, x)| \leq \frac{1}{2} \max_{x_0 \leq x \leq x_1} |(x-x_0)(x-x_1)| M_2} \quad \text{--- (7)}$$

## Lagrange Interpolating polynomials (in general):

- In general, Lagrange interpolating polynomials can be represented as:

$$f_n(x) = \sum_{i=0}^n l_i(x) f(x_i) \quad \text{--- ⑧}$$

where  $l_i(x) = \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j}$

For linear ( $n=1$ ) case:

$$f_1(x) = \left( \frac{x - x_1}{x_0 - x_1} \right) f(x_0) + \left( \frac{x - x_0}{x_1 - x_0} \right) f(x_1) \quad \text{--- ⑨}$$

and the second order version ( $n=2$ ) is:

$$f_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2) \quad \text{--- ⑩}$$

Also, note that

$$l_i(x_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

and  $\sum_{i=0}^n l_i(x) = 1$

□

## Interpolation (Contd.)

(1)

Example: Given that  $f(0) = 1$ ,  $f(1) = 3$ ,  $f(3) = 55$  find the unique polynomial of degree 2 (or less) which fits the given data:

$x$	$f(x)$
0	1
1	3
3	55

we have  $x_0 = 0$ ,  $x_1 = 1$ ,  $x_2 = 3$   
and  $f(x_0) = 1$ ,  $f(x_1) = 3$ ,  $f(x_2) = 55$

Lagrange fundamental polynomials are:

$$l_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} = \frac{(x-1)(x-3)}{(-1)(-3)} = \frac{1}{3}(x^2 - 4x + 3)$$

$$l_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} = \frac{1}{2}(3x - x^2)$$

$$l_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} = \frac{1}{6}(x^2 - x)$$

Hence, the Lagrange quadratic interpolating polynomial  $f_2(x)$  is:

$$f_2(x) = \sum_{i=0}^2 l_i(x) f(x_i) = 8x^2 - 6x + 1$$

□

Exercise:

L2

Consider the following values of the function

$$f(x) = \sin x + \cos x \text{ as:}$$

$x$	$f(x)$
$\pi/18$	1.1585
$\pi/9$	1.2817
$\pi/6$	1.3660

Construct the quadratic interpolating polynomial  $f_2(x)$  that fits the data. Hence, find  $f(\pi/2)$ . Compare with the exact value.

Error in quadratic interpolation (Lagrange):

$$E_2(f, x) = f(x) - f_2(x)$$

To get the truncation error bound, take

$$g(t) = f(t) - f_2(t) - (f(x) - f_2(x)) \frac{(t-x_0)(t-x_1)(t-x_2)}{(x-x_0)(x-x_1)(x-x_2)}$$

— Apply Rolle's theorem repeatedly for  $g(t)$ ,  $g'(t)$  and  $g''(t)$ .

$$\Rightarrow g'''(\bar{x}) = 0, \quad x_0 < \bar{x} < x_2$$

$$\Rightarrow E_2(f, x) = \frac{1}{3!} (x-x_0)(x-x_1)(x-x_2) f'''(\bar{x}) \quad \text{--- (1)}$$

3

— Now if we can find +ve real  $M_3$  such that

$$M_3 = \max_{x_0 \leq x \leq x_2} |f'''(x)|$$

then error bound is:

$$|E_2| = |f(x) - f_2(x)| \leq \frac{1}{6} M_3 \left( \max_{x_0 \leq x \leq x_2} |(x-x_0)(x-x_1)(x-x_2)| \right)$$

□

— Similarly, in general, for higher order interpolating polynomial, say of degree  $n$ :

$$E_n(f, x) = \frac{1}{(n+1)!} (x-x_0)(x-x_1) \cdots (x-x_n) f^{(n+1)}(\xi)$$

with  $x_0 < \xi < x_n$

————— (3)

and

$$|E_n| \leq \frac{1}{(n+1)!} M_n \left( \max_{x_0 \leq x \leq x_n} |(x-x_0)(x-x_1) \cdots (x-x_n)| \right)$$

Where

$$M_n = \max_{x_0 \leq x \leq x_n} |f^{(n+1)}(x)|$$

□

— Let's find the truncation error bound for Example 1.

4

$$|E_2| \leq \frac{1}{6} M_2 \left( \max_{0 \leq x \leq 3} |(x-0)(x-1)(x-3)| \right)$$

$$= \frac{1}{6} (2.1126) M_2$$

□

Find max. of  $|x(x-1)(x-3)|$ .  
Check that it occurs at  
 $x = 2.1126$

— There are two alternatives that are well suited for computer implementation

i) Lagrange form

ii) Newton's form (Newton's divided difference form)

— Let's derive the quadratic interpolating polynomial  $f_2(x)$  in Newton's form.

If there are three data points are available, we can use a second-order polynomial (or parabola) given as:

$$f_2(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) \quad \text{--- (4)}$$

at  $\underline{x=x_0}$  :  $a_0 = f(x_0)$  --- (5)

Using  $a_0$  in (4), at  $x=x_1$ :

$$a_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0} \quad \text{--- (6)}$$

Finally, using ⑤ and ⑥ in ④ at  $x=x_2$  gives:

LS

$$a_2 = \frac{\frac{f(x_2) - f(x_1)}{x_2 - x_1} - \frac{f(x_1) - f(x_0)}{x_1 - x_0}}{x_2 - x_0} \longrightarrow ⑦$$

Values of coefficients obtained in Eqs. ⑤, ⑥, ⑦ can be substituted in ④ to get  $f_2(x)$ .

### General form of Newton interpolating polynomials:

To fit an  $n^{\text{th}}$ -order polynomial to  $(n+1)$  data points, the  $n^{\text{th}}$ -order polynomial can be written as:

$$f_n(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \dots + a_n(x - x_0)(x - x_1)\dots(x - x_{n-1}) \longrightarrow ⑧$$

Here also data points can be used to evaluate the coefficients  $a_0, a_1, \dots, a_n$ , as:

$$\left. \begin{array}{l} a_0 = f(x_0) \\ a_1 = f[x_1, x_0] \\ a_2 = f[x_2, x_1, x_0] \\ \vdots \\ a_n = f[x_n, x_{n-1}, \dots, x_1, x_0] \end{array} \right\} \longrightarrow ⑨$$

Where, [ ] denotes the finite divided differences. [6]

For example, the first finite divided difference is:

$$f[x_i, x_j] = \frac{f(x_i) - f(x_j)}{x_i - x_j}$$

The second divided difference, which is the difference of two first divided differences, is expressed as:

$$f[x_i, x_j, x_k] = \frac{f[x_i, x_j] - f[x_j, x_k]}{x_i - x_k}$$

Similarly, the  $n^{\text{th}}$  finite divided difference is:

$$f[x_n, x_{n-1}, \dots, x_1, x_0] = \frac{f[x_n, x_{n-1}, \dots, x_1] - f[x_{n-1}, x_{n-2}, \dots, x_0]}{x_n - x_0}$$

$i$	$x_i$	$f(x_i)$	First	Second	Third
0	$x_0$	$f(x_0)$	$f[x_1, x_0]$		
1	$x_1$	$f(x_1)$	$f[x_2, x_1]$	$f[x_2, x_1, x_0]$	$f[x_3, x_2, x_1, x_0]$
2	$x_2$	$f(x_2)$		$f[x_3, x_2, x_1]$	
3	$x_3$	$f(x_3)$	$f[x_3, x_2]$		

Table 1: Recursive nature of finite-divided difference.

7

∴ Using ⑨ in terms of divided differences:

$$\begin{aligned}
 f_n(x) &= f(x_0) + (x-x_0)f[x, x_0] + (x-x_0)(x-x_1)f[x_2, x_1, x_0] + \dots \\
 &\quad + (x-x_0)(x-x_1)\dots(x-x_{n-1})f[x_n, x_{n-1}, \dots, x_1, x_0]
 \end{aligned}
 \tag{10}$$

Example 2:

$x_i$	$f(x_i)$
2.5	14
3.2	15
2.0	8

Use Newton's divided difference to calculate  $f_2(x)$ : Also find  $f_2(2.8)$ .

$i$	$x_i$	$f(x_i)$	$f[x_{i+1}, x_i]$	$f[x_{i+2}, x_{i+1}, x_i]$
0	2.5	14	1.42857	
1	3.2	15		-8.80952
2	2.0	8	5.03333	

∴ Second order interpolation can be implemented as:

$$f_2(x) = f(x_0) + (x-x_0)f[x, x_0] + (x-x_0)(x-x_1)f[x_2, x_1, x_0]$$

at  $x = 2.8$ :

$$f_2(2.8) = 15.48571$$

□

## Errors of Newton's Interpolating Polynomials

11

- As you have already seen that  $n^{\text{th}}$ -order interpolating polynomial (in Newton's form) is:

$$f_n(x) = f(x_0) + (x-x_0)f[x_1, x_0] + (x-x_0)(x-x_1)f[x_2, x_1, x_0] + \dots \\ \downarrow \\ + (x-x_0)(x-x_1) \dots (x-x_{n-1})f[x_n, x_{n-1}, \dots, x_0]$$

Newton's divided difference interpolating polynomial. —①

- Note that Eq. ① is similar to the Taylor series expansion, and the terms are finite divided differences.

- Similar to Taylor series expansion, here also the truncation part (error) could be expressed as:

$$R_n = f[x, x_n, x_{n-1}, \dots, x_1, x_0](x-x_0)(x-x_1) \dots (x-x_n) \quad —②$$

Where  $f[x, x_n, x_{n-1}, \dots, x_0]$  is the  $(n+1)^{\text{th}}$  finite divided difference.

- As finite divided difference contains the unknown  $f(x)$ , therefore it can not be solved for the error.

- However, if an additional data point  $f(x_{n+1})$  is available, then

$$R_n \cong f[x_{n+1}, x_n, \dots, x_0](x-x_0)(x-x_1) \dots (x-x_n) \quad —③$$

□

12

### Example 1: (Quadratic interpolation)

$$\begin{array}{ll} x_0 = 1 & f(x_0) = 0 \\ x_1 = 4 & f(x_1) = 1.386294 \\ x_2 = 6 & f(x_2) = 1.791759 \end{array}$$

check that:

$$f_2(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1)$$

where

$$a_0 = f(x_0) = 0$$

$$a_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0} = f[x_1, x_0] = 0.4620981$$

$$a_2 = \frac{\left( \frac{f(x_2) - f(x_1)}{x_2 - x_1} \right) - \left( \frac{f(x_1) - f(x_0)}{x_1 - x_0} \right)}{x_2 - x_0} = f[x_2, x_1, x_0] = -0.05187$$

$$f_2(x) = 0 + 0.4620981(x - 1) - 0.0518731(x - 1)(x - 4)$$

Now:  
 $f_2(2) = 0.5658444$

### Example 2:

Using Eq. (3) estimate the error for the quadratic interpolating polynomial obtained in Example-1.

- If ~~value~~ function  $f(x)$  is not known to us, then we can estimate the error using (3) as (provided we know the additional value at  $x_3$ ):

$$R_2 = f[x_3, x_2, x_1, x_0](x - x_0)(x - x_1)(x - x_2)$$

OR

$$R_2 = 0.007865(x-1)(x-4)(x-6)$$

136

Note that at  $x=2$ :

$$R_2 = 0.0629242$$

□

- Now let's do an example to get the idea about the significance of the position and sequence of the given data.
- This is discussed in Example 3.

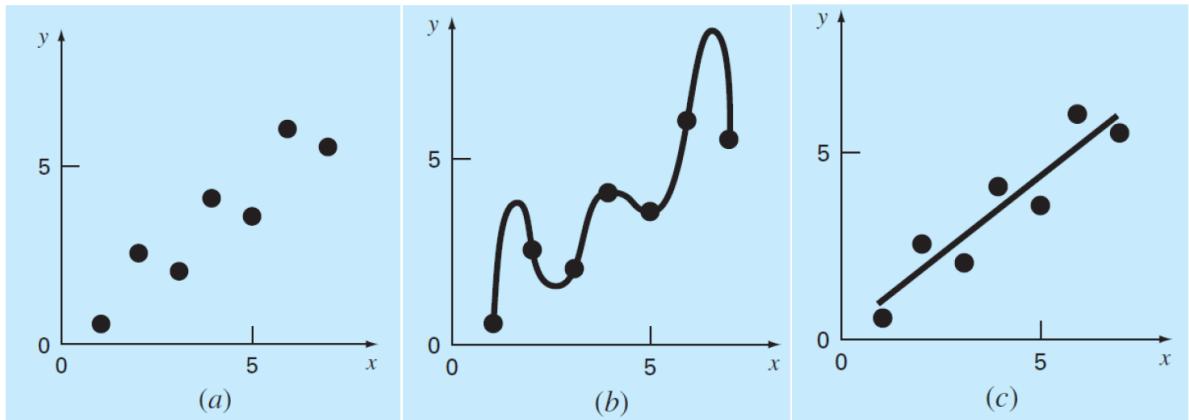


Fig. 1. Difference between interpolation and regression: (a) data points; (b) polynomial fit (interpolation); (c) using least-squares (regression) fit.

**Remarks:**

Figure 1. shows that when data is highly scattered, polynomial interpolation is inappropriate. A more appropriate technique for such cases is to derive an approximating function that fits the shape or general trend of the data (without matching its values at given data points)

**Example-1:** See the following data set for  $f(x) = \ln x$  to evaluate  $f(x)$  at  $x = 2$

<b>x</b>	<b><math>f(x) = \ln x</math></b>	<b>Order F(x)</b>	<b>Error</b>
1	0	0 0.000000	0.462098
4	1.3862944	1 0.462098	0.103746
6	1.7917595	2 0.565844	0.062924
5	1.6094379	3 0.628769	0.046953
3	1.0986123	4 0.675722	0.021792
1.5	0.4054641	5 0.697514	-0.003616
2.5	0.9162907	6 0.693898	-0.000459
3.5	1.2527630	7 0.693439	

Fig. 2 Tabulated data (left side) and different orders interpolation. Note that true value is  $\ln 2 = 0.6931472$

14

- Note from Fig. 2 that error reduces as the order increases.
- It can be concluded that the fifth order version yields a good estimate and that higher-order terms do not significantly enhance the prediction.

### Inverse interpolation:

- Interpolation: finding some value  $f(x)$  for some  $x$  that is between given data points.
- Inverse interpolation: Find the  $x$  for which  $f(x)$  is a certain value.

For example:  $f(x) = \frac{1}{x}$  (Example 3.)

$x$	1	2	3	4	5	6	7
$f(x)$	1	0.5	0.3333	0.25	0.2	0.1667	0.1429

→

$f(x)$	0.1429	0.1667	0.2	0.25	0.3333	0.5	1
$x$	7	6	5	4	3	2	1

- Now suppose you are asked to determine the value of  $x$  that corresponds to  $f(x) = 0.3$ . Such a problem is known as inverse interpolation.

- Unfortunately, when you reverse the variables, then in general the values along the new abscissa will not be evenly spaced. [5]
- Such non-uniform spacing on the abscissa often leads to oscillations in the resulting interpolating polynomial. This can occur even for lower order polynomials.
- Therefore, rather than finding an interpolation of  $x$  as a function of  $f(x)$ , it may be useful to find an equation for  $f(x)$  as a function of  $x$  using interpolation and <sup>then</sup> solve the corresponding root problem:
  - That is interpolation problem reduces to a root problem.

For the problem mentioned in Example 3:

- Fit a quadratic polynomial to the three points:

$$(2, 0.5), (3, 0.3333) \text{ and } (4, 0.25)$$

$$\Rightarrow f_2(x) = 1.08333 - 0.375x + 0.041667x^2$$

Now the answer to the inverse interpolation of finding the  $x$  corresponding to  $f(x) = 0.3$  would therefore involve determining the root of

$$0.3 = 1.08333 - 0.375x + 0.041667x^2$$

16

Using quadratic formula (for roots) :

$$x = 0.375 \pm \frac{\sqrt{(-0.375)^2 - 4(0.041667)0.78333}}{2(0.041667)} = \frac{5.704158}{3.295842}$$

Thus, the second root, 3.296, is a good approximation of the true value of 3.3333.  $\square$

## Oscillations (Interpolation)

Higher-order polynomials can not only lead to round-off errors due to ill-conditioning, but can also introduce oscillations to an interpolation problem. Consider the following function known as Runge's function:

$$f(x) = \frac{1}{1 + 25x^2}$$

In the figures below, the dashed line represents Runge's function, the circles represent data points of the function, and the solid line represents the results of a polynomial interpolation:

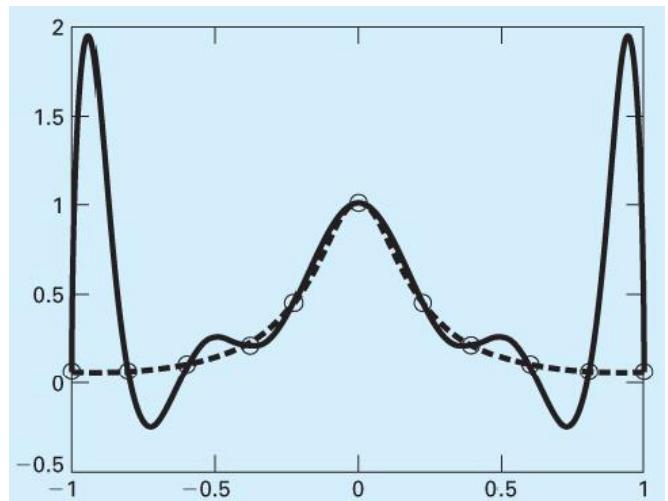
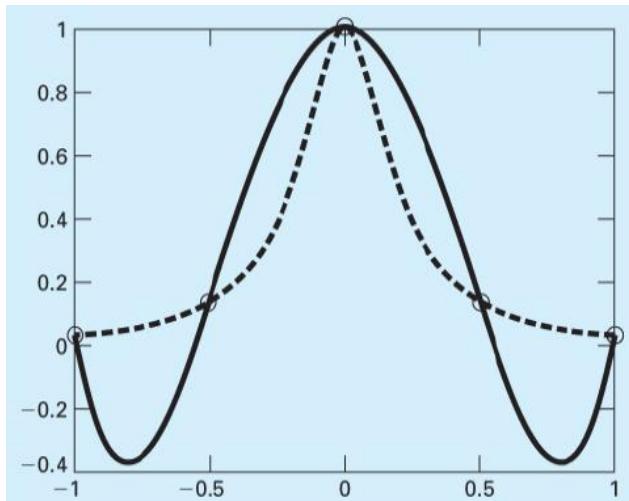


Figure 1:

a. Interpolation using fourth degree polynomial  
with 5 data points.

b. Interpolation using tenth degree polynomial  
with 11 data points

### Spline Interpolation:

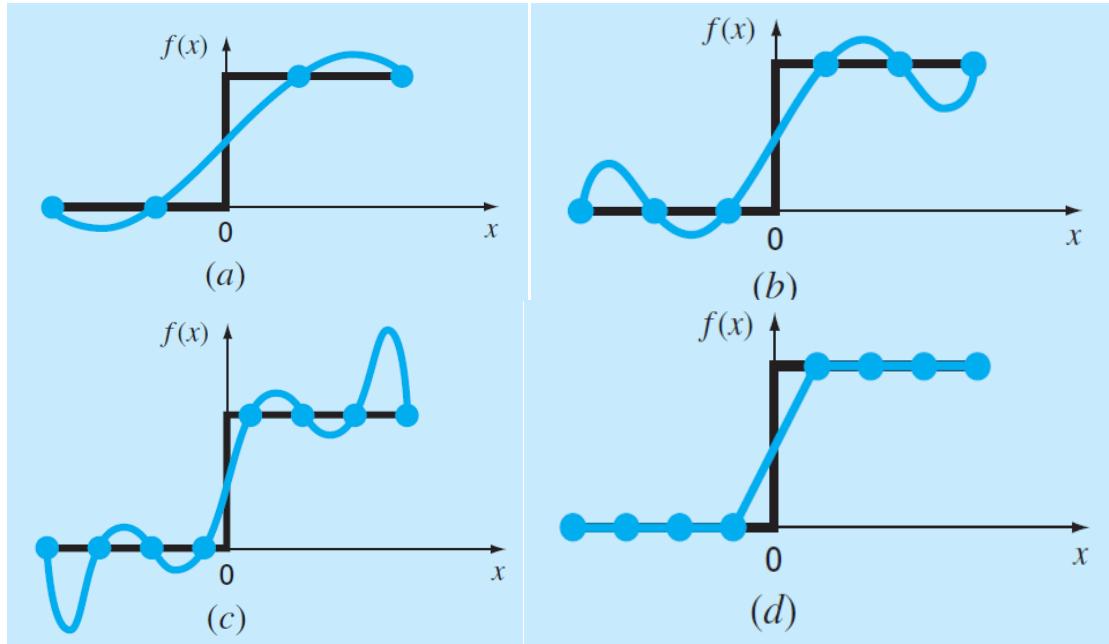


Figure 2: Representation of a situation where splines are superior to higher-order interpolating polynomials. The function to be fit undergoes an abrupt increase at  $x = 0$ . Parts (a) through (c) indicate that the abrupt change induces oscillations in interpolating polynomials. In contrast, a linear spline (d) provides a much more acceptable approximation.

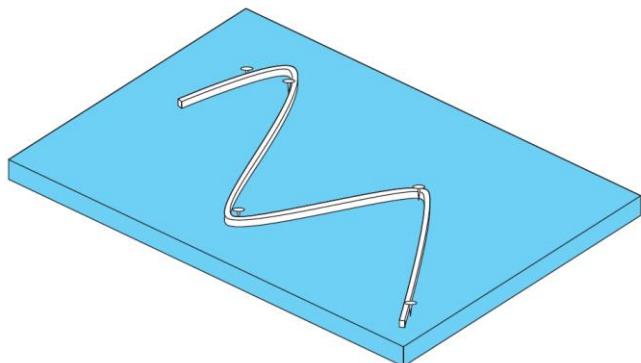


Figure 3: The concept of spline originated from the drafting technique of using a thin flexible strip (spline) to draw smooth curves through a set of points.

## Spline Interpolation

11

- An alternative approach — to apply lower-order polynomials to subset of data points. Such connecting polynomials are called Spline functions.

### Linear splines

- The simplest connection between two (data) points is a straight line
- The first-order splines for a collection of ordered data points can be defined as a set of linear functions:

$$\left. \begin{array}{l} f(x) = f(x_0) + m_0(x - x_0), \quad x_0 \leq x \leq x_1, \\ f(x) = f(x_1) + m_1(x - x_1), \quad x_1 \leq x \leq x_2 \\ \vdots \\ f(x) = f(x_{n-1}) + m_{n-1}(x - x_{n-1}), \quad x_{n-1} \leq x \leq x_n \end{array} \right\} \quad \text{--- (1)}$$

Where  $m_i$  ( $i=0, 1, \dots, n-1$ ) is the slope of straight line joining the points.

$$m_i = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} \quad \text{--- (2)}$$

- These equations can be used to evaluate the function at any point between  $x_0$  and  $x_n$  by first locating the interval within which the point lies.

Example 1:

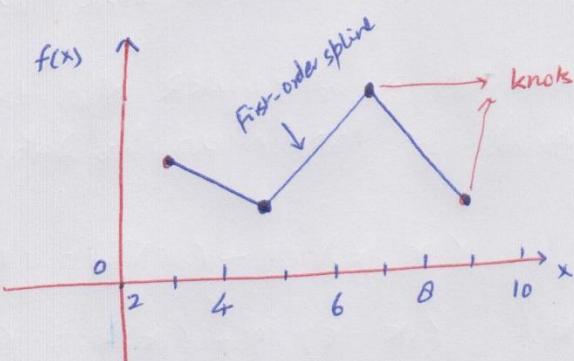
$x$	$f(x)$
3.0	2.5
4.5	1.0
7.0	2.5
9.0	0.5

Fit the data using first-order splines. Evaluate the function at  $x=5$ .

Here,  $s$  lies between  $x=4.5$  and  $x=7.0$

Therefore,

$$m = \frac{2.5 - 1}{7 - 4.5} = 0.60$$



Knots: data points where two splines meet

$$f(x) = f(4.5) + 0.60(x - 4.5)$$

$$\text{at } x=5, f(5) = 1.3 \quad \square$$

### Quadratic Splines:

- To ensure that  $n^{\text{th}}$  derivatives are continuous at the knots, a spline of atleast  $(n+1)$  order must be used.
- Third-order polynomials (cubic splines) ensures continuous first and second derivatives — Frequently used for engineering applications.

- Quadratic splines - having continuous first derivatives at the knots.
- Derive a second-order polynomial for each interval between data points. Therefore, in general:

$$f(x_i) = a_i x^2 + b_i x + c_i \quad (3)$$

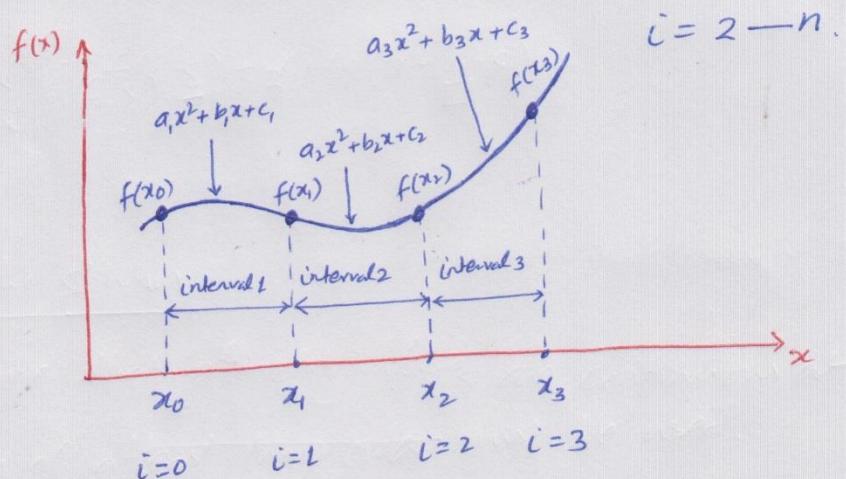
For  $(n+1)$  data points, there are  $n$  intervals and  $3n$  unknowns.

- To evaluate  $3n$  unknowns,  $3n$  conditions are required:

- a) The function values of adjacent polynomials must be equal at the interior knots:

$$a_{i-1} x_{i-1}^2 + b_{i-1} x_{i-1} + c_{i-1} = f(x_{i-1}) \quad (4)$$

$$a_i x_i^2 + b_i x_i + c_i = f(x_i) \quad (5)$$



- Eqs. (4) and (5) each provide  $n-1$  conditions for a total of  $2n-2$  conditions.

4

b) The first and last functions must pass through the end points:

$$a_1 x_0^2 + b_1 x_0 + c_1 = f(x_0) \quad \text{--- (6)}$$

$$a_n x_n^2 + b_n x_n + c_n = f(x_n) \quad \text{--- (7)}$$

— 2 conditions

Total  $2n$  conditions.

c) The first derivatives at the interior knots must be equal.

— The first derivative of (3) is

$$f'(x) = 2ax + b$$

∴ the conditions are:

$$2a_{i-1} x_{i-1} + b_{i-1} = 2a_i x_i + b_i \quad \text{--- (8)}$$

$$i=2-n$$

Eq. (8) provides another  $n-1$  conditions for a total of  $2n + n-1 = 3n-1$  conditions.

d) Assume that the second derivative is zero at the first point

$$\Rightarrow a_1 = 0 \quad \text{--- (9)}$$

That is the first two points will be connected by a straight line.

□

15

Example 2:

Fit quadratic splines for the data given in Example-1.

Also use the result to estimate the value at  $x=5$ .

— We are given with 4 data points and 3 intervals.

$\therefore 3n = 9$  unknowns must be determined.

Eqs. (4) and (5) yields 4 conditions :

$$20 \cdot 25 a_1 + 4 \cdot 5 b_1 + c_1 = 1 \cdot 0$$

$$20 \cdot 25 a_2 + 4 \cdot 5 b_2 + c_2 = 1 \cdot 0$$

$$49 a_2 + 7 b_2 + c_2 = 2 \cdot 5$$

$$49 a_3 + 7 b_3 + c_3 = 2 \cdot 5$$

Eqs. (6) and (7) gives :

$$9 a_1 + 3 b_1 + c_1 = 2 \cdot 5$$

$$81 a_3 + 9 b_3 + c_3 = 2 \cdot 5$$

Continuity of derivatives creates : (Eq. (8)) :

$$9 a_1 + b_1 = 9 a_2 + b_2$$

$$14 a_2 + b_2 = 14 a_3 + b_3$$

Finally, Eq. (9) gives:  $a_1 = 0$ . So, we have to solve (only) eight simultaneous equations :

— In matrix form :

16

$$\left[ \begin{array}{ccccccccc|c} 4.5 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & b_1 \\ 0 & 0 & 20.25 & 4.5 & 1 & 0 & 0 & 0 & 0 & c_1 \\ 0 & 0 & 49 & 7 & 1 & 0 & 0 & 0 & 0 & a_2 \\ 0 & 0 & 0 & 0 & 0 & 49 & 7 & 1 & b_2 \\ 3 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & c_2 \\ 0 & 0 & 0 & 0 & 0 & 81 & 9 & 1 & a_3 \\ 1 & 0 & -9 & -1 & 0 & 0 & 0 & 0 & b_3 \\ 0 & 0 & 14 & 1 & 0 & -14 & -1 & 0 & c_3 \end{array} \right] = \begin{bmatrix} 1 \\ 1 \\ 2.5 \\ 2.5 \\ 2.5 \\ 0.5 \\ 0 \\ 0 \end{bmatrix}$$

Solving this:

$$a_1 = 0$$

$$b_1 = -1$$

$$c_1 = 5.5$$

$$a_2 = 0.64$$

$$b_2 = -6.76$$

$$c_2 = 18.46$$

$$a_3 = -1.6$$

$$b_3 = 24.6$$

$$c_3 = -91.3$$

Substituting these ~~into~~ into original quadratic equations to develop the following relationships for each interval:

$$f_1(x) = -x + 5.5, \quad 3.0 \leq x \leq 4.5$$

$$f_2(x) = 0.64x^2 - 6.76x + 18.46, \quad 4.5 \leq x \leq 7.0$$

$$f_3(x) = -1.6x^2 + 24.6x - 91.3, \quad 7.0 \leq x \leq 9.0$$

∴ Use  $f_2$  to get the estimated value at  $x=5$ , which is:

$$f_2(5) = 0.66.$$

□

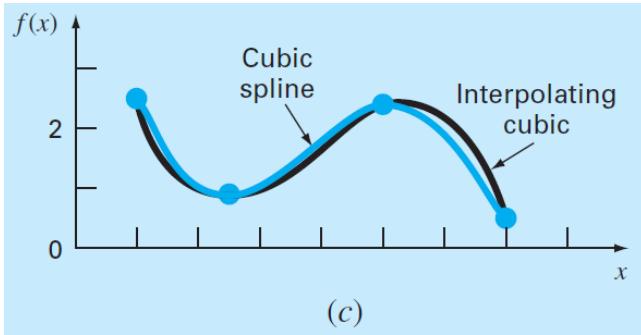
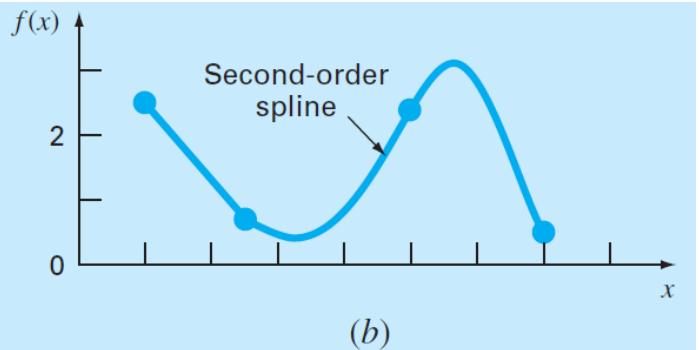
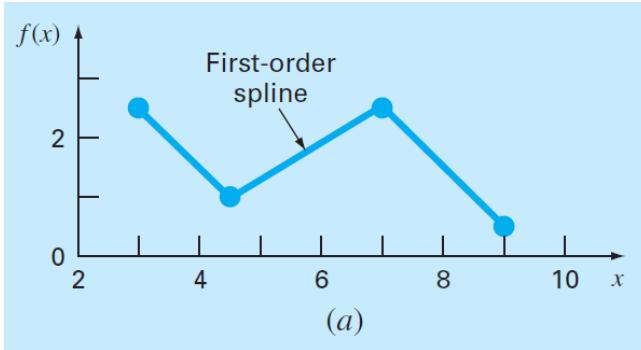


Figure 4: Spline fits of a set of four points (as discussed in Examples 1-3): (a) linear splines; (b) quadratic spline; (c) cubic spline (also with cubic interpolating polynomial)

## Cubic splines :

L7

- Derive a third-order polynomial for each interval between knots as:

$$f_i(x) = a_i x^3 + b_i x^2 + c_i x + d_i \quad \text{--- (10)}$$

∴ for  $(n+1)$  data points with  $n$  intervals, there are  $4n$  unknowns  $\Rightarrow 4n$  conditions are required.

1. The function values must be equal at the interior knots  
( $2n-2$  conditions)
2. The first and last function must pass through the end points  
(2 conditions)
3. The first derivatives at the interior knots must be equal  
( $n-1$  conditions)
4. The second derivatives at the interior knots must be equal  
( $n-1$  conditions)
5. The second derivative at the end knots are zero  
(2 conditions)

Example 3: Fit cubic splines to the same data used in Examples 1 and 2. Estimate value at  $x=5$ .

- As each pair of knots is connected by a cubic polynomial  
 $\Rightarrow$  second derivative within each interval is a straight line.  
 $\therefore$  second derivative can be represented by a first-order (linear) Lagrange interpolating polynomial as:

$$f_i''(x) = f_i''(x_{i-1}) \left( \frac{x - x_i}{x_{i-1} - x_i} \right) + f''(x_i) \left( \frac{x - x_{i-1}}{x_i - x_{i-1}} \right) \quad (1)$$

Where  $f_i''(x)$  — value of the second-order derivative at any  $x$  within the  $i$ th interval.

- Eq. (1) is a straight line connecting the second derivative at the first knot  $f''(x_{i-1})$  with the second derivative at the second knot  $f''(x_i)$ .

- Integrate Eq. (1) twice to yield an expression for  $f_i(x)$ . Note that the expression will have two constants of integration and to evaluate them use:

$$f(x) = f(x_{i-1}) \quad \text{at } x = x_{i-1}$$

$$\text{and} \quad f(x) = f(x_i) \quad \text{at } x = x_i$$

$$\Rightarrow f_i(x) = \frac{f_i''(x_{i-1})}{6(x_i - x_{i-1})} (x_i - x)^3 + \frac{f_i''(x_i)}{6(x_i - x_{i-1})} (x - x_{i-1})^3 +$$

[2]

$$\left[ \frac{f(x_{i-1})}{(x_i - x_{i-1})} - \frac{f''(x_{i-1})(x_i - x_{i-1})}{6} \right] (x_i - x) \\ + \left[ \frac{f(x_i)}{(x_i - x_{i-1})} - \frac{f''(x_i)(x_i - x_{i-1})}{6} \right] (x - x_{i-1}) \quad \text{--- (12)}$$

The second derivatives appearing in Eq. (12) can be evaluated by using the condition that the first derivatives at the knots must be continuous:

$$f'_i(x_i) = f'_{i+1}(x_i) \quad \text{--- (13)}$$

Now differentiate Eq. (12) and use condition (13) for both the  $(i-1)^{\text{th}}$  and  $i^{\text{th}}$  intervals:

$$\Rightarrow (x_i - x_{i-1}) f''(x_{i-1}) + 2(x_{i+1} - x_{i-1}) f''(x_i) \\ + (x_{i+1} - x_i) f''(x_{i+1}) \\ = \frac{6}{(x_{i+1} - x_i)} [f(x_{i+1}) - f(x_i)] + \frac{6}{(x_i - x_{i-1})} [f(x_{i-1}) - f(x_i)] \quad \text{--- (14)}$$

— Note that if Eq. (14) is written for all interior knots  $(n-1)$  then  $(n-1)$  simultaneous equations result with  $(n+1)$  unknown second derivatives.

— Since for (natural) Cubic splines, the second derivative at the end knots are zero.

$\therefore$  Problem reduces to  $(n-1)$  equations in  $(n-1)$  unknowns.  $\square$

3

In summary for cubic splines we have the following:

- Cubic equation for each interval:

$$\begin{aligned}
 f_i(x) &= \frac{f_i''(x_{i-1})}{6(x_i - x_{i-1})} (x_i - x)^3 + \frac{f_i''(x_i)}{6(x_i - x_{i-1})} (x - x_{i-1})^3 \\
 &+ \left[ \frac{f(x_{i-1})}{(x_i - x_{i-1})} - \frac{f''(x_{i-1})(x_i - x_{i-1})}{6} \right] (x_i - x) \\
 &+ \left[ \frac{f(x_i)}{(x_i - x_{i-1})} - \frac{f''(x_i)(x_i - x_{i-1})}{6} \right] (x - x_{i-1}) \quad \text{--- (A)}
 \end{aligned}$$

- Two unknowns -  $f''(x_{i-1})$  and  $f''(x_i)$  - at the end of each interval - are evaluated using:

$$\begin{aligned}
 (x_i - x_{i-1})f''(x_{i-1}) + 2(x_{i+1} - x_{i-1})f''(x_i) + (x_{i+1} - x_i)f''(x_{i+1}) \\
 = \frac{6}{(x_{i+1} - x_i)} [f(x_{i+1}) - f(x_i)] + \frac{6}{(x_i - x_{i-1})} [f(x_{i-1}) - f(x_i)] \quad \text{--- (B)}
 \end{aligned}$$

Remark:

Equation (B) generates  $(n-1)$  equations in  $(n-1)$  unknowns for interior knots (as second derivatives at the end knots are zero)

Example 3: Fit cubic spline to the same data (as in Example 1). Estimate the value at  $x=5$ .

14

$x$	$f(x)$
3.0	2.5
4.5	1.0
7.0	2.5
9.0	0.5

First step: Employ (B) to get <sup>the</sup> set of simultaneous equations to determine the second derivatives at the knots.

For the first interior knot:

$$x_0 = 3 \quad f(x_0) = 2.5$$

$$x_1 = 4.5 \quad f(x_1) = 1.0$$

$$x_2 = 7 \quad f(x_2) = 2.5$$

Using (B):

$$(4.5 - 3)f''(3) + 2(7 - 3)f''(4.5) + (7 - 4.5)f''(7) \\ = \frac{6}{(7 - 4.5)}(2.5 - 1) + \frac{6}{(4.5 - 3)}(2.5 - 1)$$

Because of the second derivative is zero at the end knot:

$$\Rightarrow f''(3) = 0$$

$$8f''(4.5) + 2.5f''(7) = 9.6 \quad \text{--- (i)}$$

Similarly, (B) can be applied to the second interior point (knot) (which is  $x_2$ )

$$\Rightarrow 2.5f''(4.5) + 9f''(7) = -9.6 \quad \text{--- (ii)}$$

15

Solving (i) and (ii) :

$$\left. \begin{aligned} f''(4.5) &= 1.67909 \\ f''(7) &= -1.53308 \end{aligned} \right]$$

These can be substituted in (A) :

$$\begin{aligned} f_1(x) &= \frac{1.67909}{6(4.5-3)} (x-3)^3 + \frac{2.5}{(4.5-3)} (4.5-x) \\ &\quad + \left[ \frac{1}{(4.5-3)} - \frac{1.67909(4.5-3)}{6} \right] (x-3) \end{aligned}$$

OR

$$f_1(x) = 0.186566(x-3)^3 + 1.666667(4.5-x) + 0.246894(x-3)$$

— which is the cubic spline for the first interval.

Similarly, we can get equations for second and third intervals :

$$f_2(x) = 0.111939(7-x)^3 - 0.102205(x-4.5)^3 - 0.29962(7-x)$$

and

$$f_3(x) = -0.127757(9-x)^3 + 1.761027(9-x) + 0.25(x-7)$$

$\therefore$  the value at  $x=5$  (which falls within the second interval)  
is :

$$f_2(5) = 1.102886$$

□

## Numerical Integration

L1

— Find

$$I = \int_a^b f(x) dx \cong \int_a^b f_n(x) dx \quad \dots \quad ①$$

where  $f_n(x)$  — a polynomial approximation of  $f(x)$

$$f_n(x) = a_0 + a_1 x + \dots + a_n x^n$$

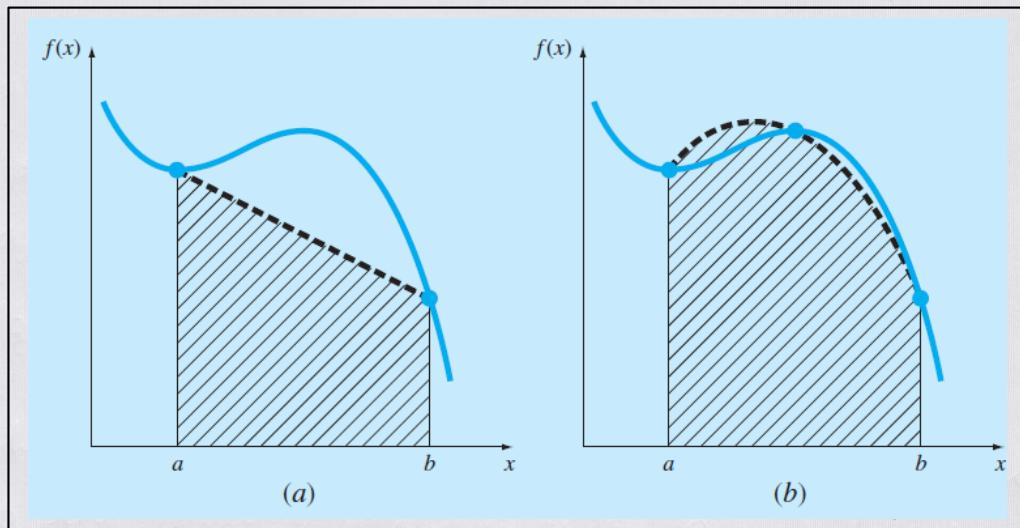


Fig. 1: a.) a first-order polynomial (a straight line) is used as an approximation of  $f(x)$ .

b.) a parabola is used for the approximation of given function  $f(x)$ .

## Newton-Cotes integration formulas:

- Based on the strategy of replacing a complicated function or tabulated data with an approximating function.
- - i) Closed form Newton-Cotes formulas.
  - ii) Open form Newton-Cotes formulas.
- The closed form are those where the data points at the beginning and end of the limits of integration are known.
- The open forms have integration limits that extend beyond the range of the data. (Generally, not preferred for definite integrals).

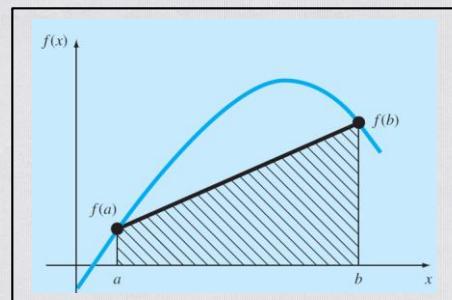
## Closed form Newton-Cotes formulas:

### A) Trapezoidal Rule:

This corresponds to the case where the polynomial in Eq.(1) is of first-order (linear).

That is:

$$\int_a^b f(x) dx \approx \int_a^b f_1(x) dx \quad \text{--- (2)}$$



3

— Straight line joining (two) data points:

$$f_1(x) = f(a) + \frac{f(b) - f(a)}{(b-a)}(x-a) \quad \text{--- (3)}$$

— Area under this straight line is an estimate of the

$$\int_a^b f(x) dx.$$

$$\Rightarrow I \cong \int_a^b \left[ f(a) + \frac{f(b) - f(a)}{(b-a)}(x-a) \right] dx$$

OR

$$I \cong (b-a) \left[ \frac{f(a) + f(b)}{2} \right] \quad \text{--- (4)}$$

Remark:

1. Geometrically, from fig. 2 :

$$\begin{aligned} I &\cong \text{width} * \text{average height} \\ &= (b-a) * \left( \frac{f(b) + f(a)}{2} \right) \end{aligned} \quad \text{--- (5)}$$

2. All the Newton-Cotes formulas can be expressed in the general format as given in Eq. (5). In fact they differ only with respect to the formulation of average height.

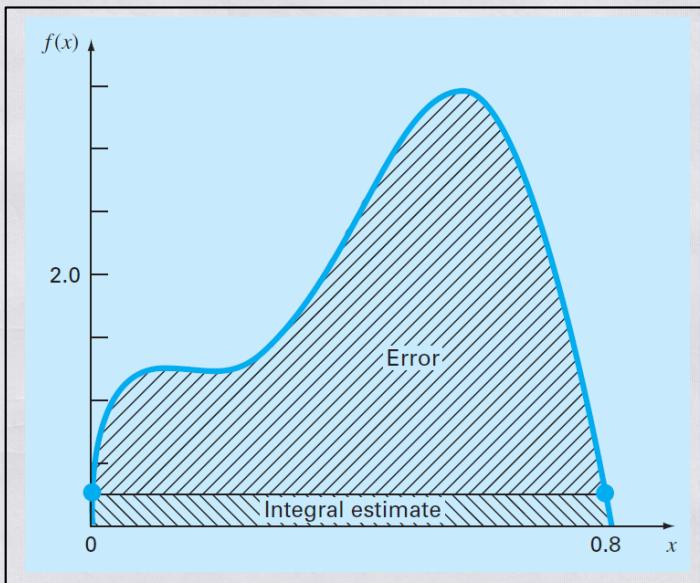
(4)

### Error of the Trapezoidal Rule.

- Since we are using linear approximation for the function, therefore, the error will be substantial.
- For example, choose a function:

$$f(x) = 0.2 + 25x - 200x^2 + 675x^3 - 900x^4 + 400x^5$$

with  $a=0$  and  $b=0.8$



- An estimate for the local truncation error for a single application of the trapezoidal rule is:

$$E_t = -\frac{1}{12} f''(\xi)(b-a)^3 \quad — (6)$$

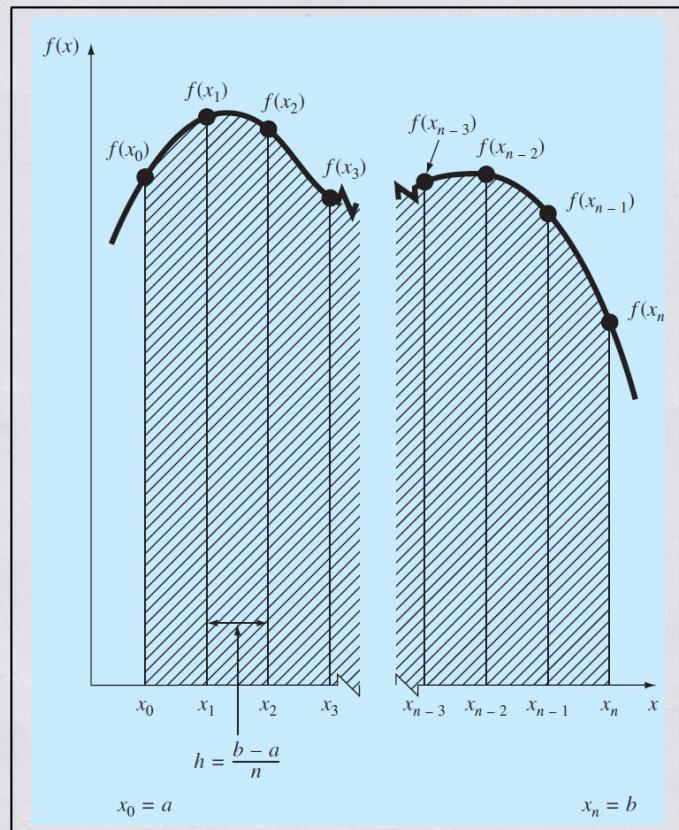
where  $\xi$  lies between  $a$  and  $b$ .

- Equation (6) indicates that trapezoidal rule is exact if  $f(x)$  is linear.

## The multiple application of Trapezoidal Rule

- To improve the accuracy of the trapezoidal rule - divide the integration interval from  $a$  to  $b$  into a number of segments and apply method to each segment:
- If there are  $(n+1)$  equally-spaced base points:  $x_0, x_1, \dots, x_n$  then there are  $n$ -segments of equal width:

$$h = \frac{b-a}{n}$$



- The total integral:

$$I = \int_{x_0}^{x_1} f(x) dx + \int_{x_1}^{x_2} f(x) dx + \dots + \int_{x_{n-1}}^{x_n} f(x) dx$$

16

— Using trapezoidal rule in each segment:

$$I \cong h \left( \frac{f(x_0) + f(x_1)}{2} \right) + h \left( \frac{f(x_1) + f(x_2)}{2} \right) + \dots + h \left( \frac{f(x_{n-1}) + f(x_n)}{2} \right)$$

OR  $= \frac{h}{2} \left[ f(x_0) + 2 \sum_{i=1}^{n-1} f(x_i) + f(x_n) \right] \quad \text{--- (7)}$

$$= (b-a) * \left( \frac{f(x_0) + 2 \sum_{i=1}^{n-1} f(x_i) + f(x_n)}{2n} \right)$$

[ ] [ ]  
 width                      average height

— Error expression for the multiple application of trapezoidal rule:

$$E_t = - \frac{(b-a)^3}{12n^3} \sum_{i=1}^n f''(\bar{x}_i) \quad \text{--- (8)}$$

Where  $\bar{x}_i$  is located in the segment  $i$ . Equation (8) can be simplified by taking the mean value of the second-derivative as:

$$\bar{f}'' \cong \frac{\sum_{i=1}^n f''(\bar{x}_i)}{n} \Rightarrow \sum_{i=1}^n f''(\bar{x}_i) \cong n\bar{f}''$$

$$\therefore \boxed{E_a = - \frac{(b-a)^3}{12n^2} \bar{f}''} \quad \text{--- (9)}$$

L7

- To improve the efficiency of the methods — Employ higher order polynomial approximation of  $f(x)$ .

### Simpson's Rule

- Aside from applying the trapezoidal rule with finer segmentation, another way to obtain a more accurate estimate of an integral is to use higher-order polynomial to connect the points:

#### A.) Simpson's $\frac{1}{3}$ Rule

a second-order interpolating polynomial is substituted in Eq. (1):

$$I = \int_a^b f(x) dx \cong \int_a^b f_2(x) dx$$

Let  $a = x_0$ ,  $b = x_2$  (as three points:  $x_0$ ,  $x_1$  &  $x_2$  are required for a second-order fit)

$$\begin{aligned} I &\cong \int_{x_0}^{x_2} \left[ \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} f(x_0) + \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} f(x_1) \right. \\ &\quad \left. + \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} f(x_2) \right] dx \end{aligned}$$

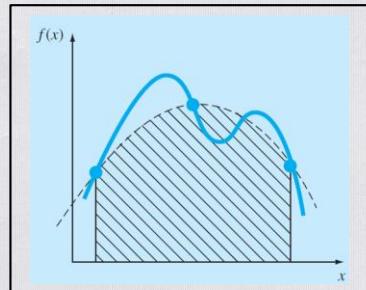
After integration:

19

$$\boxed{I \approx \frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)]} \quad (10)$$

Where  $h = \frac{b-a}{2}$

- Eq. (10) can also be represented as:



$$I \approx (b-a) \left( \frac{f(x_0) + 4f(x_1) + f(x_2)}{6} \right)$$

width                          average height

Note:

Here  $a = x_0$ ,  $b = x_2$  and  $x_1$  is the mid-point, which is  $\left(\frac{b+a}{2}\right)$ .

- Single application of Simpson's  $\frac{1}{3}$  rule has a truncation error;

$$E_T = -\frac{1}{90} h^5 f^{(4)}(\xi)$$

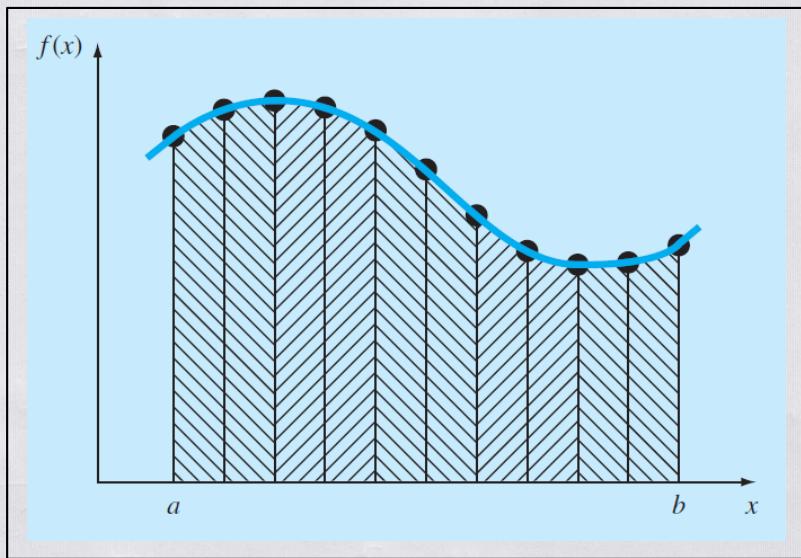
Where  $\xi$  lies in the interval  $a$  to  $b$ . As,  $h = \frac{b-a}{2}$

$$\Rightarrow \boxed{E_T = -\frac{(b-a)^5}{2880} f^{(4)}(\xi)} \quad (11)$$

[9.]

### Multiple application of Simpson's $\frac{1}{3}$ rule :

- Similar to trapezoidal rule, Simpson's rule can be improved by dividing interval into a number of segments of equal width :  
$$h = \frac{b-a}{n}$$
- Note that the method can be employed only if the number of segments is even.



— The total integral :  $I = \int_{x_0}^{x_2} f(x) dx + \int_{x_2}^{x_4} f(x) dx + \dots + \int_{x_{n-2}}^{x_n} f(x) dx$

Using Simpson's  $\frac{1}{3}$  rule for the individual integral

$$I \cong 2h \left( \frac{f(x_0) + 4f(x_1) + f(x_2)}{6} \right) + \dots + 2h \left( \frac{f(x_{n-2}) + 4f(x_{n-1}) + f(x_n)}{6} \right)$$

1/10

OR

$$I \cong (b-a) \left( f(x_0) + 4 \sum_{i=1,3,5}^{n-1} f(x_i) + 2 \sum_{j=2,4,6}^{n-2} f(x_j) + f(x_n) \right)$$

width                          average height

- Truncation error:

$$\boxed{E_a = -\frac{(b-a)^5}{180n^4} f^{(4)}}$$

□

## Simpson's 3/8 Rule

11

- As Simpson's  $\frac{1}{3}$  rule is limited to situations where we have an even number of segments (and an odd number of points).
  - Therefore, Simpson's  $\frac{1}{3}$  rule is — an odd-segment even-point formula.
- An odd-segment even-point formula — Simpson's  $\frac{3}{8}$  rule — is used in conjunction with the  $\frac{1}{3}$  rule to permit evaluation of both even and odd number of segments.

$$I = \int_a^b f(x) dx \approx \int_a^b f_3(x) dx$$

↓      → Third Lagrange interpolating polynomial (which require four data points)

$$\begin{matrix} x_0 & x_1 & x_2 & x_3 \\ \parallel & & & \parallel \\ a & & & b \end{matrix}$$

$$\begin{aligned}
 I &\stackrel{b=x_3}{\approx} \int_{x_0}^{x_3} \left[ \frac{(x-x_0)(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)(x_0-x_3)} f(x_0) + \frac{(x-x_0)(x-x_1)(x-x_3)}{(x_1-x_0)(x_1-x_2)(x_1-x_3)} f(x_1) \right. \\
 &\quad \left. + \frac{(x-x_0)(x-x_1)(x-x_3)}{(x_2-x_0)(x_2-x_1)(x_2-x_3)} f(x_2) + \frac{(x-x_0)(x-x_1)(x-x_2)}{(x_3-x_0)(x_3-x_1)(x_3-x_2)} f(x_3) \right] dx
 \end{aligned}$$

Int. w.r.t - x

12

$$\boxed{I \cong \frac{3h}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)]} \quad \text{--- } ①$$

Where  $h = \frac{(b-a)}{3}$

OR

$$I \cong (b-a) \left( \frac{f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)}{8} \right)$$

width                      average height

— Simpson 3/8 rule has an error of

$$E_I = -\frac{3}{80} h^5 f^{(4)}(\bar{x}) \quad , \quad a < \bar{x} < b$$

OR

$$\boxed{E_I = -\frac{(b-a)^5}{6480} f^{(4)}(\bar{x})} \quad \text{--- } ②$$

## Simpson's 3/8 Rule

### Problem Statement.

- (a) Use Simpson's 3/8 rule to integrate

$$f(x) = 0.2 + 25x - 200x^2 + 675x^3 - 900x^4 + 400x^5$$

from  $a = 0$  to  $b = 0.8$ .

- (b) Use it in conjunction with Simpson's 1/3 rule to integrate the same function for five segments.

### Solution.

- (a) A single application of Simpson's 3/8 rule requires four equally spaced points:

$$f(0) = 0.2 \quad f(0.2667) = 1.432724$$

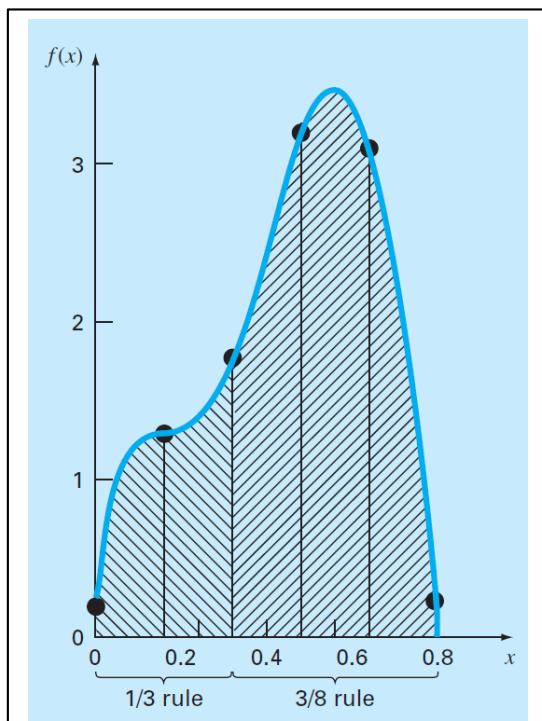
$$f(0.5333) = 3.487177 \quad f(0.8) = 0.232$$

Using Eq. (1):

$$I \cong 0.8 \frac{0.2 + 3(1.432724 + 3.487177) + 0.232}{8} = 1.519170$$

$$E_t = 1.640533 - 1.519170 = 0.1213630 \quad \epsilon_t = 7.4\%$$

$$E_a = -\frac{(0.8)^5}{6480}(-2400) = 0.1213630$$



(b) The data needed for a five-segment application ( $h = 0.16$ ) is

$$\begin{array}{ll} f(0) = 0.2 & f(0.16) = 1.296919 \\ f(0.32) = 1.743393 & f(0.48) = 3.186015 \\ f(0.64) = 3.181929 & f(0.80) = 0.232 \end{array}$$

The integral for the first two segments is obtained using Simpson's 1/3 rule:

$$I \cong 0.32 \frac{0.2 + 4(1.296919) + 1.743393}{6} = 0.3803237$$

For the last three segments, the 3/8 rule can be used to obtain

$$I \cong 0.48 \frac{1.743393 + 3(3.186015 + 3.181929) + 0.232}{8} = 1.264754$$

The total integral is computed by summing the two results:

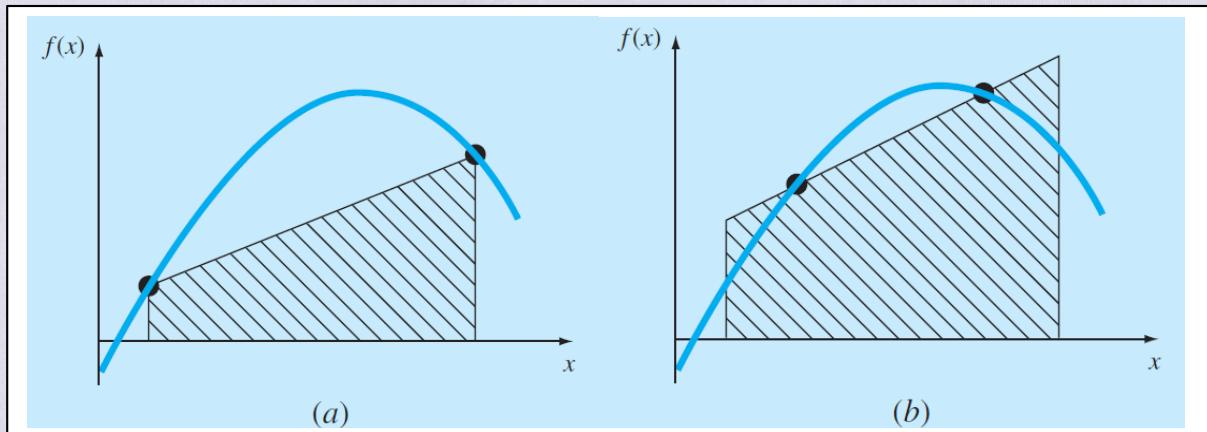
$$I = 0.3803237 + 1.264753 = 1.645077$$

$$E_t = 1.640533 - 1.645077 = -0.00454383 \quad \varepsilon_t = -0.28\%$$

13

## Gauss- Quadrature

- So far we have studied numerical integration or quadrature formulas known as the Newton-Cotes formulas. In these formulas, the integral estimate was based on evenly spaced function values. Therefore, the location of the base points used in these equations was fixed.
- Now, suppose that the constraint of fixed base points was removed and we were free to evaluate the area under a straight line joining any two points on the curve.
  - That is can we improve Trapezoidal rule ?



14

- Gauss-quadrature is the class of technique to implement such a strategy.
- We will discuss about Gauss-Legendre formulas.

- Method of Undetermined Coefficients:

↳ To derive the Trapezoidal rule ( $I \approx (b-a) \left( \frac{f(b)+f(a)}{2} \right)$ )

Express Eq-(i) as:

$$I \approx C_0 f(a) + C_1 f(b) \quad \text{--- (ii)} \quad C_i's \text{ constant}$$

- As trapezoidal rule yields exact result if the integrand is a constant or straight line.

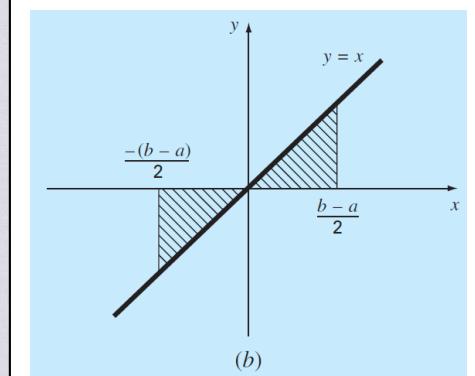
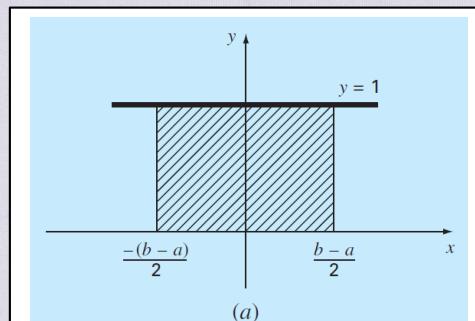
- Simple equations that represent these cases are:  $y=1$  and  $y=x$ .

$$C_0 + C_1 = \int_1 dx \\ - \left( \frac{b-a}{2} \right)$$

and

$$-C_0 \left( \frac{b-a}{2} \right) + C_1 \left( \frac{b-a}{2} \right) = \int x dx \\ - \left( \frac{b-a}{2} \right)$$

$$\Rightarrow \begin{cases} C_0 + C_1 = (b-a) \\ -C_0 \left( \frac{b-a}{2} \right) + C_1 \left( \frac{b-a}{2} \right) = 0 \end{cases}$$



$$\Rightarrow C_0 = C_1 = \frac{b-a}{2}$$

15

Substitute in (ii) :

$$I \equiv \left( \frac{b-a}{2} \right) f(a) + \left( \frac{b-a}{2} \right) f(b)$$

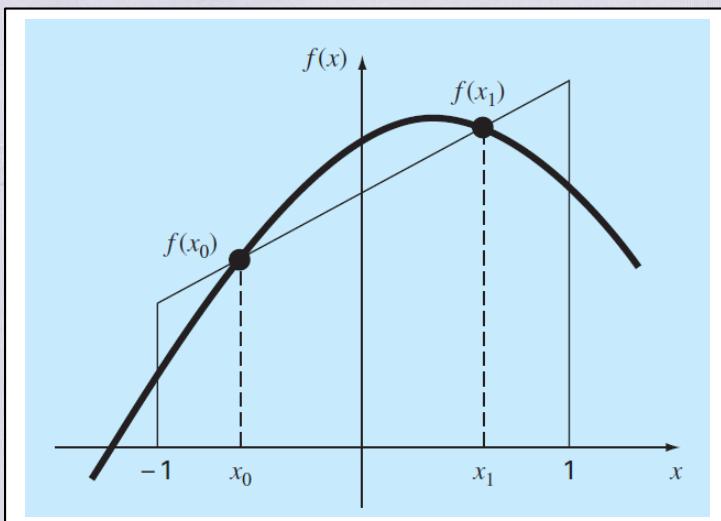
$$= (b-a) \left( \frac{f(b) + f(a)}{2} \right)$$

Derivation of two-point Gauss-Legendre formula:

— (Similar to the above derivation of trapezoidal rule) The object of Gauss-quadrature is to determine the coefficient of an equation of the form:

$$I \equiv C_0 f(x_0) + C_1 f(x_1) \quad \text{--- (3)}$$

Here, unlike the trapezoidal rule that used fixed end points  $a$  and  $b$ , the function argument  $x_0$  and  $x_1$  are not fixed at the end points, but are unknowns.



16

— Here, we have total of four unknowns.

$\Rightarrow$  We require four conditions to determine them exactly.

— Here we assume that the integrals of a constant, linear, parabolic and cubic functions are exact.

$$y=1, \quad y=x, \quad y=x^2, \quad y=x^3$$

$$\begin{aligned} C_0 f(x_0) + C_1 f(x_1) &= \int_{-1}^1 1 dx = 2 & C_0 + C_1 = 2 \\ C_0 f(x_0) + C_1 f(x_1) &= \int_{-1}^1 x dx = 0 & C_0 x_0 + C_1 x_1 = 0 \\ C_0 f(x_0) + C_1 f(x_1) &= \int_{-1}^1 x^2 dx = \frac{2}{3} & C_0 x_0^2 + C_1 x_1^2 = \frac{2}{3} \\ C_0 f(x_0) + C_1 f(x_1) &= \int_{-1}^1 x^3 dx = 0 & C_0 x_0^3 + C_1 x_1^3 = 0 \end{aligned}$$

On solving these four equations:

$$\left. \begin{aligned} C_0 &= C_1 = 1 \\ x_0 &= -\frac{1}{\sqrt{3}}, \quad x_1 = \frac{1}{\sqrt{3}} \end{aligned} \right\} \text{Substitute in (3)}$$

$$\boxed{I \approx f\left(-\frac{1}{\sqrt{3}}\right) + f\left(\frac{1}{\sqrt{3}}\right)}$$

— Gauss-Legendre Two-point formula □

## Gauss-quadrature

L1

- Note that for the Gauss-Legendre quadrature rule integration limits in eqns. are from -1 to 1.
- A simple change of variables can be used to translate other limits of integration into this form.
- Let  $x_d$  be the new variable and  $x$  old variable.  $x_d$  and  $x$  are related as:

$$x = a_0 + a_1 x_d \quad \text{--- (1)}$$

such that  $x=a$  then  $x_d=-1 \Rightarrow a = a_0 + a_1(-1) \quad \text{--- (2)}$

And

$$x=b \text{ then } x_d=1 \Rightarrow b = a_0 + a_1(1) \quad \text{--- (3)}$$

From (2) and (3):  $a_0 = \frac{b+a}{2}$  and  $a_1 = \frac{b-a}{2}$

∴ From (1)

$$x = \frac{(b+a) + (b-a)x_d}{2}$$

$$dx = \left(\frac{b-a}{2}\right) dx_d$$

→  $x$  and  $dx$   
can be substituted  
before applying  
Gauss-Legendre rule

$$\int_a^b f(x) dx$$

## Two-Point Gauss-Legendre Formula

**Problem Statement.** Use G-L rule to evaluate the integral of

$$f(x) = 0.2 + 25x - 200x^2 + 675x^3 - 900x^4 + 400x^5$$

between the limits  $x = 0$  to  $0.8$ . Recall that this was the same problem that we solved in Chap. 21 using a variety of Newton-Cotes formulations. The exact value of the integral is  $1.640533$ .

**Solution.** Before integrating the function, we must perform a change of variable so that the limits are from  $-1$  to  $+1$ . To do this, we substitute  $a = 0$  and  $b = 0.8$  into Eq. (1) to yield

$$x = 0.4 + 0.4x_d$$

The derivative of this relationship is

$$dx = 0.4 dx_d$$

Both of these can be substituted into the original equation to yield

$$\begin{aligned} & \int_0^{0.8} (0.2 + 25x - 200x^2 + 675x^3 - 900x^4 + 400x^5) dx \\ &= \int_{-1}^1 [0.2 + 25(0.4 + 0.4x_d) - 200(0.4 + 0.4x_d)^2 + 675(0.4 + 0.4x_d)^3 \\ &\quad - 900(0.4 + 0.4x_d)^4 + 400(0.4 + 0.4x_d)^5] 0.4 dx_d \end{aligned}$$

Therefore, the right-hand side is in the form that is suitable for evaluation using Gauss quadrature. The transformed function can be evaluated at  $-1/\sqrt{3}$  to be equal to  $0.516741$

and at  $1/\sqrt{3}$  to be equal to  $1.305837$ . Therefore, the integral according to G-L rule is

$$I \cong 0.516741 + 1.305837 = 1.822578$$

## Error in Gauss-quadrature :

[2]

- The error in the Gauss-Legendre rule is :

$$E_T = \frac{2^{2n+3} [(n+1)!]^4}{(2n+3) [(2n+2)!]^3} f^{(2n+2)}(\bar{z}_n) \quad \text{--- (4)}$$

$n$  = the number of points minus one.

and  $f^{(2n+2)}(\bar{z}_n)$  = the  $(2n+2)^{\text{th}}$  derivative of the function after the change of variable with  $\bar{z}_n$  located somewhere between -1 to 1.

- Eq. (4) indicates the superiority of Gauss quadrature to Newton-Cotes formula.

- From (4), for two-point-Gauss-Legendre rule:

$E_T = \frac{1}{135} f^{(4)}(\bar{z}_n), \quad -1 < \bar{z}_n < 1 \quad \text{--- (5)}$

- Note: For many functions confronted in engineering practice, Gauss-quadrature provides an efficient means for evaluating integrals.

□

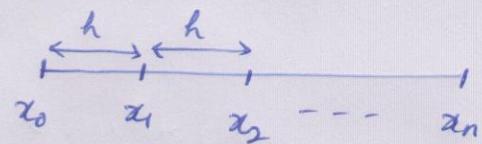
## Numerical Differentiation

3

- Base on Taylor's series expansion.

$$f(x+h) = f(x) + h f'(x) + \frac{h^2}{2!} f''(x) + \dots$$

- Suppose we have equi-distant points:



$$x_{i+1} = x_i + h$$

$$i = 0, \dots, n-1.$$

$$\Rightarrow f(x_{i+1}) = f(x_i + h) = f(x_i) + h f'(x_i) + \frac{h^2}{2!} f''(x_i) + \dots$$

OR

$$f'(x_i) = \frac{f(x_{i+1}) - f(x_i)}{h} - \frac{h}{2} f''(x_i) + O(h^2) \quad \text{--- A}$$

OR

$$f'(x_i) = \frac{f(x_{i+1}) - f(x_i)}{h} + O(h) \quad \text{--- ⑥}$$

- First Forward difference formula  
and  $h$  is the step-size.

- Similarly,

$$f(x_{i-1}) = f(x_i - h) = f(x_i) - h f'(x_i) + \frac{h^2}{2!} f''(x_i) + \dots$$

--- B

14

OR

$$\boxed{f'(x_i) = \frac{f(x_i) - f(x_{i-1})}{h} + O(h)} \quad \text{--- (7)}$$

— First Backward difference formula.

— Again, subtracting (A) from (B) :

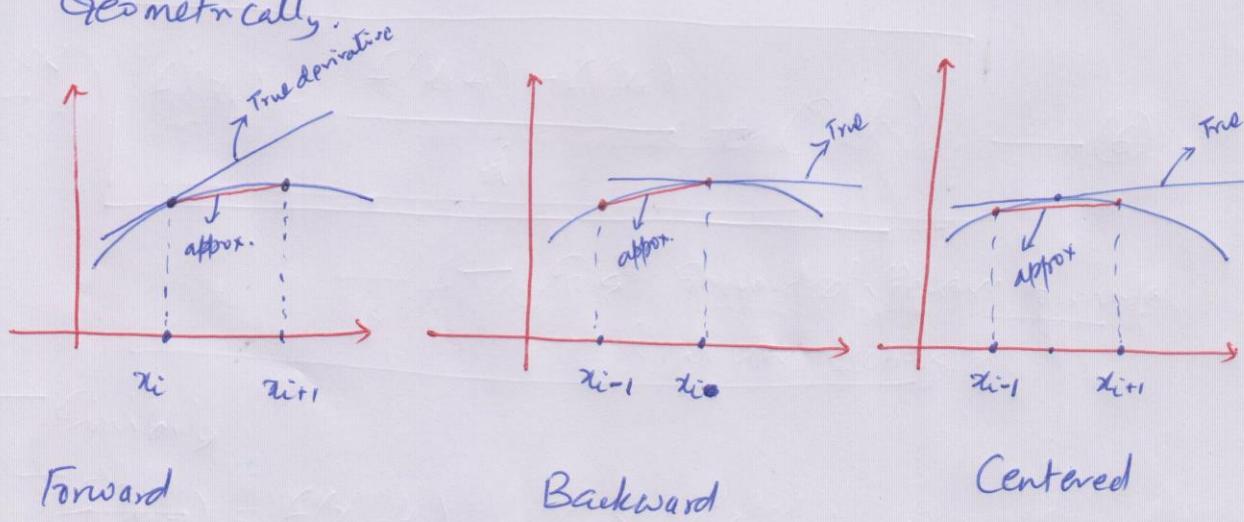
$$f(x_{i+1}) - f(x_{i-1}) = 2hf'(x_i) + \frac{2h^3}{3!} f'''(x_i) + \dots$$

OR

$$\boxed{f'(x_i) = \frac{f(x_{i+1}) - f(x_{i-1})}{2h} + O(h^2)} \quad \text{--- (8)}$$

— Centered difference approximation to the first derivative.

— Geometrically :



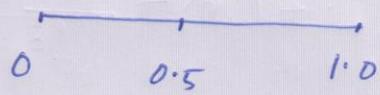
□

Example:  $f(x) = -0.1x^4 - 0.15x^3 - 0.5x^2 - 0.25x + 1.2$

15

Estimate first derivative of  $f(x)$  at  $x=0.5$  with step-size  $h=0.5$  using Forward, Backward & Centred formulas:

— For  $h=0.5$ :



$$x_{i-1} = 0 \Rightarrow f(x_{i-1}) = 1.2$$

$$x_i = 0.5 \Rightarrow f(x_i) = 0.925$$

$$x_{i+1} = 1.0 \Rightarrow f(x_{i+1}) = 0.2$$

Forward:

$$f'(0.5) \cong \frac{f(x_{i+1}) - f(x_i)}{h} = \frac{0.2 - 0.925}{0.5} = -1.45$$

Backward:

$$f'(0.5) \cong \frac{f(x_i) - f(x_{i-1})}{h} = -0.55$$

Centred:

$$f'(0.5) \cong \frac{f(x_{i+1}) - f(x_{i-1})}{2h} = -1.0$$

— Exact value  $f'(0.5) = -0.9125$

□

## Curve Fitting

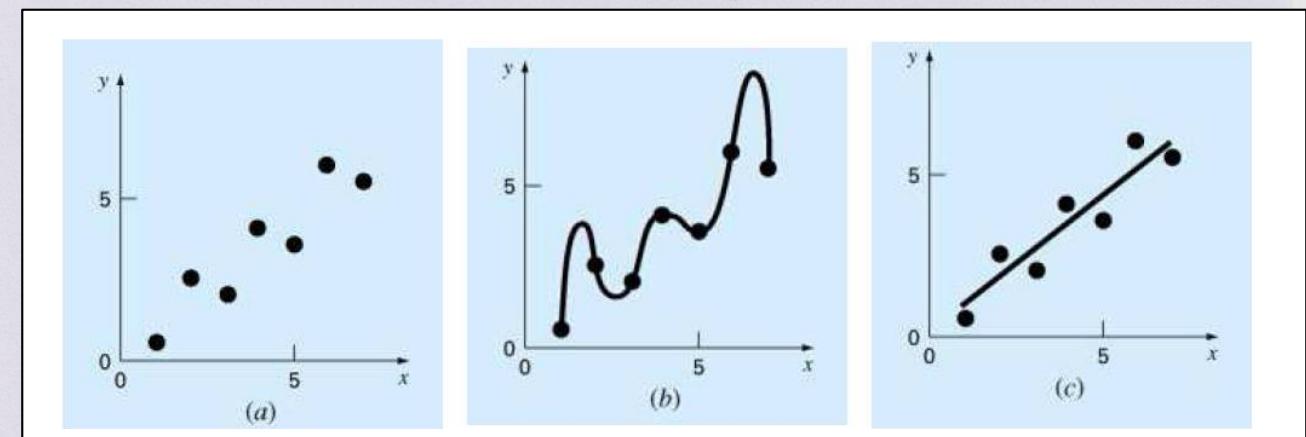
- Two types of curve fitting:

A.) Least square regression :

- For a given data set, derive a single curve that represent the general trend of the data.
- Used when the given data exhibit a significant degree of error or noise.

B.) Interpolation :

- For a given data set, fit a curve or a series of curves that pass directly through each of the data points.
- Used when data are very precise.



12

## - Linear Regression :

- Fitting a straight line to a set of paired observations:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ .

Mathematical expression for the straight line (model)

$$y = a_0 + a_1 x, \quad a_0 - \text{intercept}, a_1 - \text{slope}$$

Define:

$$e_i = y_i|_{\text{measured}} - y_i|_{\text{model}} = y_i - (a_0 + a_1 x_i)$$

Criterion for best fit:

$$\min S_r = \min_{a_0, a_1} \sum_{i=1}^n e_i^2 = \min_{a_0, a_1} \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

Find  $a_0$  and  $a_1$ ?

$$\frac{\partial S_r}{\partial a_0} = -2 \sum_{i=1}^n (y_i - a_0 - a_1 x_i) = 0 \quad \textcircled{1}$$

$$\frac{\partial S_r}{\partial a_1} = -2 \sum_{i=1}^n [(y_i - a_0 - a_1 x_i) x_i] = 0 \quad \textcircled{2}$$

13

From Eq. (1):

$$\sum_{i=1}^n y_i - \sum_{i=1}^n a_0 - \sum_{i=1}^n a_1 x_i = 0$$

OR

$$n a_0 + \sum_{i=1}^n x_i a_1 = \sum_{i=1}^n y_i \quad \text{--- (3)}$$

From Eq. (2):

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n a_0 x_i - \sum_{i=1}^n a_1 x_i^2 = 0$$

OR

$$\sum_{i=1}^n x_i a_0 + \sum_{i=1}^n x_i^2 a_1 = \sum_{i=1}^n x_i y_i \quad \text{--- (4)}$$

- Eqns. (3) and (4) are called normal equations.

From Eq. (3):

$$\begin{aligned} a_0 &= \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n x_i a_1 \\ &= \bar{y} - \bar{x} a_1 \end{aligned}$$

Where,  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  — means of  $x$  &  $y$ .

14

From Eq. (4):

$$\sum_{i=1}^n x_i \left( \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n x_i a_1 \right) + \sum_{i=1}^n x_i^2 a_1 = \sum_{i=1}^n x_i y_i$$

OR

$$a_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2}$$

OR

$$a_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

Definitions :

$$S_y = \sqrt{\sum_{i=1}^n e_i^2} = \sqrt{\sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2}$$

Standard error of the estimate :

$$S_{yx} = \sqrt{\frac{S_y}{n-2}}$$

— spread around the regression line.

15

Standard deviation of data points:

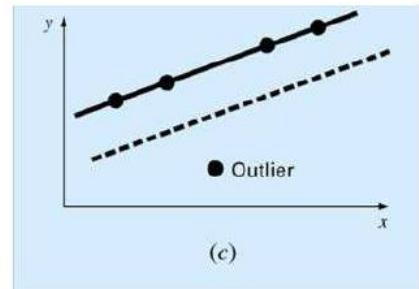
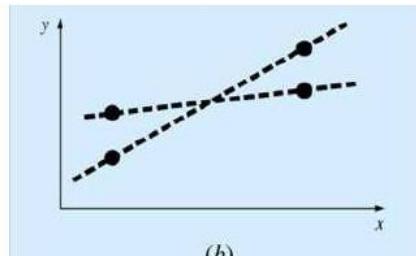
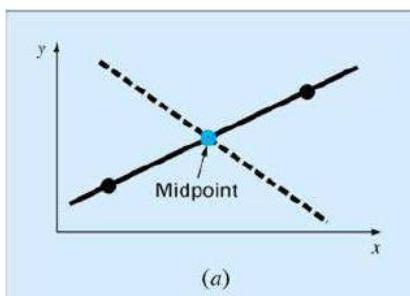
$$S_y = \sqrt{\frac{S_t}{n-1}} = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

where,  $S_t = \sum_{i=1}^n (y_i - \bar{y})^2$ , that is, spread around the mean value  $\bar{y}$ .

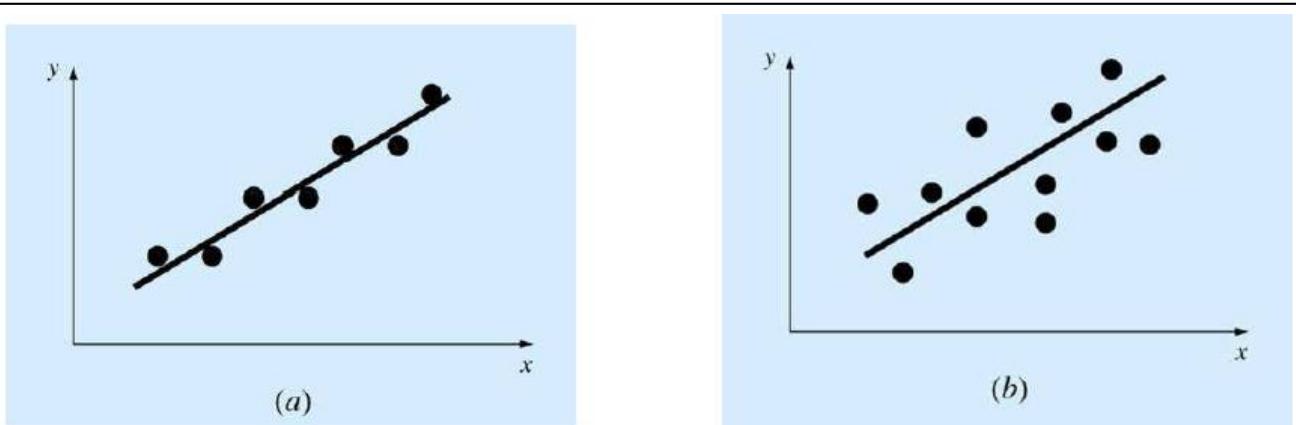
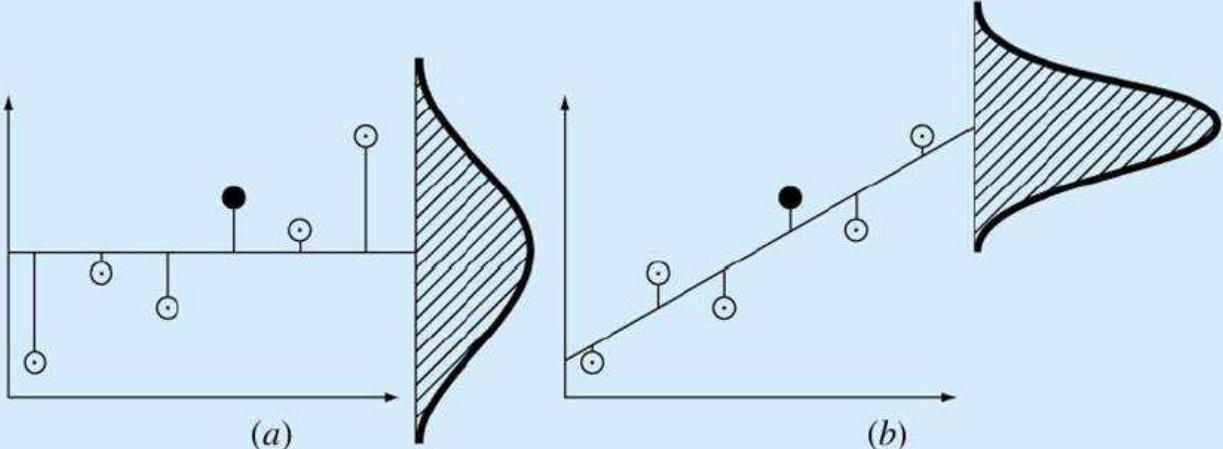
Correlation coefficient:

$$\gamma = \sqrt{\frac{S_t - S_y}{S_t}}$$

- Improvement or error reduction due to describing the data in terms of a straight line rather than as an average value.



Examples of some other criteria for best fit that are inadequate for regression: (a)  $\min \sum e_i$ , (b)  $\min \sum |e_i|$  and (c)  $\min \max e_i$ .



16

Example:

$x$	1	2	3	4	5	6	7
$y$	0.5	2.5	2.0	4.0	3.5	6.0	5.5

$$\sum x_i = 1 + 2 + \dots + 7 = 28$$

$$\sum y_i = 0.5 + 2.5 + \dots + 5.5 = 24$$

$$\sum x_i^2 = 1^2 + 2^2 + \dots + 7^2 = 140$$

$$\sum x_i y_i = 1 \times 0.5 + 2 \times 2.5 + \dots + 7 \times 5.5 = 119.5$$

$$a_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \quad (\text{Here, } n=7)$$

$$= 0.8393$$

$$a_0 = \bar{y} - \bar{x} a_1 = \frac{1}{n} \sum_{i=1}^n y_i - a_1 \frac{1}{n} \sum_{i=1}^n x_i = 0.07143$$

Model:

$$y = 0.07143 + 0.8393 x$$

(7)

$$S_r = \sum_{i=1}^n e_i^2, \quad e_i = y_i - (a_0 + a_1 x_i)$$

$$e_1 = 0.5 - 0.07143 - 0.8393 \times 1 = -0.410$$

$$e_2 = 2.5 - 0.07143 - 0.8393 \times 2 = 0.750$$

$$e_3 = 2.0 - 0.07143 - 0.8393 \times 3 = -0.589$$

⋮

⋮

$$e_7 = 5.5 - 0.07143 - 0.8393 \times 7 = -0.446$$

$$\Rightarrow S_r = (-0.410)^2 + 0.750^2 + (-0.589)^2 + \dots + (-0.446)^2 \\ = 2.9911$$

and

$$S_t = \sum_{i=1}^n (y_i - \bar{y})^2 = 22.714$$

Standard deviation of the data points :

$$S_y = \sqrt{\frac{S_t}{n-1}} = \sqrt{\frac{22.714}{6}} = 1.946$$

Standard error of the estimate :

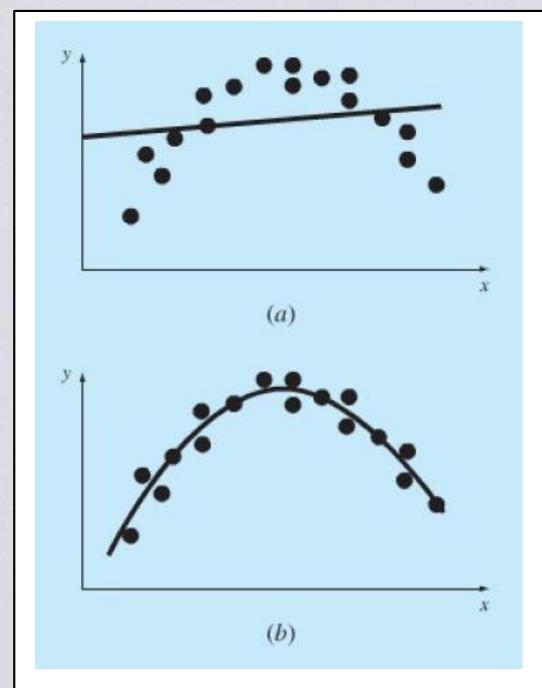
$$S_{y/x} = \sqrt{\frac{S_r}{n-2}} = \sqrt{\frac{2.9911}{5}} = 0.774$$

$$\boxed{\gamma^2 = \frac{S_t - S_r}{S_t}} \\ = 0.868$$

$$S_{y/x} < S_y, S_r < S_t. \text{ Correlation coefficient } \gamma = \sqrt{\frac{S_t - S_r}{S_t}} = 0.932 \quad \square$$

- Polynomial Regression :

- Developed a strategy to derive the equation of straight line using the least-squares criterion.
- Some practical data, although exhibiting a marked pattern (as shown in the fig. below), is poorly represented by a straight line. For such cases, a curve would be better suited to fit the data.
- The alternative is to fit polynomials to the data using polynomial regression. And the least-squares approach can be extended to fit the data to a higher-order polynomial.



- Polynomial Regression:

Given data:  $(x_i, y_i), i=1, 2, \dots, n$

Fit a second-order polynomial:

$$y = a_0 + a_1 x + a_2 x^2 \quad \text{--- (model)}$$

$$e_i = y_i \Big|_{\text{measured}} - y_i \Big|_{\text{model}}$$

$$= y_i - (a_0 + a_1 x_i + a_2 x_i^2)$$

Criterion for best fit:

$$\min S_y = \min_{a_0, a_1, a_2} \sum_{i=1}^n e_i^2$$

$$= \min_{a_0, a_1, a_2} \sum_{i=1}^n (y_i - a_0 - a_1 x_i - a_2 x_i^2)^2$$

Find  $a_0, a_1$  and  $a_2$ :

$$\frac{\partial S_y}{\partial a_0} = -2 \sum_{i=1}^n (y_i - a_0 - a_1 x_i - a_2 x_i^2) = 0 \quad \text{--- (1)}$$

$$\frac{\partial S_y}{\partial a_1} = -2 \sum_{i=1}^n [(y_i - a_0 - a_1 x_i - a_2 x_i^2) x_i] = 0 \quad \text{--- (2)}$$

$$\frac{\partial S_y}{\partial a_2} = -2 \sum_{i=1}^n [(y_i - a_0 - a_1 x_i - a_2 x_i^2) x_i^2] = 0 \quad \text{--- (3)}$$

[3]

OR

$$\left\{ \begin{array}{l} n a_0 + \sum_{i=1}^n x_i a_1 + \sum_{i=1}^n x_i^2 a_2 = \sum_{i=1}^n y_i \quad \text{--- (1)} \\ \sum_{i=1}^n x_i a_0 + \sum_{i=1}^n x_i^2 a_1 + \sum_{i=1}^n x_i^3 a_2 = \sum_{i=1}^n x_i y_i \quad \text{--- (2)} \\ \sum_{i=1}^n x_i^2 a_0 + \sum_{i=1}^n x_i^3 a_1 + \sum_{i=1}^n x_i^4 a_2 = \sum_{i=1}^n x_i^2 y_i \quad \text{--- (3)} \end{array} \right.$$

Three (linear) equations in three unknowns:  $a_0, a_1$ , and  $a_2$

$a_0, a_1$ , and  $a_2$  can be calculated from (1) - (3).

Note: In general, to fit  $p^{\text{th}}$  degree polynomial:

$$y = a_0 + a_1 x + a_2 x^2 + \dots + a_p x^p$$

with least-squares regression approach, one need to solve a system of  $(p+1)$  linear equations.

Standard error:  $S_{y/x} = \sqrt{\frac{S_\sigma}{n - (p+1)}}$

K

## Multiple Linear Regression :

— Used when  $y$  (dependent variable) is a linear function of two or more independent variables.

$$\text{Model: } y = a_0 + a_1 x_1 + a_2 x_2$$

Given data:  $(x_{1i}, x_{2i}, y_i), i = 1, 2, \dots, n$ .

$$e_i = y_i|_{\text{measured}} - y_i|_{\text{model}}$$

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_{1i} - a_2 x_{2i})^2$$

Find  $a_0, a_1$  and  $a_2$  (which minimize  $S_r$ ) :

$$\frac{\partial S_r}{\partial a_0} = -2 \sum_{i=1}^n (y_i - a_0 - a_1 x_{1i} - a_2 x_{2i}) = 0 \quad \text{--- (A)}$$

$$\frac{\partial S_r}{\partial a_1} = -2 \sum_{i=1}^n [(y_i - a_0 - a_1 x_{1i} - a_2 x_{2i}) x_{1i}] = 0 \quad \text{--- (B)}$$

$$\frac{\partial S_r}{\partial a_2} = -2 \sum_{i=1}^n [(y_i - a_0 - a_1 x_{1i} - a_2 x_{2i}) x_{2i}] = 0 \quad \text{--- (C)}$$

OR

$$\begin{cases} n a_0 + \sum_{i=1}^n x_{1i} a_1 + \sum_{i=1}^n x_{2i} a_2 = \sum_{i=1}^n y_i \\ \sum x_{1i} a_0 + \sum x_{1i}^2 a_1 + \sum x_{1i} x_{2i} a_2 = \sum x_{1i} y_i \\ \sum x_{2i} a_0 + \sum x_{1i} x_{2i} a_1 + \sum x_{2i}^2 a_2 = \sum x_{2i} y_i \end{cases}$$

15

OR

$$\begin{bmatrix} n & \sum x_{1i} & \sum x_{2i} \\ \sum x_{1i} & \sum x_{1i}^2 & \sum x_{1i}x_{2i} \\ \sum x_{2i} & \sum x_{1i}x_{2i} & \sum x_{2i}^2 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_{1i}y_i \\ \sum x_{2i}y_i \end{bmatrix}$$

- In general for  $m$ -dimensional case:

Standard error:  $S_{y/x} = \sqrt{\frac{s_r}{n-(m+1)}}$

— \* —

- General Linear Least Squares:

Model:

$$Y = a_0 Z_0 + a_1 Z_1 + a_2 Z_2 + \dots + a_m Z_m$$

Where:  $Z_0, Z_1, \dots, Z_m$  are  $(m+1)$  different functions.

- Simple linear LSR:  $Z_0 = 1, Z_1 = x, Z_i = 0$  for  $i \geq 2$ .
- Polynomial LSR:  $Z_i = x^i$  ( $Z_0 = 1, Z_1 = x, Z_2 = x^2, \dots$ )
- Multiple linear LSR:  $Z_0 = 1, Z_i = x_i$  for  $i \geq 1$ .

Here, "linear" refers only to the model's dependence on its parameters  
— that is  $a$ 's.

K

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n \left( y_i \Big|_{\text{measured}} - y_i \Big|_{\text{model}} \right)^2$$

Given data:  $(z_{0i}, z_{1i}, z_{2i}, \dots, z_{mi}, y_i), i=1, 2, \dots, n$

$$S_r = \sum_{i=1}^n \left( y_i - \sum_{j=0}^m a_j z_{ji} \right)^2$$

Find  $a_j, j=0, 1, 2, \dots, m$  to minimize  $S_r$ .

$$\therefore \frac{\partial S_r}{\partial a_k} = -2 \sum_{i=1}^n \left( y_i - \sum_{j=0}^m a_j z_{ji} \right) \cdot z_{ki} = 0$$

OR

$$\sum_{i=1}^n y_i z_{ki} = \sum_{i=1}^n \sum_{j=0}^m z_{ki} z_{ji} a_j, \quad k=0, 1, 2, \dots, m$$

OR

$$\sum_{j=0}^m \sum_{i=1}^n z_{ki} z_{ji} a_j = \sum_{i=1}^n y_i z_{ki}$$

Normal Equations.

$$\boxed{Z^T Z \{A\} = Z^T \{Y\}},$$

$$Z = \begin{bmatrix} z_{01} & z_{11} & \cdots & z_{m1} \\ z_{02} & z_{12} & \cdots & z_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ z_{0n} & z_{1n} & \cdots & z_{mn} \end{bmatrix}$$

Where

$$\{A\}^T = (a_0, a_1, \dots, a_m) \text{ and } \{Y\}^T = (y_1, y_2, \dots, y_n)$$

$\downarrow$   
unknown coefficients

$\downarrow$   
observed values of the dep. variable

## Numerical Methods

11

Example-1:

Fit  $y = \alpha_2 x^{\beta_2}$  to the following data (using a logarithmic transformation of the data) ①

$x$	$y$	$\log x$	$\log y$
1	0.5	0	-0.301
2	1.7	0.301	0.226
3	3.4	0.477	0.534
4	5.7	0.602	0.753
5	8.4	0.699	0.922

Original data

log transformed data (from Eq. (1))

We can linearize power equation (power law) by taking log (base 10):

$$\log y = \beta_2 \log x + \log \alpha_2 \quad \text{--- (2)}$$

log transformed data is given in the above table.

A linear regression of the log transformed data yields:

$$\log y = 1.75 \log x - 0.300 \quad (\text{Check it!})$$

∴ Intercept  $\log \alpha_2 = -0.300$ , and by taking antilog  $\Rightarrow \boxed{\alpha_2 = 10^{-0.3}} \\ = 0.5$

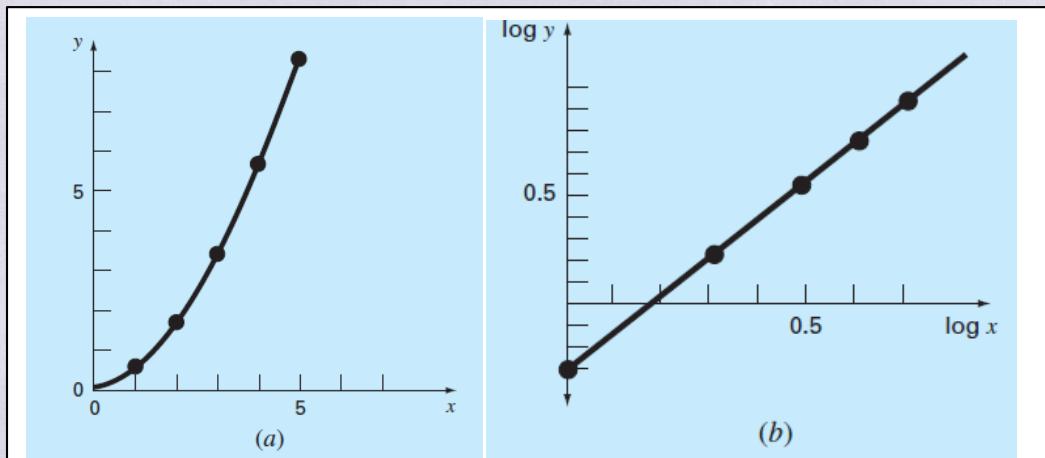
→ slope  $\beta_2 = 1.75$ .

L2

∴ The power equation (fit) is:

$$y = 0.5x^{1.75}$$

□



(a) Plot of untransformed data with the power equation that fits given data. (b) Plot of transformed data used to determine the coefficients of the power equation.

Example-2: (Polynomial regression) Fit a second-order poly. to the following data:

$x_i$	$y_i$	$(y_i - \bar{y})^2 = S_T$	$\frac{(y_i - a_0 - a_1 x_i - a_2 x_i^2)}{= S_r}$
0	2.1		
1	7.7		
2	13.6		
3	27.2		
4	40.9		
5	61.1		

→ Here,  $m=2$ ,  $\sum_{i=1}^6 x_i = 15$ ,  $\sum x_i^4 = 979$

13

$$n=6, \quad \sum y_i = 152.6, \quad \sum x_i y_i = 585.6$$

$$\bar{x} = 2.5, \quad \sum x_i^2 = 55, \quad \sum x_i^2 y_i = 2480.8$$

$$\bar{y} = 25.433, \quad \sum x_i^3 = 225$$

∴ System of simultaneous linear equations:

$$\begin{bmatrix} 6 & 15 & 55 \\ 15 & 55 & 225 \\ 55 & 225 & 979 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 152.6 \\ 585.6 \\ 2480.8 \end{bmatrix}$$

Solving these :  $a_0 = 2.47857, a_1 = 2.35929$   
 (Gauss elimination)  $a_2 = 1.86071$

∴ The least-squares quadratic equation :

$$y = 2.47857 + 2.35929x + 1.86071x^2 \quad (2)$$

$S_{yx}$  (standard error of the estimate based on regression poly (2))

$$= \sqrt{\frac{S_r}{n-(m+1)}} = \sqrt{\frac{3.74657}{6-3}} = 1.12$$

14

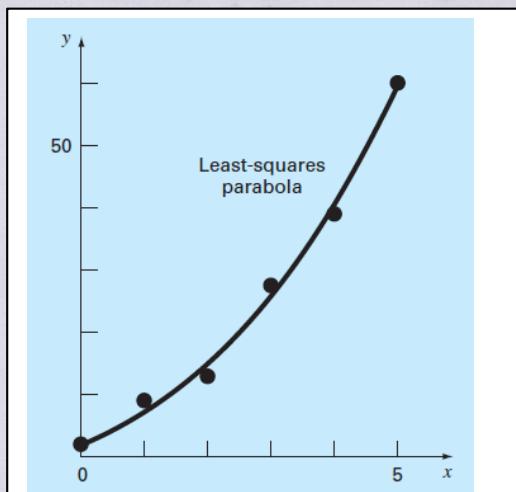
→ Coefficient of determination

$$r^2 = \left( \frac{S_t - S_r}{S_t} \right) = \frac{2513.39 - 3.74657}{2513.39} \\ = 0.99851$$

and correlation coefficient:  $r = \sqrt{0.99851} = 0.99925$

Note: That is 99.851% of the original uncertainty has been explained by the model.

□



Fit of a second-order polynomial.

Example-3: (multilinear regression) Following data was calculated from the equation:

$$y = 5 + 4x_1 - 3x_2$$

$x_1$	$x_2$	$y$
0	0	5
2	1	10
2.5	2	9
1	3	0
4	6	3
7	2	27

Use multiple linear regression to fit this data:

Required computations to get the normal equations are given below:

$y$	$x_1$	$x_2$	$x_1^2$	$x_2^2$	$x_1x_2$	$x_1y$	$x_2y$
5	0	0	0	0	0	0	0
10	2	1	4	1	2	20	10
9	2.5	2	6.25	4	5	22.5	18
0	1	3	1	9	3	0	0
3	4	6	16	36	24	12	18
27	7	2	49	4	14	189	54
$\Sigma$	54	16.5	14	76.25	54	243.5	100

Table: Computations required to develop the normal equations.

System of equations:

$$\begin{bmatrix} 6 & 16.5 & 14 \\ 16.5 & 76.25 & 54 \\ 14 & 54 & 100 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 54 \\ 243.5 \\ 100 \end{bmatrix}$$

16

$$\Rightarrow a_0 = 5, a_1 = 4, a_2 = -3$$

That is consistent with original data.  $\square$

$\square$

### Nonlinear Regression:

How to fit nonlinear models to the data?

Here, models have a nonlinear dependence on their parameters:

For example:  $f(x) = a_0(1 - e^{-a_1 x}) + e$  — (i)

Nonlinear regression is based on determining the values of the parameters that minimize the sum of squares of residuals.

Being nonlinear case, sol' procedure requires an iterative process.

Use Gauss-Newton method for minimizing the sum of the squares of the residuals between data and nonlinear opt.

乙

## — Game - Newton Method

- Taylor series expansion is used to linearize the original nonlinear equation.
  - Then, least-squares technique can be used to obtain new estimate of the parameters that move in the direction of minimizing the residual.

### Steps:

A.) Relationship between the nonlinear equation and the data can be expressed as:

$$y_i = f(x_i; a_0, a_1, \dots, a_m) + e_i$$

measured
model
error

OR

$$y = f(x_i) + e_i \quad \text{--- (ii)}$$

Using Taylor series expansion: (for a two parameter case)

$$f(x_i) \Big|_{j+1} = f(x_i) \Big|_j + \frac{\partial f(x_i) \Big|_j}{\partial a_0} \Delta a_0 + \frac{\partial f(x_i) \Big|_j}{\partial a_1} \Delta a_1 \quad (iii)$$

$$-\Delta a_0 = a_{0,j+1} - a_{0,j} \quad \begin{cases} j = \text{initial guess} \\ j+1 = \text{prediction} \end{cases} \quad \boxed{18}$$

$$\Delta a_1 = a_{1,j+1} - a_{1,j}$$

From (ii) and (iii) :

$$y_i - f(x_i)|_j = \frac{\partial f(x_i)|_j}{\partial a_0} \Delta a_0 + \frac{\partial f(x_i)|_j}{\partial a_1} \Delta a_1 + e_i$$

OR, in matrix notations,

$$\{D\} = [Z_j] \{\Delta A\} + \{E\} \quad \longrightarrow \textcircled{iv}$$

Where

$$[Z_j] = \begin{bmatrix} \frac{\partial f_1}{\partial a_0} & \frac{\partial f_1}{\partial a_1} \\ \frac{\partial f_2}{\partial a_0} & \frac{\partial f_2}{\partial a_1} \\ \vdots & \vdots \\ \frac{\partial f_n}{\partial a_0} & \frac{\partial f_n}{\partial a_1} \end{bmatrix}$$

$n$  = no. of data points

$\frac{\partial f_i}{\partial a_k}$  = partial derivative  
of the function wrt  
the  $k^{\text{th}}$  parameter

evaluated at the  
 $i^{\text{th}}$  data point.

$$\{D\} = \left\{ \begin{array}{l} y_1 - f(x_1) \\ y_2 - f(x_2) \\ \vdots \\ y_n - f(x_n) \end{array} \right\}$$

↓  
Diff. between  
measurement  
and function  
values.

[9]

$$\{\Delta A\} = \begin{Bmatrix} \Delta a_0 \\ \Delta a_1 \\ \vdots \\ \Delta a_m \end{Bmatrix}$$

↓  
changes in  
the parameter  
values.

C.) Apply linear least-squares to Eq.(iv) in the following normal equations:

$$\left( [z_j]^T [z_j] \right) \{\Delta A\} = [z_j]^T \{D\} \quad \text{--- (v)}$$

That is solve Eq.(v) for  $\{\Delta A\}$ , which can be used to get the improved values for the parameters, as:

$$a_0|_{j+1} = a_0|_j + \Delta a_0$$

$$\text{and } a_i|_{j+1} = a_i|_j + \Delta a_i$$

This procedure will be repeated until the sol.

Converges, that is, until:

$$|\epsilon_k|_k = \left| \frac{a_k|_{j+1} - a_k|_j}{a_k|_{j+1}} \right| \leq 100 \% \quad \begin{array}{l} \text{falls below} \\ \text{acceptable} \\ \text{stopping criterion} \end{array} \quad \square$$

## Numerical Methods

11

Example-1 (Nonlinear regression using Gauss-Newton Method )

Fit the function  $f(x; a_0, a_1) = a_0(1 - e^{-a_1 x})$  to the data:

X	0.25	0.75	1.25	1.75	2.25
Y	0.28	0.57	0.68	0.74	0.79

Use initial guesses of  $a_0 = 1.0$  and  $a_1 = 1.0$  for the parameters.

Sol: :  $f = a_0(1 - e^{-a_1 x})$

$$\begin{aligned} \frac{\partial f}{\partial a_0} &= 1 - e^{-a_1 x} \\ \frac{\partial f}{\partial a_1} &= a_0 x e^{-a_1 x} \end{aligned}$$

$$\therefore [Z_0] = \begin{bmatrix} \frac{\partial f_1}{\partial a_0} & \frac{\partial f_1}{\partial a_1} \\ \frac{\partial f_2}{\partial a_0} & \frac{\partial f_2}{\partial a_1} \\ \frac{\partial f_3}{\partial a_0} & \frac{\partial f_3}{\partial a_1} \\ \frac{\partial f_4}{\partial a_0} & \frac{\partial f_4}{\partial a_1} \\ \frac{\partial f_5}{\partial a_0} & \frac{\partial f_5}{\partial a_1} \end{bmatrix} = \begin{bmatrix} 0.2212 & 0.1947 \\ 0.5276 & 0.3543 \\ 0.7135 & 0.3581 \\ 0.8262 & 0.3041 \\ 0.8946 & 0.2371 \end{bmatrix}$$

12

$$[Z_0]^T [Z_0] = \begin{bmatrix} 2.3193 & 0.9489 \\ 0.9489 & 0.4404 \end{bmatrix}$$

Vector  $D$  - differences between the measurements and the model predictions.

$$\{D\} = \left\{ \begin{array}{c} \\ y_i - f(x_i) \\ \end{array} \right\}_{1 \leq i \leq 5} = \begin{bmatrix} 0.0588 \\ 0.0424 \\ -0.0335 \\ -0.0862 \\ -0.1046 \end{bmatrix}$$

$$[Z_0]^T \{D\} = \begin{bmatrix} -0.1533 \\ -0.0365 \end{bmatrix}$$

Vector  $\Delta A$  - changes in the parameter values

$$\{\Delta A\} = \begin{Bmatrix} \Delta a_0 \\ \Delta a_1 \end{Bmatrix} = \begin{Bmatrix} -0.2714 \\ 0.5019 \end{Bmatrix}$$

∴ updated values of free parameters:

$$\begin{Bmatrix} a_0 \\ a_1 \end{Bmatrix} = \begin{Bmatrix} 1.0 \\ 1.0 \end{Bmatrix} + \begin{Bmatrix} -0.2714 \\ 0.5019 \end{Bmatrix} = \begin{Bmatrix} 0.7286 \\ 1.5019 \end{Bmatrix}$$

— (\*)

13

- Error/residual:

$$e_i = \underbrace{f(x_i)}_{\text{measured}} - \underbrace{a_0(1-e^{-a_1 x})}_{\text{model}}$$

- Sum of squares of the residuals:

$$S_r = \sum_{i=1}^5 e_i^2$$

$$S_r \Big|_{IG} = 0.0248$$

and  $S_r \Big|_{\oplus} = 0.0242$

- Repeat the computations until stopping criterion is satisfied.

□

14

## Ordinary Differential Equations:

Q. Given  $y' = \frac{dy}{dx} = f(x, y)$ , find  $y(x)$ .

### Euler's Method:

- An iterative method.

Given  $(x_i, y_i)$ ,  $x_{i+1} = x_i + h$ ,  $y_{i+1} = y_i + \phi h$   
where,  $\phi$  is the estimated slope.

- In Euler's Method, first derivative is used to estimate the function slope, as:

$$\phi = f(x_i, y_i) \text{ and } y_{i+1} = y_i + f(x_i, y_i) \cdot h$$

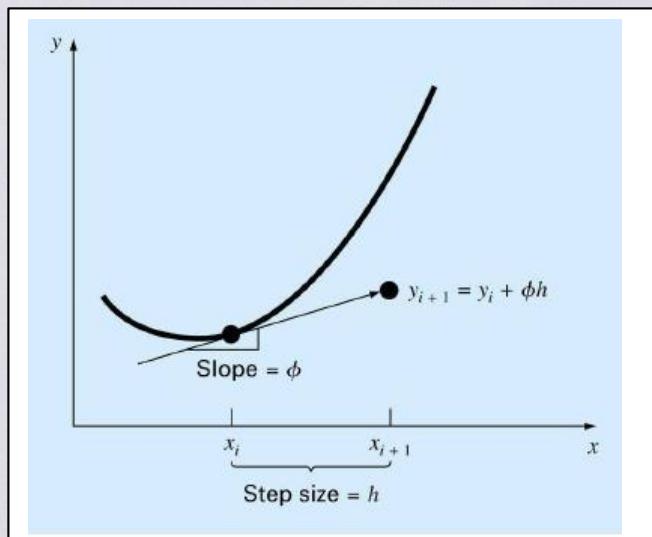


Figure 1: Illustration of iterative methods

15

- Use Taylor series expansion:

$$y_{i+1} = y_i + h y'_i + \frac{h^2}{2!} y''_i + \cdots + \frac{h^n}{n!} y_i^{(n)} + R_n$$

with  $h = x_{i+1} - x_i$  and  $R_n = \frac{y_{(n+1)}^{(\alpha)}}{(n+1)!} h^{n+1} = O(h^{n+1})$

$$x_i < \alpha < x_{i+1}$$

$y' = f(x, y)$ , we have  $y'_i = f(x_i, y_i)$

$$y'' = f'(x_i, y_i)$$

$$\vdots$$

$$y_i^{(n)} = f^{(n-1)}(x_i, y_i)$$

$$\begin{aligned} y_{i+1} &= y_i + f(x_i, y_i)h + \frac{1}{2!} f'(x_i, y_i)h^2 + \cdots \\ &\quad + \frac{1}{n!} f^{(n-1)}(x_i, y_i)h^n + O(h^{n+1}) \end{aligned}$$

Using Euler's method:

$$\boxed{y_{i+1} = y_i + f(x_i, y_i)h}$$

16

$\therefore$  The true local truncation error in using Euler's method:

$$E_t = \frac{1}{2} f'(x_i, y_i) h^2 + \dots + \frac{1}{n!} f^{(n-1)}(x_i, y_i) h^n + O(h^{n+1})$$

Here,  $h$  is sufficiently small, higher powers of  $h$  can be neglected. Thus, approximated local truncation error is:

$$E_a = \frac{1}{2} f'(x_i, y_i) h^2$$

- $E_a$  is proportional to  $h^2$  and  $f'(x_i, y_i)$
- Taylor series provides only the local truncation error.
- Global truncation error using Euler's method is proportional to the step-size,  $O(h)$ .
- Truncation error can be reduced by decreasing the step-size ( $h$ )
- Euler's method provides error free solution/prediction if the function  $f(x)$  (that is sol. of diff. eq<sup>n</sup>) is linear.

□.

Example:

Integrate the eqn:  $\frac{dy}{dx} = -2x^3 + 12x^2 - 20x + 8.5$

17

from  $x=0$  to  $x=1$

Using Euler's method with a step-size  
 $h=0.5$  and  $0.25$ .

The initial condition: at  $x=0, y=1$

Euler's method with  $h=0.5$ :

$$x_1 = x_0 + h = 0.5$$

$$y_1 = y_0 + f(x_0, y_0)h = 1 + f(0, 1)*0.5 = 5.25$$

$$x_2 = x_1 + h = 1.0$$

$$y_2 = y_1 + f(x_1, y_1)h = 5.25 + f(0.5, 5.25)*0.5 = 5.875 \quad \left. \right\} \text{---A}$$

Euler's method with  $h=0.25$ :

$$x_1 = x_0 + h = 0.25$$

$$y_1 = y_0 + f(x_0, y_0)h = 1 + f(0, 1)*0.25 = 3.1250 \quad \left. \right\} \text{---B}$$

$$x_2 = x_1 + h = 0.5$$

$$y_2 = y_1 + f(x_1, y_1)h = 3.1250 + f(0.25, 3.1250)*0.25 = 4.1797 \quad \left. \right\} \text{---B}$$

$$x_3 = x_2 + h = 0.75 \quad | \quad x_4 = x_3 + h = 1.0$$

$$y_3 = 4.4922 \quad | \quad y_4 = 4.3438$$

(B)

→ Exact soln.

$$y = \int (-2x^3 + 12x^2 - 20x + 8.5) dx$$

$$= -\frac{1}{2}x^4 + 4x^3 - 10x^2 + 8.5x + C$$

Using EC:  $x=0, y=1$

$$\Rightarrow y = -\frac{1}{2}x^4 + 4x^3 - 10x^2 + 8.5x + 1 \quad (C)$$

$$\therefore \boxed{y(1) = 3}$$

Percent relative error (with  $h=0.5$ ):

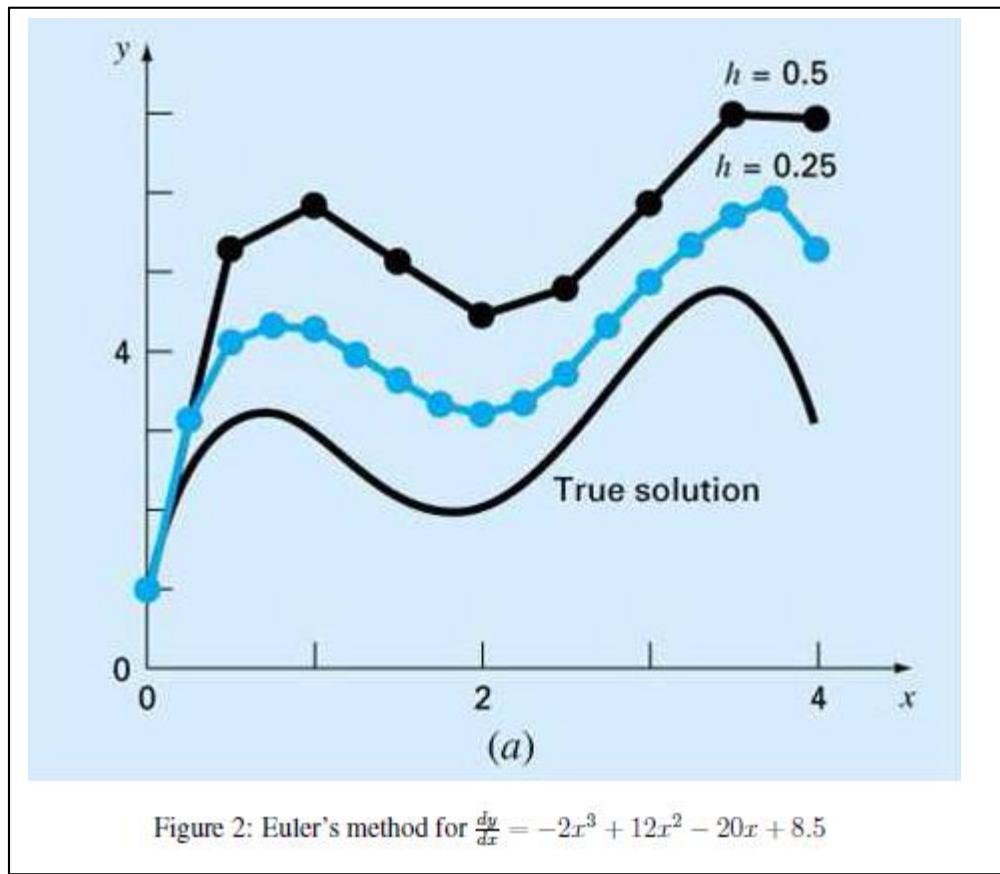
$$\epsilon_t = \left| \frac{\text{true value} - \text{approximate value}}{\text{true value}} \right| * 100 \%$$

$$\cancel{26.3\%} = 95.8\%$$

Percent relative error (with  $h=0.25$ )

$$\epsilon_t = 44.79\%$$

⇒ Reducing step-size can reduce the approximation error.



## Numerical Methods

11

### Higher-order Taylor Series Method:

Using second-order Taylor Series:

$$y_{i+1} = y_i + y'_i h + \frac{y''_i}{2} h^2 = y_i + f(x_i, y_i)h + \frac{1}{2} f'(x_i, y_i)h^2$$

Where

$$f'(x, y) = \frac{\partial f(x, y)}{\partial x} + \frac{\partial f(x, y)}{\partial y} \frac{dy}{dx} \quad (\text{chain-rule})$$

Local Truncation error:

$$E_a = \frac{1}{3!} f''(x_i, y_i) h^3$$

$$\text{Where } f''(x, y) = \frac{\partial f'}{\partial x} + \frac{\partial f'}{\partial y} \frac{dy}{dx}$$

Note:  $f'(x, y)$  and  $f''(x, y)$  may be difficult to evaluate for complicated functions.

□

12

## Runge - Kutta methods :

Runge - Kutta (RK) methods can achieve the accuracy of higher Taylor series BUT avoid evaluating the higher order derivatives.

The general form of RK method is

$$y_{i+1} = y_i + \phi(x_i, y_i, h)h$$

Where,  $\phi(x_i, y_i, h)$  is called an increment function and is written in general form as:

$$\phi = a_1 k_1 + a_2 k_2 + \dots + a_n k_n$$

Where,

$$k_1 = f(x_i, y_i)$$

$$k_2 = f(x_i + p_1 h, y_i + q_{11} h k_1)$$

$$k_3 = f(x_i + p_2 h, y_i + q_{21} k_1 h + q_{22} k_2 h)$$

⋮

$$k_n = f(x_i + p_{n-1} h, y_i + q_{n-1,1} k_1 h + q_{n-2,2} k_2 h + \dots + q_{n,n-1} k_{n-1} h)$$

L3

- Various types of RK methods can be devised by choosing different number of terms in  $\phi$  and different values of the parameters  $a$ 's,  $p$ 's and  $q$ 's.

### First-Order RK methods :

$n=1$ , letting  $a_1 = 1$ , we have,  $\phi(x_i, y_i, h) = a_1 k_1 = k_1$ ,

$$\Rightarrow y_{i+1} = y_i + f(x_i, y_i)h \quad \underbrace{\text{--- Euler's Method}}_{\downarrow} \quad \text{first-order RK method.}$$

### Second-order RK Methods :

$$y_{i+1} = y_i + (a_1 k_1 + a_2 k_2)h$$

where

$$k_1 = f(x_i, y_i)$$

$$k_2 = f(x_i + p_1 h, y_i + q_1 k_1, h)$$

Q: How to find constants  $a_1$ ,  $a_2$ ,  $p_1$  and  $q_1$ ?

4

Use Taylor series:

$$\begin{aligned}
 y_{i+1} &= y_i + y'_i h + \frac{1}{2} y''_i h^2 + \underbrace{H_o T}_{\downarrow \text{(ignore them)}} \\
 &= y_i + f(x_i, y_i)h + \frac{1}{2} f'(x_i, y_i)h^2 \\
 &= y_i + f(x_i, y_i)h + \frac{1}{2} \left[ \frac{\partial f(x_i, y_i)}{\partial x} + \frac{\partial f(x_i, y_i)}{\partial y} y'_i \right] h^2
 \end{aligned}
 \quad \text{--- (1)}$$

Using 2nd-order RK method:

$$\begin{aligned}
 y_{i+1} &= y_i + a_1 k_1 h + a_2 k_2 h \\
 &= y_i + a_1 f(x_i, y_i)h + a_2 k_2 h
 \end{aligned}
 \quad \text{--- (2)}$$

Where,

$$\begin{aligned}
 k_2 &= f(x_i + p_1 h, y_i + q_1 k_1 h) \\
 &= f(x_i, y_i) + \frac{\partial f(x_i, y_i)}{\partial x} p_1 h + \frac{\partial f(x_i, y_i)}{\partial y} q_1 k_1 h + H_o T
 \end{aligned}$$

(3)

$$= f(x_i, y_i) + \frac{\partial f(x_i, y_i)}{\partial x} p_1 h + \frac{\partial f(x_i, y_i)}{\partial y} q_1 f(x_i, y_i) h$$

Substitute  $k_2$  in ② by ③:

15

$$\begin{aligned}y_{i+1} = y_i + (a_1 + a_2)f(x_i, y_i)h + a_2 \frac{\partial f(x_i, y_i)}{\partial x} p_i h^2 \\+ a_2 \frac{\partial f(x_i, y_i)}{\partial y} q_{11} f(x_i, y_i) h^2\end{aligned}\quad \text{--- } ④$$

Comparing ④ and ①:

$$\left. \begin{array}{l} a_1 + a_2 = 1 \\ a_2 p_i = \frac{1}{2} \\ a_2 q_{11} = \frac{1}{2} \end{array} \right\} \begin{array}{l} \text{Three equations in} \\ \text{Four unknowns} \\ \Rightarrow \text{infinite no. of} \\ \text{sols are possible.} \end{array}$$

A) Heun Method:

$$a_1 = \frac{1}{2}, a_2 = \frac{1}{2}, \text{ and } p_i = q_{11} = 1$$

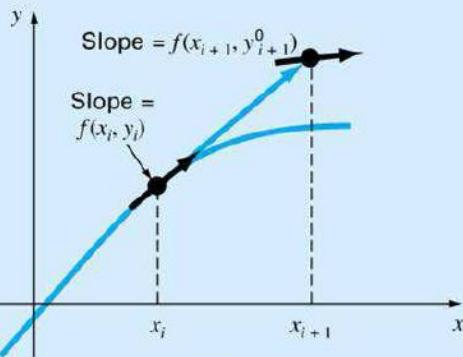
$$\boxed{\begin{aligned}y_{i+1} &= y_i + \frac{1}{2}(k_1 + k_2) \\k_1 &= f(x_i, y_i) \\k_2 &= f(x_i + h, y_i + k_1 h)\end{aligned}}$$

→ Heun method can be represented in a predictor-corrector approach:

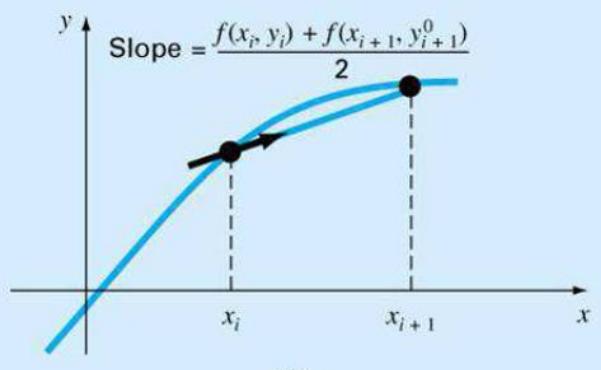
16

- Predictor:  $y_{i+1}^0 = y_i + f(x_i, y_i)h$

- Corrector:  $y_{i+1} = y_i + \left( \frac{f(x_i, y_i) + f(x_{i+1}, y_{i+1}^0)}{2} \right) h$



(a)



(b)

Heun method: (a) Predictor; (b) Corrector

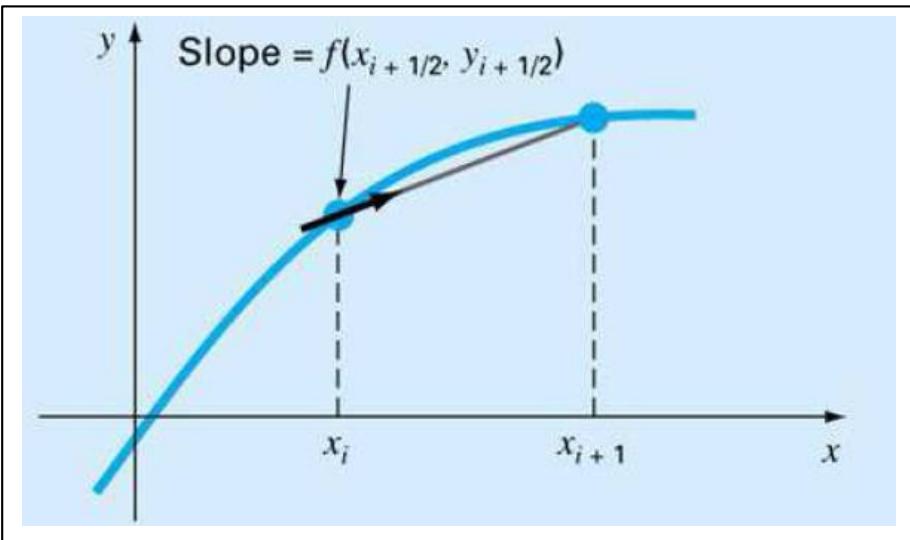
### B) The Midpoint method

$$a_1 = 0, a_2 = 1, \text{ and } p_1 = q_1 = \frac{1}{2}$$

$$y_{i+1} = y_i + k_2 h$$

$$k_1 = f(x_i, y_i)$$

$$k_2 = f\left(x_i + \frac{h}{2}, y_i + \frac{k_1 h}{2}\right)$$



Midpoint Method

## Numerical Methods

11

### - Fourth-order Runge-Kutta (RK) Methods:

— Fourth-order RK methods have the following general form:

$$y_{i+1} = y_i + (a_1 k_1 + a_2 k_2 + a_3 k_3 + a_4 k_4) h$$

Similar to second-order RK methods, there are infinite number of versions of 4<sup>th</sup>-order RK methods. The most commonly used form is:

$$y_{i+1} = y_i + \frac{1}{6} (k_1 + 2k_2 + 2k_3 + k_4) h$$

Where,

$$k_1 = f(x_i, y_i)$$

$$k_2 = f\left(x_i + \frac{1}{2}h, y_i + \frac{1}{2}k_1 h\right)$$

$$k_3 = f\left(x_i + \frac{1}{2}h, y_i + \frac{1}{2}k_2 h\right)$$

$$k_4 = f(x_i + h, y_i + k_3 h)$$

↳ RK<sub>4</sub> method.

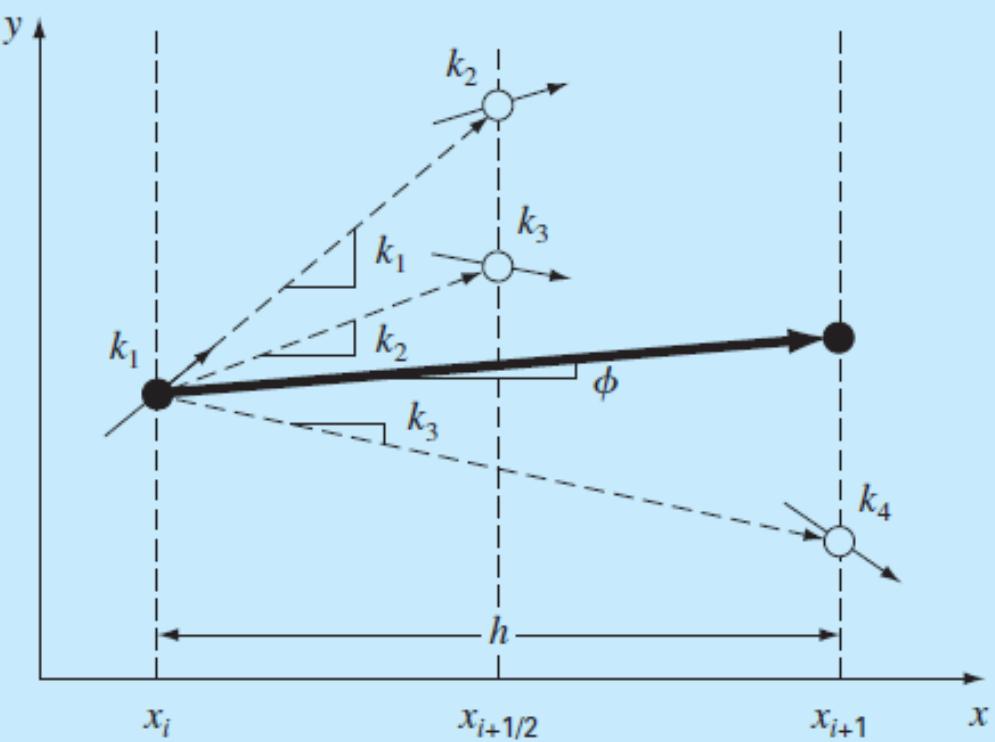


Fig. Slope estimates in the fourth-order RK method.

12

Example:  $f(x, y) = -2x^3 + 12x^2 - 20x + 0.5$

Solve  $\begin{cases} y' = f(x, y) \\ \text{with } y(0) = 1 \end{cases}$

Using RK<sub>4</sub> method with  $h = 0.5$ .

$i = 0, x_0 = 0, y_0 = 1$

$k_1 = f(x_0, y_0) = f(0, 1) = 0.5$

$k_2 = f\left(x_0 + \frac{h}{2}, y_0 + \frac{1}{2}k_1 h\right) = f(0.25, 3.125) = 4.21875$

$k_3 = f\left(x_0 + \frac{1}{2}h, y_0 + \frac{1}{2}k_2 h\right) = f(0.25, 2.0547) = 4.21875$

$k_4 = f(x_0 + h, y_0 + k_3 h) = f(0.25, 3.10937) = 1.25$

$\Rightarrow x_1 = x_0 + h = 0.5$

$y_1 = y_0 + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4) = 3.21875$

□

Also,  $y|_{x=0.5} = 3.21875$  — Exact value.

This is because  $y(x)$  is a 4<sup>th</sup>-order polynomial  
 $\Rightarrow$  4<sup>th</sup>-order RK methods give exact solution.

L3

→ In general:  $n^{\text{th}}$ -order RK methods

- A.) Accurate to  $n^{\text{th}}$ -order polynomial.
- B.) Equivalent to  $n^{\text{th}}$ -order Taylor Series approx.
- C.) Does not require to evaluate derivatives.

— System of ODEs:

For a system of simultaneous ODEs like:

$$\left. \begin{array}{l} \frac{dy_1}{dx} = f_1(x, y_1, y_2, \dots, y_n) \\ \frac{dy_2}{dx} = f_2(x, y_1, y_2, \dots, y_n) \\ \vdots \\ \frac{dy_n}{dx} = f_n(x, y_1, y_2, \dots, y_n) \end{array} \right\}$$

sol. requires  
 $n$  initial conditions.

— RK methods discussed for a single ODE can be extended to system of ODEs.

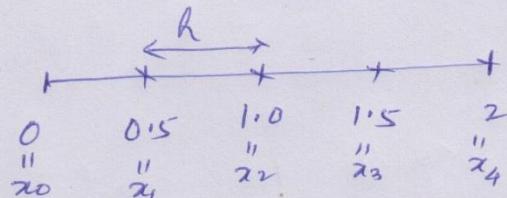
Example 2: Use Euler's method to solve

14

$$\left. \begin{array}{l} \frac{dy_1}{dx} = -0.5 y_1 \\ \frac{dy_2}{dx} = 4 - 0.3 y_2 - 0.1 y_1 \end{array} \right\}$$

with  $y_1(0) = 4$  and  $y_2(0) = 6$ .

Find solution at  $x=2$  with a step-size ( $h$ ) = 0.5.



$$y_{i+1,1} = y_{i,1} + f_1(x_i, y_{i,1}, y_{i,2})$$

$$y_{i+1,2} = y_{i,2} + f_2(x_i, y_{i,1}, y_{i,2})$$

Here,  $f_1(x, y_1, y_2) = -0.5 y_1$

$$f_2(x, y_1, y_2) = 4 - 0.3 y_2 - 0.1 y_1$$

-  $i=0, x_1 = x_0 + h = 0.5$

$$\left. \begin{array}{l} y_{1,1} = y_{0,1} + f_1(x_0, y_{0,1}, y_{0,2})h = 4 + f_1(0, 4, 6) = 3 \\ y_{1,2} = y_{0,2} + f_2(x_0, y_{0,1}, y_{0,2})h = 6 + f_2(0, 4, 6) = 6.9 \end{array} \right\} \quad \square$$

15

$$i=1, \quad x_2 = x_1 + h = 1.0$$

$$\begin{cases} y_{1,1} = y_{1,1} + f_1(x_1, y_{1,1}, y_{1,2})h = 3 + f_1(0.5, 3, 6.9) = 2.25 \\ y_{1,2} = y_{1,2} + f_2(x_1, y_{1,1}, y_{1,2})h = 6.9 + f_2(0.5, 3, 6.9) = 7.715 \end{cases}$$

↓  
↓  
↓

$$y_{4,1} = \dots$$

$$y_{4,2} = \dots$$

□

Q. Solve the prob. in Example 2 using RK<sub>4</sub> method.

$$\begin{cases} y_{i+1,1} = y_{i,1} + \frac{1}{6}(k_{1,1} + 2k_{2,1} + 2k_{3,1} + k_{4,1})h \\ y_{i+1,2} = y_{i,2} + \frac{1}{6}(k_{1,2} + 2k_{2,2} + 2k_{3,2} + k_{4,2})h \end{cases}$$

where,

$$k_{1,1} = f_1(x_i, y_{i,1}, y_{i,2})$$

$$k_{2,1} = f_1\left(x_i + \frac{1}{2}h, y_{i,1} + \frac{1}{2}k_{1,1}h, y_{i,2} + \frac{1}{2}k_{1,2}h\right)$$

$$k_{3,1} = f_1\left(x_i + \frac{1}{2}h, y_{i,1} + \frac{1}{2}k_{2,1}h, y_{i,2} + \frac{1}{2}k_{2,2}h\right)$$

$$k_{4,1} = f_1(x_i + h, y_{i,1} + k_{3,1}h, y_{i,2} + k_{3,2}h)$$

$$k_{1,2} = f_2(x_i, y_{i,1}, y_{i,2})$$

$$k_{2,2} = f_2\left(x_i + \frac{1}{2}h, y_{i,1} + \frac{1}{2}k_{1,1}h, y_{i,2} + \frac{1}{2}k_{1,2}h\right)$$

$$k_{3,2} = f_2\left(x_i + \frac{1}{2}h, y_{i,1} + \frac{1}{2}k_{2,1}h, y_{i,2} + \frac{1}{2}k_{2,2}h\right)$$

$$k_{4,2} = f_2(x_i + h, y_{i,1} + k_{3,1}h, y_{i,2} + k_{3,2}h)$$

**Solution.** First, we must solve for all the slopes at the beginning of the interval:

$$k_{1,1} = f_1(0, 4, 6) = -0.5(4) = -2$$

$$k_{1,2} = f_2(0, 4, 6) = 4 - 0.3(6) - 0.1(4) = 1.8$$

where  $k_{i,j}$  is the  $i$ th value of  $k$  for the  $j$ th dependent variable. Next, we must calculate the first values of  $y_1$  and  $y_2$  at the midpoint:

$$y_1 + k_{1,1} \frac{h}{2} = 4 + (-2) \frac{0.5}{2} = 3.5$$

$$y_2 + k_{1,2} \frac{h}{2} = 6 + (1.8) \frac{0.5}{2} = 6.45$$

which can be used to compute the first set of midpoint slopes,

$$k_{2,1} = f_1(0.25, 3.5, 6.45) = -1.75$$

$$k_{2,2} = f_2(0.25, 3.5, 6.45) = 1.715$$

These are used to determine the second set of midpoint predictions,

$$y_1 + k_{2,1} \frac{h}{2} = 4 + (-1.75) \frac{0.5}{2} = 3.5625$$

$$y_2 + k_{2,2} \frac{h}{2} = 6 + (1.715) \frac{0.5}{2} = 6.42875$$

which can be used to compute the second set of midpoint slopes,

$$k_{3,1} = f_1(0.25, 3.5625, 6.42875) = -1.78125$$

$$k_{3,2} = f_2(0.25, 3.5625, 6.42875) = 1.715125$$

These are used to determine the predictions at the end of the interval

$$y_1 + k_{3,1} h = 4 + (-1.78125)(0.5) = 3.109375$$

$$y_2 + k_{3,2} h = 6 + (1.715125)(0.5) = 6.857563$$

which can be used to compute the endpoint slopes,

$$k_{4,1} = f_1(0.5, 3.109375, 6.857563) = -1.554688$$

$$k_{4,2} = f_2(0.5, 3.109375, 6.857563) = 1.631794$$

The values of  $k$  can then be used to compute [Eq. (25.40)]:

$$y_1(0.5) = 4 + \frac{1}{6}[-2 + 2(-1.75 - 1.78125) - 1.554688]0.5 = 3.115234$$

$$y_2(0.5) = 6 + \frac{1}{6}[1.8 + 2(1.715 + 1.715125) + 1.631794]0.5 = 6.857670$$

Proceeding in a like manner for the remaining steps yields

<b>x</b>	<b>y<sub>1</sub></b>	<b>y<sub>2</sub></b>
0	4	6
0.5	3.115234	6.857670
1.0	2.426171	7.632106
1.5	1.889523	8.326886
2.0	1.471577	8.946865

\*\*\*

## Numerical Methods

11

### - Fourth-order Runge-Kutta (RK) methods:

— Fourth-order RK methods have the following general form:

$$y_{i+1} = y_i + (a_1 k_1 + a_2 k_2 + a_3 k_3 + a_4 k_4) h$$

Similar to second-order RK methods, there are infinite number of versions of 4<sup>th</sup>-order RK methods. The most commonly used form is:

$$y_{i+1} = y_i + \frac{1}{6} (k_1 + 2k_2 + 2k_3 + k_4) h$$

Where,

$$k_1 = f(x_i, y_i)$$

$$k_2 = f\left(x_i + \frac{1}{2}h, y_i + \frac{1}{2}k_1 h\right)$$

$$k_3 = f\left(x_i + \frac{1}{2}h, y_i + \frac{1}{2}k_2 h\right)$$

$$k_4 = f(x_i + h, y_i + k_3 h)$$

↳ RK<sub>4</sub> method.

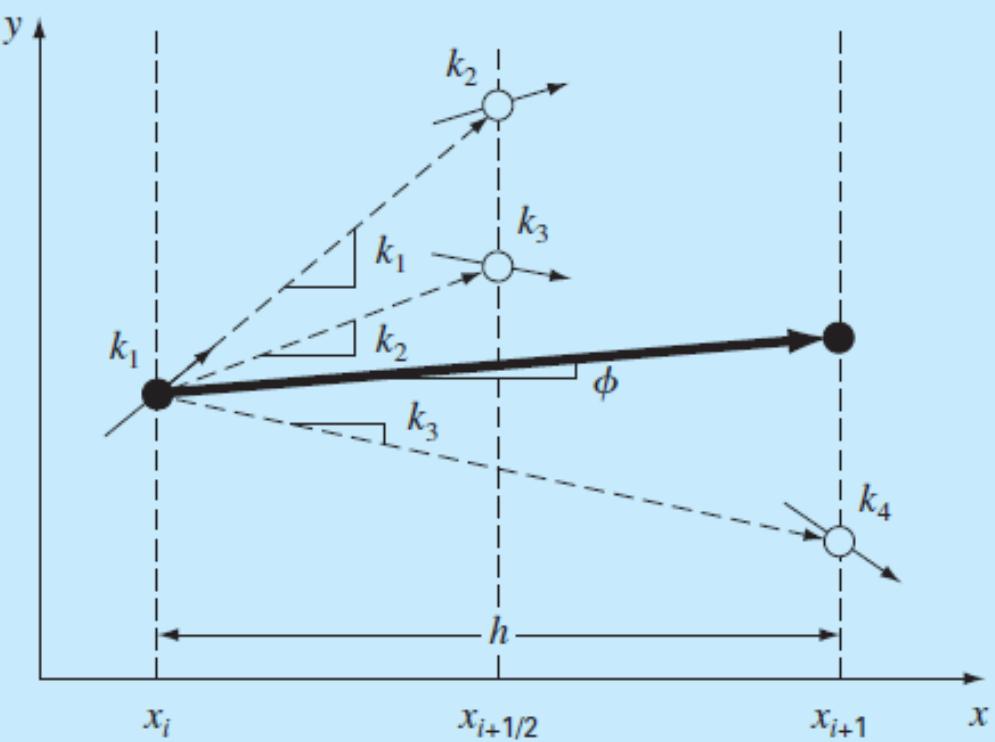


Fig. Slope estimates in the fourth-order RK method.

12

Example:  $f(x, y) = -2x^3 + 12x^2 - 20x + 8.5$

Solve  $\begin{cases} y' = f(x, y) \\ \text{with } y(0) = 1 \end{cases}$

Using RK<sub>4</sub> method with  $h = 0.5$ .

$i = 0, x_0 = 0, y_0 = 1$

$k_1 = f(x_0, y_0) = f(0, 1) = 8.5$

$k_2 = f\left(x_0 + \frac{h}{2}, y_0 + \frac{1}{2}k_1 h\right) = f(0.25, 3.125) = 4.21875$

$k_3 = f\left(x_0 + \frac{1}{2}h, y_0 + \frac{1}{2}k_2 h\right) = f(0.25, 2.0547) = 4.21875$

$k_4 = f(x_0 + h, y_0 + k_3 h) = f(0.25, 3.10937) = 1.25$

$\Rightarrow x_1 = x_0 + h = 0.5$

$y_1 = y_0 + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4) = 3.21875$

□

Also,  $y|_{x=0.5} = 3.21875$  — Exact value.

This is because  $y(x)$  is a 4<sup>th</sup>-order polynomial  
 $\Rightarrow$  4<sup>th</sup>-order RK methods give exact solution.

→ In general:       $n^{\text{th}}$ -order RK methods

13

- A.) Accurate to  $n^{\text{th}}$ -order polynomial.
- B.) Equivalent to  $n^{\text{th}}$ -order Taylor series approx.
- C.) Does not require to evaluate derivatives.

— System of ODEs:

For a system of simultaneous ODEs like:

$$\left. \begin{array}{l} \frac{dy_1}{dx} = f_1(x, y_1, y_2, \dots, y_n) \\ \frac{dy_2}{dx} = f_2(x, y_1, y_2, \dots, y_n) \\ \vdots \\ \frac{dy_n}{dx} = f_n(x, y_1, y_2, \dots, y_n) \end{array} \right\}$$

sol. requires  
 $n$  initial conditions.

— RK methods discussed for a single ODE can be extended to system of ODEs.

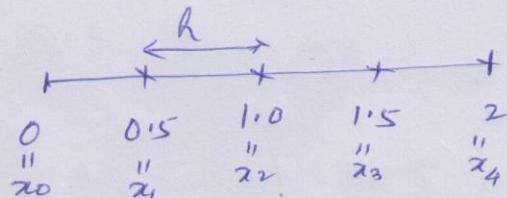
Example 2: Use Euler's method to solve

14

$$\left. \begin{aligned} \frac{dy_1}{dx} &= -0.5 y_1 \\ \frac{dy_2}{dx} &= 4 - 0.3 y_2 - 0.1 y_1 \end{aligned} \right\}$$

with  $y_1(0) = 4$  and  $y_2(0) = 6$ .

Find solution at  $x=2$  with a step-size ( $h$ ) = 0.5.



$$y_{i+1,1} = y_{i,1} + f_1(x_i, y_{i,1}, y_{i,2})$$

$$y_{i+1,2} = y_{i,2} + f_2(x_i, y_{i,1}, y_{i,2})$$

Here,  $f_1(x, y_1, y_2) = -0.5 y_1$

$$f_2(x, y_1, y_2) = 4 - 0.3 y_2 - 0.1 y_1 \quad \left. \right\}$$

-  $i=0, x_1 = x_0 + h = 0.5$

$$\left. \begin{aligned} y_{0,1} &= y_{0,1} + f_1(x_0, y_{0,1}, y_{0,2})h = 4 + f_1(0, 4, 6) = 3 \\ y_{0,2} &= y_{0,2} + f_2(x_0, y_{0,1}, y_{0,2})h = 6 + f_2(0, 4, 6) = 6.9 \end{aligned} \right. \quad \square$$

15

$$i=1, \quad x_2 = x_1 + h = 1.0$$

$$\begin{cases} y_{1,1} = y_{1,1} + f_1(x_1, y_{1,1}, y_{1,2})h = 3 + f_1(0.5, 3, 6.9) = 2.25 \\ y_{1,2} = y_{1,2} + f_2(x_1, y_{1,1}, y_{1,2})h = 6.9 + f_2(0.5, 3, 6.9) = 7.715 \end{cases}$$

↓  
↓  
↓

$$y_{4,1} = \dots$$

$$y_{4,2} = \dots$$

□

Q. Solve the prob. in Example 2 using RK<sub>4</sub> method.

$$\begin{cases} y_{i+1,1} = y_{i,1} + \frac{1}{6}(k_{1,1} + 2k_{2,1} + 2k_{3,1} + k_{4,1})h \\ y_{i+1,2} = y_{i,2} + \frac{1}{6}(k_{1,2} + 2k_{2,2} + 2k_{3,2} + k_{4,2})h \end{cases}$$

where,

$$\begin{cases} k_{1,1} = f_1(x_i, y_{i,1}, y_{i,2}) \\ k_{2,1} = f_1\left(x_i + \frac{h}{2}, y_{i,1} + \frac{k_{1,1}h}{2}, y_{i,2} + \frac{k_{1,2}h}{2}\right) \\ k_{3,1} = f_1\left(x_i + \frac{h}{2}, y_{i,1} + \frac{k_{2,1}h}{2}, y_{i,2} + \frac{k_{2,2}h}{2}\right) \\ k_{4,1} = f_1\left(x_i + h, y_{i,1} + \frac{k_{3,1}h}{2}, y_{i,2} + \frac{k_{3,2}h}{2}\right) \end{cases}$$

$$\begin{cases} k_{1,2} = f_2(x_i, y_{i,1}, y_{i,2}) \\ k_{2,2} = f_2\left(x_i + \frac{h}{2}, y_{i,1} + \frac{k_{1,1}h}{2}, y_{i,2} + \frac{k_{1,2}h}{2}\right) \\ k_{3,2} = f_2\left(x_i + \frac{h}{2}, y_{i,1} + \frac{k_{2,1}h}{2}, y_{i,2} + \frac{k_{2,2}h}{2}\right) \\ k_{4,2} = f_2\left(x_i + h, y_{i,1} + \frac{k_{3,1}h}{2}, y_{i,2} + \frac{k_{3,2}h}{2}\right) \end{cases}$$

**Solution.** First, we must solve for all the slopes at the beginning of the interval:

$$k_{1,1} = f_1(0, 4, 6) = -0.5(4) = -2$$

$$k_{1,2} = f_2(0, 4, 6) = 4 - 0.3(6) - 0.1(4) = 1.8$$

where  $k_{i,j}$  is the  $i$ th value of  $k$  for the  $j$ th dependent variable. Next, we must calculate the first values of  $y_1$  and  $y_2$  at the midpoint:

$$y_1 + k_{1,1} \frac{h}{2} = 4 + (-2) \frac{0.5}{2} = 3.5$$

$$y_2 + k_{1,2} \frac{h}{2} = 6 + (1.8) \frac{0.5}{2} = 6.45$$

which can be used to compute the first set of midpoint slopes,

$$k_{2,1} = f_1(0.25, 3.5, 6.45) = -1.75$$

$$k_{2,2} = f_2(0.25, 3.5, 6.45) = 1.715$$

These are used to determine the second set of midpoint predictions,

$$y_1 + k_{2,1} \frac{h}{2} = 4 + (-1.75) \frac{0.5}{2} = 3.5625$$

$$y_2 + k_{2,2} \frac{h}{2} = 6 + (1.715) \frac{0.5}{2} = 6.42875$$

which can be used to compute the second set of midpoint slopes,

$$k_{3,1} = f_1(0.25, 3.5625, 6.42875) = -1.78125$$

$$k_{3,2} = f_2(0.25, 3.5625, 6.42875) = 1.715125$$

These are used to determine the predictions at the end of the interval

$$y_1 + k_{3,1} h = 4 + (-1.78125)(0.5) = 3.109375$$

$$y_2 + k_{3,2} h = 6 + (1.715125)(0.5) = 6.857563$$

which can be used to compute the endpoint slopes,

$$k_{4,1} = f_1(0.5, 3.109375, 6.857563) = -1.554688$$

$$k_{4,2} = f_2(0.5, 3.109375, 6.857563) = 1.631794$$

The values of  $k$  can then be used to compute [Eq. (25.40)]:

$$y_1(0.5) = 4 + \frac{1}{6}[-2 + 2(-1.75 - 1.78125) - 1.554688]0.5 = 3.115234$$

$$y_2(0.5) = 6 + \frac{1}{6}[1.8 + 2(1.715 + 1.715125) + 1.631794]0.5 = 6.857670$$

Proceeding in a like manner for the remaining steps yields

<b>x</b>	<b>y<sub>1</sub></b>	<b>y<sub>2</sub></b>
0	4	6
0.5	3.115234	6.857670
1.0	2.426171	7.632106
1.5	1.889523	8.326886
2.0	1.471577	8.946865

\*\*\*

## Numerical Methods

11

- All previous methods discussed so far to solve the IVP are one-step methods which utilize information at a single point  $x_i$  to predict a value of dependent variable  $y_{i+1}$  at a future point  $x_{i+1}$ . The multistep methods are based on the insight that, once the computation has begun, information from previous points can be used to estimate the function at a future point.

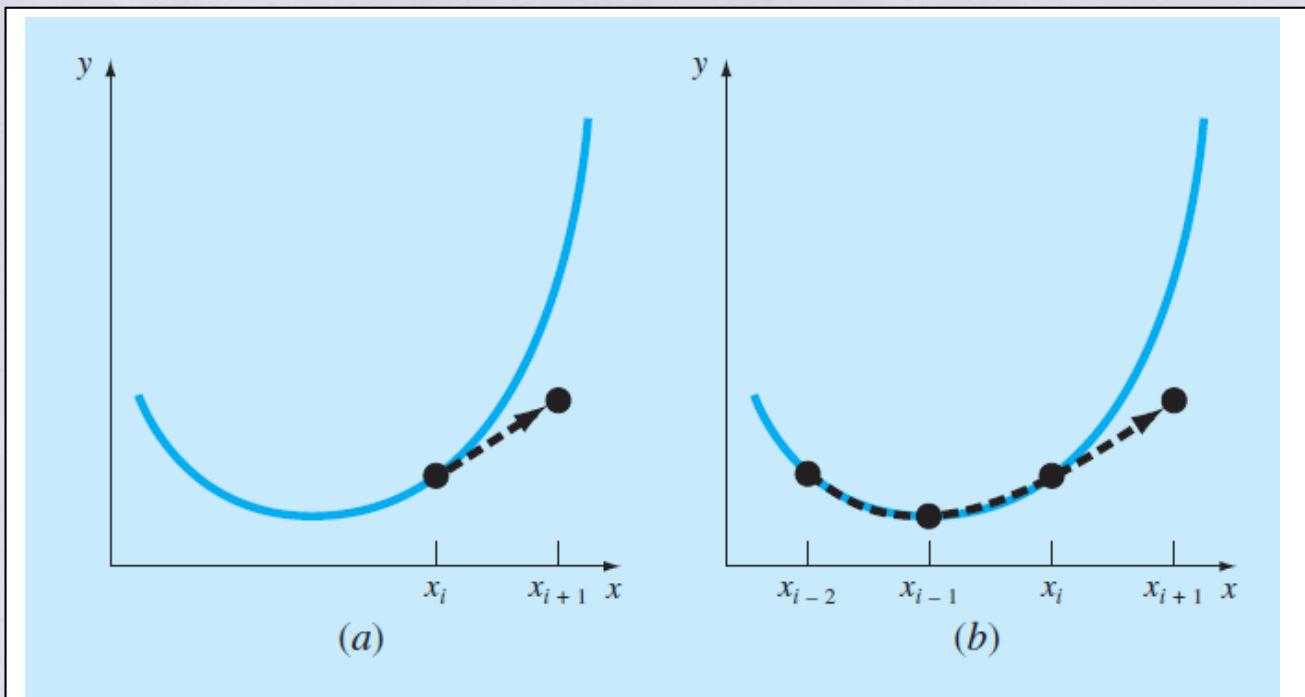


Fig. 1: Graphical depiction of the fundamental difference between (a) one-step and (b) multistep methods for solving ODEs.

## The non-self-starting Heun Method:

- Recall (from Lecture-21) that Heun approach uses Euler's method as a predictor:

$$y_{i+1}^0 = y_i + f(x_i, y_i)h \quad \dots \quad (1)$$

And the trapezoidal rule as a corrector:

$$y_{i+1} = y_i + \left[ \frac{f(x_i, y_i) + f(x_{i+1}, y_{i+1}^0)}{2} \right] h \quad \dots \quad (2)$$

Note that here, predictor step has local truncation error of  $\mathcal{O}(h^2)$  and corrector step has local trunc. error of  $\mathcal{O}(h^3)$

$\Rightarrow$  Predictor is weak.

- We can improve Heun's method by devising a predictor that has a local truncation error of  $\mathcal{O}(h^3)$

New predictor:

$$y_{i+1}^0 = y_{i-1} + f(x_i, y_i)2h \quad \dots \quad (3)$$

which have trunc. error of  $\mathcal{O}(h^3)$ .

- 13
- Note that Eq. (3) is not self-starting, as it involves  $y_{i-1}$  (which would not be available). Due to this reason, Eqs. (2) and (3) are called as non-self-starting Heun method:

$$\text{Predictor: } \hat{y}_{i+1}^0 = y_i^m + f(x_i, y_i^m) \cdot 2h$$

$$\text{Corrector: } \hat{y}_{i+1}^j = y_i^m + \left[ \frac{f(x_i, y_i^m) + f(x_{i+1}, \hat{y}_{i+1}^{j-1})}{2} \right] h$$

$$- j = 1, 2, \dots, m$$

Where the corrector is applied iteratively from  $j=1$  to  $m$  to obtain refined solutions.

- The approximate percentage relative error is:

$$|\varepsilon_a| = \left| \frac{\hat{y}_{i+1}^j - \hat{y}_{i+1}^{j-1}}{\hat{y}_{i+1}^j} \right| \times 100\%$$

- Termination criterion:  $|\varepsilon_a| < \varepsilon_s$  (Tolerance)

□

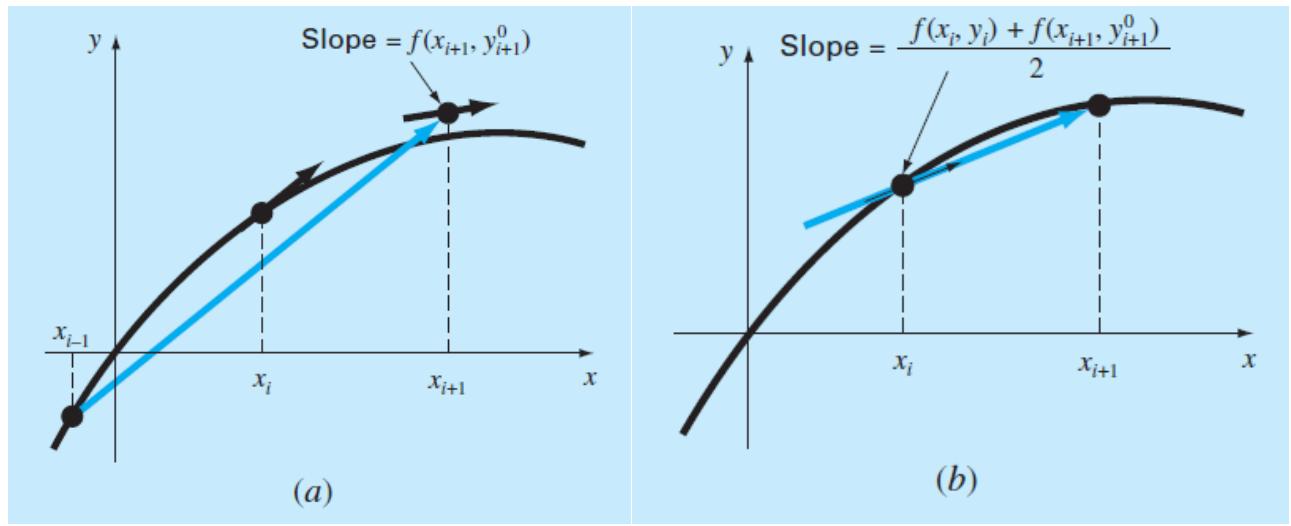


Fig. 2: Graphical depiction of the non-self-starting Heun method. (a) The midpoint method that is used as a predictor. (b) The trapezoidal rule that is employed as a corrector.

Example: Use the non-self-starting Heun method to integrate:

$$\left. \begin{aligned} y' &= 4e^{0.8x} - 0.5y \\ y(0) &= 2 \end{aligned} \right\}$$

with step-size  $h = 1.0$ . Additional information is required for the multistep method:  $y = -0.3929953$  at  $x = 1.0$ .

Sol:  $x_{-1} = -1, y_{-1} = -0.3929953, x_0 = 0, y_0 = 2$

Step-I:  $x_1 = x_0 + h$

Predictor step is used to extrapolate linearly from  $x_{-1}$  to  $x_1$ :

$$\begin{aligned} y_1^0 &= y_{-1} + f(x_0, y_0) 2h \\ &= -0.3929953 + (4e^{0.8 \times 0} - 0.5 \times 2) \times 2 \times 1 \\ &= 5.607005 \end{aligned}$$

The corrector is then used to compute the value:

When  $j = 1$ :

14

LS

$$y_1' = y_0 + \left[ \frac{f(x_0, y_0) + f(x_1, y_1^0)}{2} \right] h$$

$$= 6.549331$$

Approximate percentage relative error:

$$\epsilon_a = \left| \frac{y_1' - y_1^0}{y_1'} \right| \times 100 \% = 14.39 \%$$

- When  $j=2$ :

$$y_1^2 = y_0 + \left[ \frac{f(x_0, y_0) + f(x_1, y_1^1)}{2} \right] h$$

$$= 6.313749$$

Approximate percentage relative error:

$$\epsilon_a = \left| \frac{y_1^2 - y_1^1}{y_1^2} \right| \times 100 \% = 3.73 \%$$

Step 2:

□

## Numerical Methods

11

### - Numerical sol. of Partial Differential Equations

- Three types of fundamental PDEs:

A) Elliptic:  $\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0$ ,  $T = T(x, y)$

Laplace Equation — ①

B) Parabolic:  $\frac{\partial T}{\partial t} = \lambda \frac{\partial^2 T}{\partial x^2}$ ,  $T = T(x, t)$

Heat Conduction Equation — ②

C) Hyperbolic:  $\frac{\partial^2 U}{\partial t^2} = c^2 \frac{\partial^2 U}{\partial x^2}$ ,  $U = U(x, t)$

Wave Equation — ③

### Elliptic PDEs: Numerical Sol.

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0 \quad (\text{Laplace})$$

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = f(x, y) \quad (\text{Poisson})$$

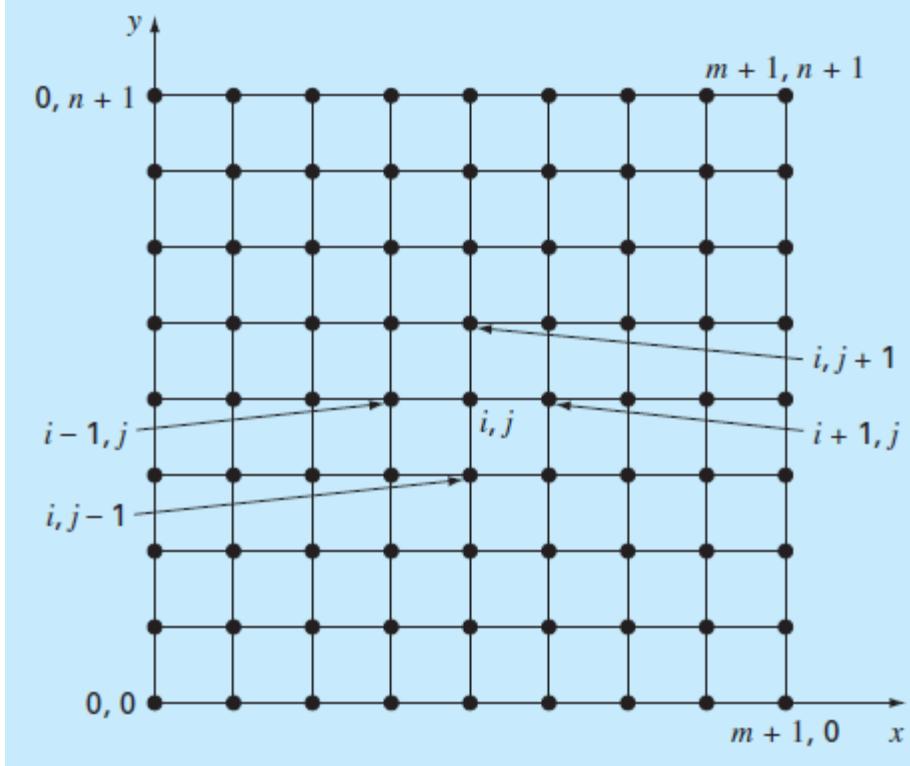


Fig. 1: A grid used for the finite-difference solution of elliptic PDEs in two independent variables- Laplace equation.

Consider 1D Laplace equation:

12

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0 \quad \dots \quad (1)$$

Use second-order centered difference formula to discretize eqn. (1)

$$\frac{\partial^2 T}{\partial x^2} = \frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{\Delta x^2} + O(\Delta x^2)$$

and  $\frac{\partial^2 T}{\partial y^2} = \frac{T_{i,j+1} - 2T_{i,j} + T_{i,j-1}}{\Delta y^2} + O(\Delta y^2)$

Substitute in Eq. (1):

$$\frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{\Delta x^2} + \frac{T_{i,j+1} - 2T_{i,j} + T_{i,j-1}}{\Delta y^2} = 0$$

If  $\Delta x = \Delta y$

$$\Rightarrow T_{i+1,j} + T_{i-1,j} + T_{i,j+1} + T_{i,j-1} - 4T_{i,j} = 0 \quad \dots \quad (4)$$

- Difference Equation.

## Boundary Conditions:

(3)

Boundary conditions along the domain boundaries must be specified to obtain a unique sol.

- If boundary condition is specified on the unknown ( $T$ )
  - Dirichlet Boundary Condition. ✓
- If boundary condition is specified on the derivative of the unknown ( $T$ )
  - Neumann Boundary Condition.

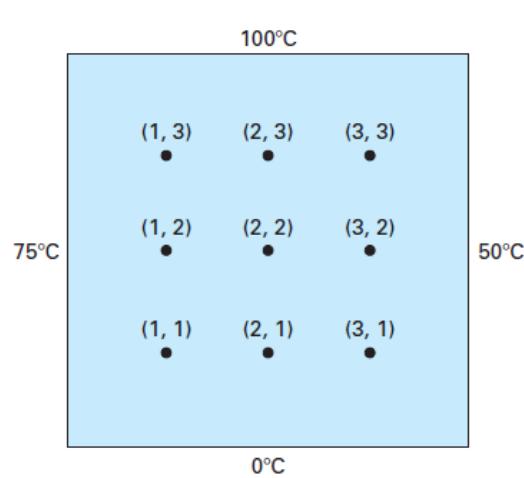


Fig. 2: A heated plate where boundary temperatures are held at constant levels- Dirichlet boundary condition.

∴ For node (1,1) : (from (4)) :

$$T_{2,1} + T_{0,1} + T_{1,2} + T_{1,0} - 4T_{1,1} = 0 \quad \boxed{5}$$

$$\therefore T_{0,1} = 75 \text{ and } T_{1,0} = 0 \Rightarrow -4T_{1,1} + T_{1,2} + T_{2,1} = -75$$

(4)

- Similar equations can be developed for the other interior points.
- The resultant system of 9 simultaneous equations with 9 unknowns is given as:

$$\begin{array}{ccccccccc}
 4T_{11} & -T_{21} & -T_{12} & & & & & = & 75 \\
 -T_{11} & +4T_{21} & -T_{31} & -T_{22} & & & & = & 0 \\
 & -T_{21} & +4T_{31} & & -T_{32} & & & = & 50 \\
 -T_{11} & & +4T_{12} & -T_{22} & -T_{13} & & & = & 75 \\
 & -T_{21} & -T_{12} & +4T_{22} & -T_{32} & -T_{23} & & = & 0 \\
 & & -T_{31} & -T_{22} & +4T_{32} & & -T_{33} & = & 50 \\
 & & -T_{12} & & +4T_{13} & -T_{23} & & = & 175 \\
 & & -T_{22} & & -T_{13} & +4T_{23} & -T_{33} & = & 100 \\
 & & -T_{32} & & -T_{23} & & +4T_{33} & = & 150
 \end{array}$$

OR

$$A \begin{bmatrix} T_{1,1} \\ T_{2,1} \\ T_{3,1} \\ T_{1,2} \\ T_{2,2} \\ T_{3,2} \\ T_{1,3} \\ T_{2,3} \\ T_{3,3} \end{bmatrix} = \begin{bmatrix} b \end{bmatrix} \quad \text{--- (6)}$$

- The Liebmann method (To solve system ⑥)

- based on Gauss-Seidel Method

- Here, Eq.(4) is expressed as:

$$T_{i,j} = \frac{T_{i+1,j} + T_{i-1,j} + T_{i,j+1} + T_{i,j-1}}{4} \quad \text{--- (7)}$$

and (7) is solved iteratively, for  $j = 1$  to  $n$   
 $i = 1$  to  $m$

As, system (4) is diagonally dominant  $\Rightarrow$  convergence of iterative process.

- Sometimes Overrelaxation is employed to accelerate the rate of convergence by applying the following formula after each iteration:

$$\bar{T}_{i,j}^{\text{new}} = \alpha T_{i,j}^{\text{new}} + (1-\alpha) T_{i,j}^{\text{old}} \quad \text{--- (8)}$$

where  $T_{i,j}^{\text{new}}$  - present iteration value.

$T_{i,j}^{\text{old}}$  - previous iteration value.

$\alpha$  - relaxation parameter and  $\alpha \in (0, 2)$

 overrelaxation

15

— Stopping criterion:

$$|(\varepsilon_a)_{i,j}| = \left| \frac{T_{i,j}^{\text{new}} - T_{i,j}^{\text{old}}}{T_{i,j}^{\text{new}}} \right| * 100 \%$$

If  $|(\varepsilon_a)_{i,j}| \leq \varepsilon_s$

↓  
user defined tolerance

□

## Liebmann Method

7

Example 1: Use Liebmann's method (Gauss-Seidel) to solve the temperature of the heated plate (given in Fig. 2).

Use overrelaxation with a value of 1.5 for the weighting factor and iterate to  $\epsilon_s = 1\%$ .

From Eq. (7), at  $(i,j) = (1,1)$ :

$$T_{1,1} = \frac{0 + 75 + 0 + 0}{4} = 18.75$$

Apply overrelaxation (given in Eq. (8)):

$$T_{1,1} = 1.5(18.75) + (1-1.5)0 = 28.125$$

Similarly, at  $(i,j) = (2,1)$ :

$$T_{2,1} = \frac{0 + 28.125 + 0 + 0}{4} = 7.03125$$

$$T_{2,1} = 1.5(7.03125) + (1-1.5)0 = 10.54688$$

For  $(i,j) = (3,1)$ :

$$T_{3,1} = \frac{50 + 10.54688 + 0 + 0}{4} = 15.13672$$

$$T_{3,1} = 1.5(15.13672) + (1-1.5)0 = 22.70508$$

The computation is repeated for the other rows to give

$$\begin{array}{lll} T_{12} = 38.67188 & T_{22} = 18.45703 & T_{32} = 34.18579 \\ T_{13} = 80.12696 & T_{23} = 74.46900 & T_{33} = 96.99554 \end{array}$$

Because all the  $T_{i,j}$ 's are initially zero, all  $\varepsilon_a$ 's for the first iteration will be 100%.

For the second iteration the results are

$$\begin{array}{lll} T_{11} = 32.51953 & T_{21} = 22.35718 & T_{31} = 28.60108 \\ T_{12} = 57.95288 & T_{22} = 61.63333 & T_{32} = 71.86833 \\ T_{13} = 75.21973 & T_{23} = 87.95872 & T_{33} = 67.68736 \end{array}$$

The error for  $T_{1,1}$  can be estimated as [Eq. (29.13)]

$$|(\varepsilon_a)_{1,1}| = \left| \frac{32.51953 - 28.12500}{32.51953} \right| 100\% = 13.5\%$$

Because this value is above the stopping criterion of 1%, the computation is continued. The ninth iteration gives the result

$$\begin{array}{lll} T_{11} = 43.00061 & T_{21} = 33.29755 & T_{31} = 33.88506 \\ T_{12} = 63.21152 & T_{22} = 56.11238 & T_{32} = 52.33999 \\ T_{13} = 78.58718 & T_{23} = 76.06402 & T_{33} = 69.71050 \end{array}$$

where the maximum error is 0.71%.

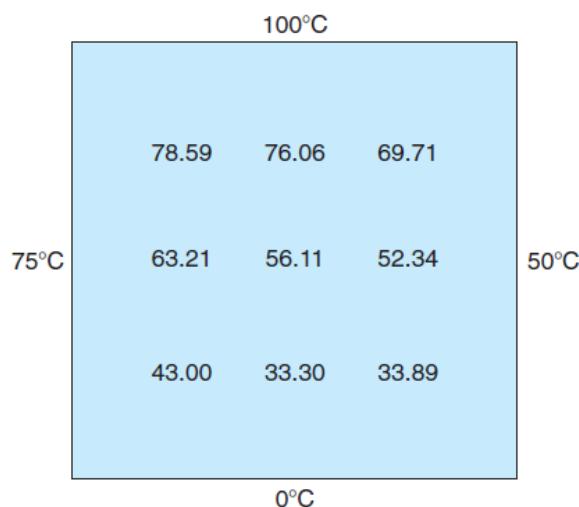


Fig. 3: Temperature distribution for a heated plate subject to fixed boundary conditions.

## Numerical Methods

11

### - Numerical sol. of Partial Differential Equations

#### - Three types of fundamental PDEs:

A) Elliptic:  $\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0$ ,  $T = T(x, y)$

Laplace Equation — ①

B) Parabolic:  $\frac{\partial T}{\partial t} = \lambda \frac{\partial^2 T}{\partial x^2}$ ,  $T = T(x, t)$

Heat Conduction Equation — ②

C) Hyperbolic:  $\frac{\partial^2 U}{\partial t^2} = c^2 \frac{\partial^2 U}{\partial x^2}$ ,  $U = U(x, t)$

Wave Equation — ③

### Elliptic PDEs: Numerical Sol.

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0 \quad (\text{Laplace})$$

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = f(x, y) \quad (\text{Poisson})$$

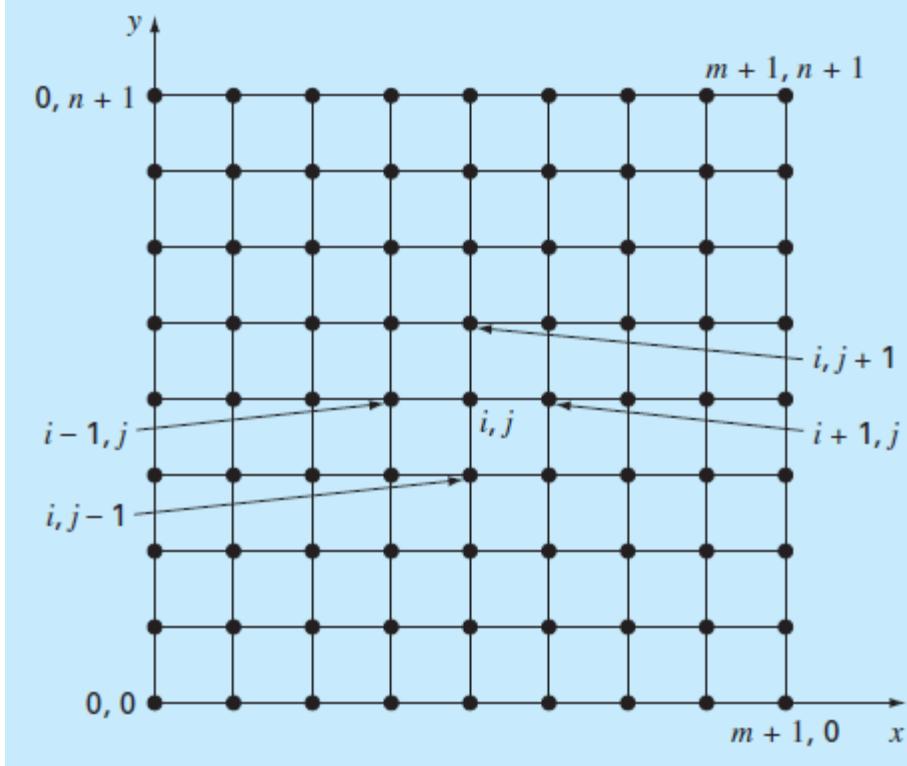


Fig. 1: A grid used for the finite-difference solution of elliptic PDEs in two independent variables- Laplace equation.

Consider 1D Laplace equation:

12

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0 \quad \dots \quad (1)$$

Use second-order centered difference formula to discretize eqn. (1)

$$\frac{\partial^2 T}{\partial x^2} = \frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{\Delta x^2} + O(\Delta x^2)$$

and  $\frac{\partial^2 T}{\partial y^2} = \frac{T_{i,j+1} - 2T_{i,j} + T_{i,j-1}}{\Delta y^2} + O(\Delta y^2)$

Substitute in Eq. (1):

$$\frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{\Delta x^2} + \frac{T_{i,j+1} - 2T_{i,j} + T_{i,j-1}}{\Delta y^2} = 0$$

If  $\Delta x = \Delta y$

$$\Rightarrow T_{i+1,j} + T_{i-1,j} + T_{i,j+1} + T_{i,j-1} - 4T_{i,j} = 0 \quad \dots \quad (4)$$

- Difference Equation

## Boundary Conditions:

(3)

Boundary conditions along the domain boundaries must be specified to obtain a unique sol.

- If boundary condition is specified on the unknown ( $T$ )
  - Dirichlet Boundary Condition. ✓
- If boundary condition is specified on the derivative of the unknown ( $T$ )
  - Neumann Boundary Condition.

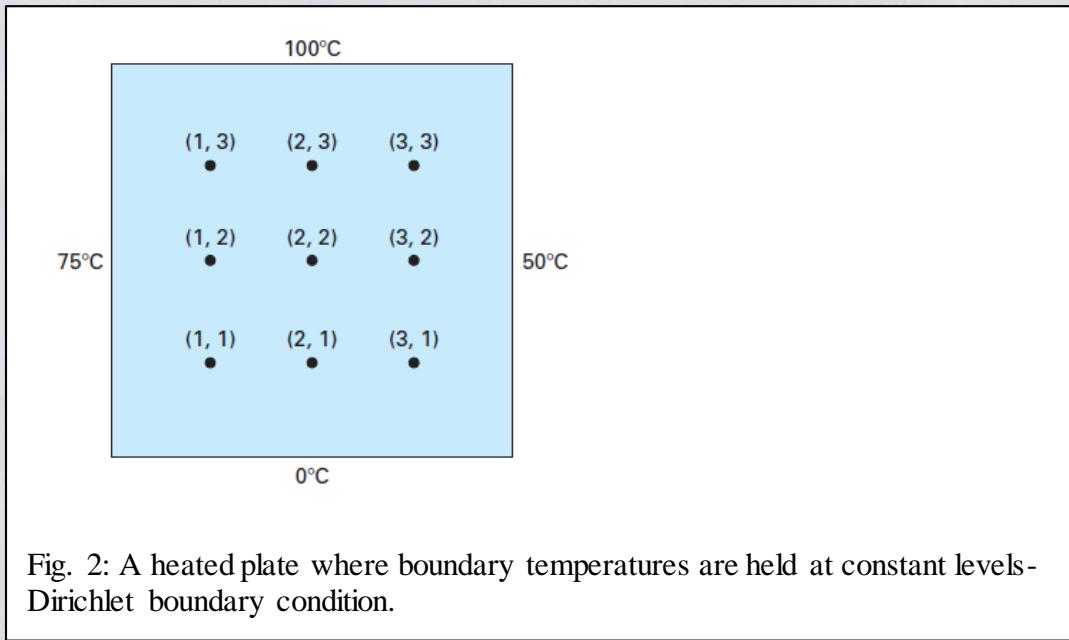


Fig. 2: A heated plate where boundary temperatures are held at constant levels- Dirichlet boundary condition.

$\therefore$  For node  $(1, 1)$  : (from (4)) :

$$T_{2,1} + T_{0,1} + T_{1,2} + T_{1,0} - 4T_{1,1} = 0 \quad \boxed{5}$$

$$\therefore T_{0,1} = 75 \text{ and } T_{1,0} = 0 \Rightarrow -4T_{1,1} + T_{1,2} + T_{2,1} = -75$$

(4)

- Similar equations can be developed for the other interior points.
- The resultant system of 9 simultaneous equations with 9 unknowns is given as:

$$\begin{array}{ccccccccc}
 4T_{11} & -T_{21} & -T_{12} & & & & & = & 75 \\
 -T_{11} & +4T_{21} & -T_{31} & -T_{22} & & & & = & 0 \\
 & -T_{21} & +4T_{31} & & -T_{32} & & & = & 50 \\
 -T_{11} & & +4T_{12} & -T_{22} & -T_{13} & & & = & 75 \\
 & -T_{21} & -T_{12} & +4T_{22} & -T_{32} & -T_{23} & & = & 0 \\
 & & -T_{31} & -T_{22} & +4T_{32} & & -T_{33} & = & 50 \\
 & & -T_{12} & & +4T_{13} & -T_{23} & & = & 175 \\
 & & -T_{22} & & -T_{13} & +4T_{23} & -T_{33} & = & 100 \\
 & & -T_{32} & & -T_{23} & & +4T_{33} & = & 150
 \end{array}$$

OR

$$A \begin{bmatrix} T_{1,1} \\ T_{2,1} \\ T_{3,1} \\ T_{1,2} \\ T_{2,2} \\ T_{3,2} \\ T_{1,3} \\ T_{2,3} \\ T_{3,3} \end{bmatrix} = \begin{bmatrix} b \end{bmatrix} \quad \text{--- (6)}$$

- The Liebmann method (To solve system ⑥)

- based on Gauss-Seidel Method

- Here, Eq.(4) is expressed as :

$$T_{i,j} = \frac{T_{i+1,j} + T_{i-1,j} + T_{i,j+1} + T_{i,j-1}}{4} \quad \text{--- (7)}$$

and (7) is solved iteratively, for  $j = 1$  to  $n$   
 $i = 1$  to  $m$

As, system (4) is diagonally dominant  $\Rightarrow$  convergence of iterative process.

- Sometimes Overrelaxation is employed to accelerate the rate of convergence by applying the following formula after each iteration :

$$\bar{T}_{i,j}^{\text{new}} = \alpha T_{i,j}^{\text{new}} + (1-\alpha) T_{i,j}^{\text{old}} \quad \text{--- (8)}$$

where  $T_{i,j}^{\text{new}}$  - present iteration value.

$T_{i,j}^{\text{old}}$  - previous iteration value.

$\alpha$  - relaxation parameter and  $\alpha \in (0, 2)$

 overrelaxation

16

— Stopping criterion:

$$|(\varepsilon_a)_{i,j}| = \left| \frac{T_{i,j}^{\text{new}} - T_{i,j}^{\text{old}}}{T_{i,j}^{\text{new}}} \right| * 100 \%$$

If  $|(\varepsilon_a)_{i,j}| \leq \varepsilon_s$

↓  
user defined tolerance

□

4

## - Heat Conduction Equation :

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

k - thermal conductivity.

---

————— ①

## - Explicit Numerical Methods

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

## Discretization :

$$\left. \frac{\partial^2 T}{\partial x^2} \right|_i^l = \frac{T_{i+1}^l - 2T_i^l + T_{i-1}^l}{\Delta x^2} + O(\Delta x^2) \quad \xrightarrow{\text{---}} \quad (2)$$

$$\frac{\partial T}{\partial t} = \frac{T_i^{l+1} - T_i^l}{\Delta t} + O(\Delta t) \quad \text{--- (3)}$$

Then from ① "

$$\frac{T_i^{l+1} - T_i^l}{\Delta t} = k \left( \frac{T_{i+1}^l - 2T_i^l + T_{i-1}^l}{\Delta x^2} \right)$$

OR

$$T_i^{l+1} = T_i^l + d \left( T_{i+1}^l - 2T_i^l + T_{i-1}^l \right) \quad \text{--- (4)}$$

where  
 $d = k \Delta t$   
 $\frac{\Delta x^2}{\lambda}$

12

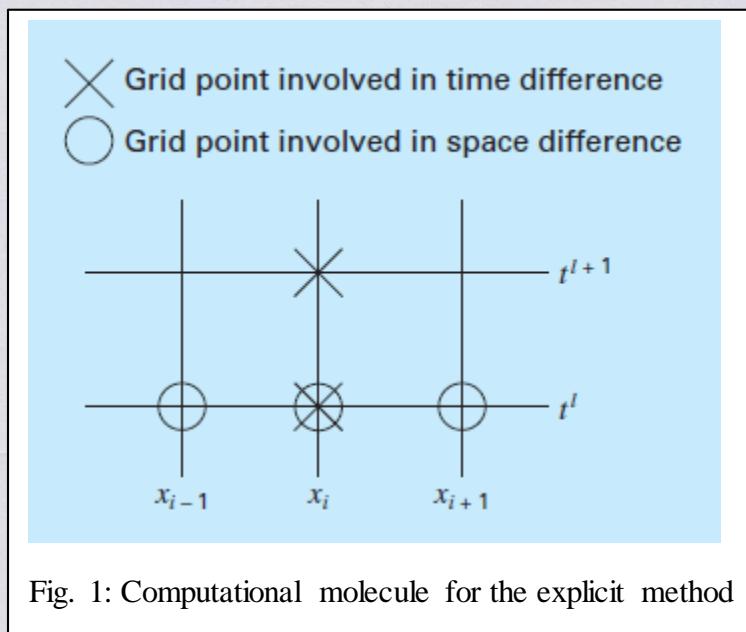


Fig. 1: Computational molecule for the explicit method

Example 2: Use explicit method to solve for the temperature distribution of a long, thin rod with a length 10 cm and the following values:

$$k = 0.035 \text{ cm}^2/\text{s}, \quad \Delta t = 0.1 \text{ s}, \quad \Delta x = 2 \text{ cm}$$

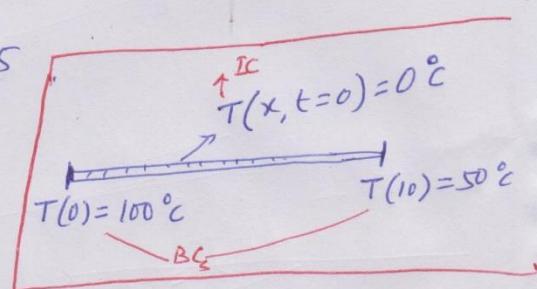
$$\Rightarrow \alpha = \frac{k \Delta t}{\Delta x^2} = 0.020875$$

Use Eq. (4):

$$T_i^{l+1} = T_i^l + \alpha \left( T_{i+1}^l - 2T_i^l + T_{i-1}^l \right)$$

$$T_1^l = T_0^0 + \alpha \left( T_2^0 - 2T_1^0 + T_0^0 \right)$$

$$= 0 + 0.020875(0 - 2(0) + 100) = 2.0875$$



At the other interior points,  $x = 4$ ,  $6$ , and  $8$  cm, the results are

$$T_2^1 = 0 + 0.020875[0 - 2(0) + 0] = 0$$

$$T_3^1 = 0 + 0.020875[0 - 2(0) + 0] = 0$$

$$T_4^1 = 0 + 0.020875[50 - 2(0) + 0] = 1.0438$$

At  $t = 0.2$  s, the values at the four interior nodes are computed as

$$T_1^2 = 2.0875 + 0.020875[0 - 2(2.0875) + 100] = 4.0878$$

$$T_2^2 = 0 + 0.020875[0 - 2(0) + 2.0875] = 0.043577$$

$$T_3^2 = 0 + 0.020875[1.0438 - 2(0) + 0] = 0.021788$$

$$T_4^2 = 1.0438 + 0.020875[50 - 2(1.0438) + 0] = 2.0439$$

- - -

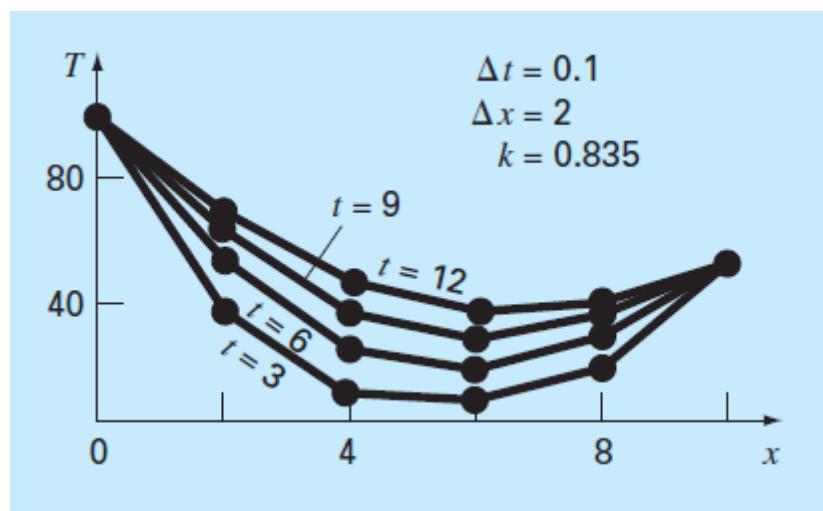


Fig. 2: Temperature distribution in a long, thin rod as computed with the explicit method.

13

Implicit Method :

$$\frac{\partial T}{\partial t} = k \frac{\partial^2 T}{\partial x^2} \quad \text{--- (5)}$$

$$\left. \begin{aligned} \frac{\partial^2 T}{\partial x^2} &= \frac{T_{i+1}^{l+1} - 2T_i^{l+1} + T_{i-1}^{l+1}}{\Delta x^2} \\ \frac{\partial T}{\partial t} &= \frac{T_i^{l+1} - T_i^l}{\Delta t} \end{aligned} \right\} \quad \text{--- (6)}$$

From (5) :

$$-\alpha T_{i-1}^{l+1} + (1+2\alpha) T_i^{l+1} - \alpha T_{i+1}^{l+1} = T_i^l \quad \text{--- (7)}$$

$$\text{where } \alpha = \frac{k \Delta t}{\Delta x^2}$$

At left end of rod ( $i=0$ ) .

$$T_0^{l+1} = \underbrace{f_0(t^{l+1})}_{\downarrow}$$

describes how the boundary  $t_{\text{amb}}$  changes with time

$\therefore \underline{i=1}$

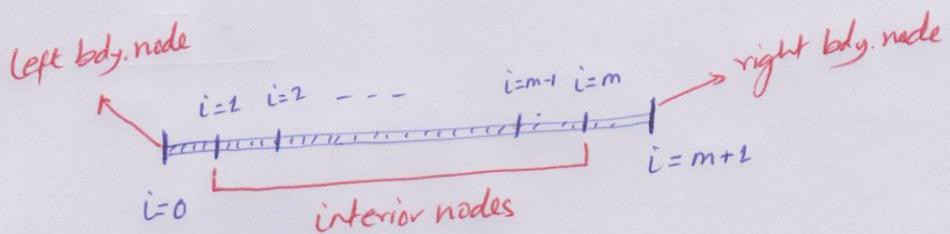
$$(1+2\alpha) T_1^{l+1} - \alpha T_2^{l+1} = T_1^l + \alpha f_0(t^{l+1}) \quad \text{--- (8)}$$

14

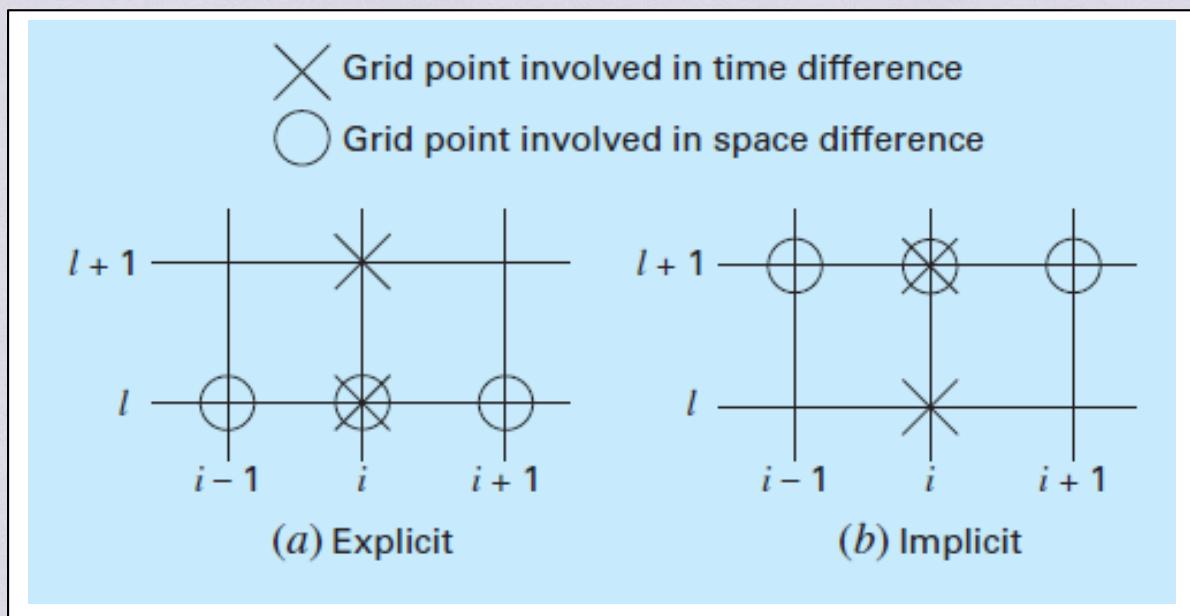
Similarly, at last interior node ( $i=m$ ):

$$-dT_{m-1}^{l+1} + (1+2\alpha)T_m^{l+1} = T_m^l + \alpha f_{m+1}(t^{l+1}) \quad \text{--- (9)}$$

Where,  $f_{m+1}(t^{l+1})$  describes the specified temp. change at the right boundary of the nod.

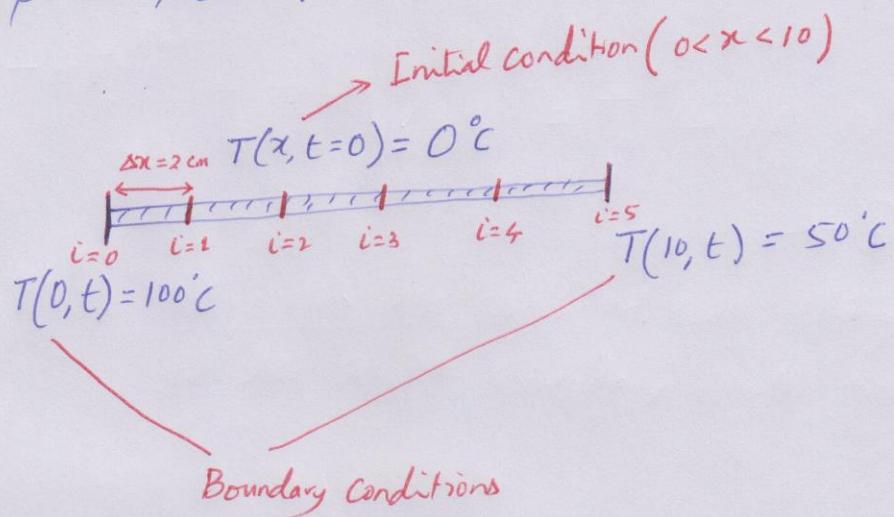


Equations (7), (8) and (9) are for interior nodes  $i$ ,  $i=1, 2, \dots, m$ . That is  $m$  linear algebraic equations in  $m$  unknowns.



15

Example 2: Solve the problem given in Example 1 using implicit method.



$$\lambda = \frac{k \Delta t}{\Delta x^2} = 0.020875$$

Using Eq. ⑧ at  $t=0$ : That is for 1<sup>st</sup> interior node ( $i=1$ )

$$1.04175 T_1^1 - 0.020875 T_2^1 = 0 + 0.020875(100)$$

OR

$$1.04175 T_1^1 - 0.020875 T_2^1 = 2.0875$$

Similarly, using Eqs. ⑦ and ⑨, we have system of 4 equations in 4 unknowns.

$$\begin{bmatrix} 1.04175 & -0.020875 & & \\ -0.020875 & 1.04175 & -0.020875 & \\ & -0.020875 & 1.04175 & -0.020875 \\ & & -0.020875 & 1.04175 \end{bmatrix} \begin{Bmatrix} T_1^1 \\ T_2^1 \\ T_3^1 \\ T_4^1 \end{Bmatrix} = \begin{Bmatrix} 2.0875 \\ 0 \\ 0 \\ 1.04375 \end{Bmatrix}$$

which can be solved for the temperature at  $t = 0.1$  s:

$$T_1^1 = 2.0047$$

$$T_2^1 = 0.0406$$

$$T_3^1 = 0.0209$$

$$T_4^1 = 1.0023$$

To solve for the temperatures at  $t = 0.2$ , the right-hand-side vector must be modified to account for the results of the first step, as in

$$\begin{Bmatrix} 4.09215 \\ 0.04059 \\ 0.02090 \\ 2.04069 \end{Bmatrix}$$

The simultaneous equations can then be solved for the temperatures at  $t = 0.2$  s:

$$T_1^2 = 3.9305$$

$$T_2^2 = 0.1190$$

$$T_3^2 = 0.0618$$

$$T_4^2 = 1.9653$$