

# CSE 330

# Numerical Methods

A word cloud centered around the word "polynomial". Other prominent words include "matrix", "number", "interpolation", "system", "error", "method", "nodes", and "formula". The words are colored in various shades of brown, tan, and blue, and their sizes vary to represent frequency or importance.

<u>Direct Method</u> $x + 2y = 0 \quad \dots \quad (1)$ $2x - \pi y = 1 \quad \dots \quad (2)$ $2x + 4y = 0 \quad (3)$ $-2x + \pi y = -1 \quad (4)$	<u>Iterative Method</u> $\sin(\pi) = 1/2 - \frac{\pi^3}{3!} + \frac{(\pi^5)^2}{5!} - \frac{(\pi^7)^2}{7!} + \dots$ $\sin(\pi) = \pi - \frac{\pi^3}{3!} + \frac{\pi^5}{5!} - \frac{\pi^7}{7!} + \dots$ $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$ $1 \quad 2 \quad 3 \quad 4$ $\text{Term}$ $1^{\text{st}} \text{ term} \rightarrow 1/2$ $2^{\text{nd}} \text{ term} \rightarrow 0.912$ $3^{\text{rd}} \text{ term} \rightarrow 0.9999996$
---	--

Due to computer limitations, the value is approximated.  
The part that is left out, is called Truncation Error

# Back in Babylonian civilization history

Eliminating terms from an infinite no. of terms

\* Different computers have different capabilities

e.g. Computer 1:  $\pi \approx 3.14159$   
Computer 2:  $\pi \approx 3.14$

### Floating Point Approximation

e.g. 2.5762 (actual value), 2.58 (approximated by a computer)

$$\begin{array}{r} 2.58 \\ - 2.5762 \\ \hline 0.0038 \end{array}$$

Rounding Error

<u>Floating point arithmetic</u>	<u>Fixed point</u>
# Floating point	# Fixed point
$(10.215)_{10}$	$10.215$
point	$0.10215 \times 10^1$ exponent
$r = \pm(d_1 d_2 d_3 \dots d_{n-1} \underbrace{d_n \dots d_n}_{\text{Fraction}}) \beta$	Base

e.g.  $(10.1)_2 = 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1}$

$\beta = (2.5)_{10}$

$\beta = 10$

<u>Floating point Numbers</u>	<u>F C R</u> $(\infty, +\infty)$
$f = \pm(0.d_1 d_2 \dots d_m) \beta^e$	$\frac{d_1}{\beta} \dots \frac{d_m}{\beta}$ fraction or mantissa or significand
sign	exponent
$\underbrace{\dots}_{\text{Underflow}}$	$\overbrace{\dots}^{\text{overflow}}$
$0 \quad \frac{1}{4} \quad \frac{1}{2}$	$0 \quad \frac{1}{4} \quad \frac{1}{2}$

# Total no. of values:

$\# \beta=2, m=3, e_{\min}=-1, e_{\max}=2$   
 $e = \{-1, 0, 1, 2\}$

$0. \square \square \times 2$

$\# (0.111)_2 \times 2^2 = \frac{7}{2},$  (highest possible value)

$(0.100)_2 \times 2^{-1} = \frac{1}{4}$  (lowest possible value)

$0.100 \times 2^1 \Rightarrow \frac{1}{4}$   
 $0.111 \times 2^{-1} \Rightarrow \frac{5}{16}$   
 $0.110 \times 2^{-1} \Rightarrow \frac{3}{8}$   
 $0.111 \times 2^{-1} \Rightarrow \frac{7}{16}$

For  $e = -1, 0, 1, 2 \rightarrow$  total 16  
 16 for +ve value  
 16 for -ve value  
 i.e. Total 32 values

### Floating Point Arithmetic - 2

# IEEE Standard: To represent floating point number we'll use 64 bit architecture

<u>Sign</u> 1 bit	<u>Exponent</u> 11 bit	<u>Fraction</u> 52 bits
$(0.d_1 d_2 \dots d_m) \beta^e$		
$0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$		
$\beta = 2^{10}$		
$2^{10} = 1024$		
$e_{\min} = 0$		
$e_{\max} = 2047$		
Total 2048 (including 0)		

### Fixed-point numbers

$$x = \pm(d_1 d_2 \dots d_{k-1} d_k \dots d_n)_{\beta}$$

where,

$$d_i = \text{digit}, \beta = \text{number system's base}, d_1, d_2, \dots, d_n \in \{0, 1, 2, \dots, \beta-1\}$$

Example:

$$(10.1)_2 \Rightarrow (d_1 d_2 d_3)_{\beta} = 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1} = (2.5)_{10}$$

$$(-3.14)_2 \Rightarrow -(d_1 d_2 d_3)_{\beta} = -(3 \times 2^0 + 1 \times 2^{-1} + 4 \times 2^{-2}) = (-4.5)_{10}$$

### Floating-point numbers

$$F = \{\pm(0.d_1 d_2 \dots d_m)_{\beta} \times \beta^e \mid \beta, d_i, e \in \mathbb{Z}, 0 \leq d_i \leq \beta-1, e_{\min} \leq e \leq e_{\max}\}$$

Here,  $(0.d_1 d_2 \dots d_m)_{\beta}$  is called the fraction or significand or mantissa.

$\beta$  is the base, and  $e$  is the exponent.

Example:

$$(123.45)_2 = 0.12345 \times 10^3$$

#  $2^3 = 8 \therefore 0 \text{ to } 7$   
 Ans if  $\beta=10 \quad 0 \leq d_i \leq 9$

<u>Convention 1: Standard Form:</u> $F = \pm(0.d_1 d_2 \dots d_m)_{\beta} \times \beta^e$ where $d_1 \neq 0$	<u>IEEE</u>
<u>Convention 2: Normalized:</u> $F = \pm(0.1 d_2 \dots d_m)_{\beta} \times \beta^e$	
<u>IEEE</u> fraction/significand/mantissa (m)	
<u>Convention 3: Denormalized:</u> $F = \pm(1.d_1 d_2 \dots d_m)_{\beta} \times \beta^e$	

In previous semesters normalized & denormalized used to have the name of the other

### Cost of using Floating Point Number

1) Numbers are not equally spaced

2) We can only represent a finite set of numbers

## #IEEE standard (1985) for double-precision (64-bit) arithmetic

Here  $\beta = 2$ , and there are 52 bits for the fraction, 11 for the exponent, and 1 for the sign. The actual format used is -

(Denormalized Form)

here, the e (exponent) range is [0, 2047] since  $2^{11} = 2048$

The largest possible number is

The smallest possible number is

The largest possible number is  $(1.000 \dots 0)_2 \times 2^{2047}$  (not small; it cannot express 0.001)

If we want numbers less than 1, we can do **exponent bias**

### Example:

$$\pm (1.d_1d_2d_3 \dots d_{52})_2 \times 2^{e-1023} \{e = 1023\} \text{ (exponent bias)}$$

$$= \pm [0.d_1d_2d_3 \dots d_{52}]_2 \times 2^{e-1023} \text{ (Normalized form)}$$

Now, the e (exponent) range is [-1022, 1025]

The largest possible number is

$$= (0.1111 \dots 1)_2 \times 2^{1023} \approx \infty$$

The smallest possible number is

$$= (0.1000 \dots 0)_2 \times 2^{-1023} \approx 0$$

Actually, the exponents -1022 and 1025 are used to store  $\pm 0 \pm \infty$ , respectively.

So,

The largest possible number (except 0) is

$$= (0.1111 \dots 1)_2 \times 2^{1023} \approx 1.798 \times 2^{308}$$

The smallest possible number (except 0) is

$$= (0.1000 \dots 0)_2 \times 2^{-1023} \approx 2.225 \times 2^{-308}$$

### Significant figures

The concept of significant figures (s.f.) helps us express numbers with the right level of precision.

. How to Count Significant Figures?

✓ Start counting from the first non-zero digit.

✓ Count all digits after that, including zeros if they come after the decimal point.

### Example:

3.1056  $\rightarrow$  5 s.f. (all digits count)

31.050  $\rightarrow$  5 s.f. (the final zero counts because it's after the decimal)

0.031056  $\rightarrow$  5 s.f. (leading zeros don't count, but everything after '3' does)

0.031050  $\rightarrow$  5 s.f. (final zero still counts)

3105.0  $\rightarrow$  5 s.f. (the zero counts because of the decimal)

. How to Round to Significant Figures?

To round a number to n significant figures, find the first n digits, then round off.

### Example:

3.1415 rounded to 3 s.f.

$\rightarrow 3.14$

0.007823 rounded to 2 s.f.

$\rightarrow 0.0078$

45678 rounded to 3 s.f.  $\rightarrow$  45,700 (notice the zeros replace rounded-off digits)

However, significant figures help keep calculations accurate without unnecessary detail. For example, if a measurement is 2.36 cm, writing 2.3600 cm suggests extra precision that wasn't actually measured.

## Rounding

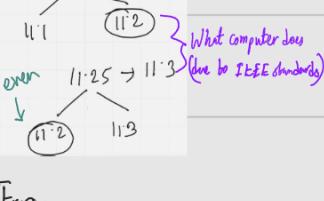
If number is 5, the number is rounded to nearest even number.

$$11.16 \rightarrow 11.2$$

$$11.14 \rightarrow 11.1$$

$$11.15 \rightarrow 11.2$$

what we actually do



## Rounding



For binary, if the number ends in 0  $\rightarrow$  even and if the number ends in 1  $\rightarrow$  odd

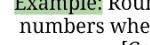
Example: Round  $(0.1001)_2 \times 2^3$  and

numbers where  $\beta = 2, m = 3, e_{min} = -1, e_{max} = 2$  [Convention 01]



Example: Round  $(0.1100100100)_2 \times 2^2$  numbers where

$\beta = 2, m = 4, e_{min} = -1, e_{max} = 2$  [Convention 01]



Example: Round  $(0.1100100100)_2 \times 2^2$  numbers where

$\beta = 2, m = 6, e_{min} = -1, e_{max} = 2$  [Convention 01]



## Rounding Error

$$x = 21.531 \Rightarrow$$

$$f(x) = 21.53 \Rightarrow$$

Error 0.001

$$\text{Error} = |f(x) - x| \quad \text{for percentage error}$$

$$\text{Relative error } (\delta) = \frac{|f(x) - x|}{|x|} \times 100$$

Scale Invariant error

$$\text{For maximum scale invariant error, } \delta_{\max} = \frac{|f(x) - x|}{|x|_{\max}}$$

$$\text{eg } x = 21.535 \quad \text{computer standard as } x = 21.55$$

$$f(x) = 21.51$$

$$0.005 \Rightarrow 0.5 \times 10^{-2}$$

$$(\text{digit} \dots) \times \beta^e \Rightarrow \frac{1}{2} \times \beta^{-m}$$

$$\frac{21.53}{21.5}$$

$$\frac{21.54}{21.5}$$

$$\frac{21.57}{21.6}$$

$$\text{Max}^m \text{ Error} = \frac{1}{2} \beta^{-m}$$

$$\text{Max}^m \text{ Possible Total Error} = \frac{1}{2} \times \beta^{-m} \times \beta^e$$

$$\frac{0.21535 \times 10^2}{0.21540 \times 10^2} \rightarrow x$$

$$\frac{0.21540 \times 10^2}{0.21545 \times 10^2} \rightarrow f(x)$$

$$\text{error} = 0.00005 \times 10^2 = 0.5 \times 10^{-4}$$

$$|f(x) - x| \leq \frac{1}{2} \beta^{-m} \beta^e$$

$$\frac{|f(x) - x|}{|x|} \leq \frac{1}{2} \beta^{-m} \beta^e$$

$$\delta \leq \frac{1}{2} \beta^{-m} \beta^e$$

$$\frac{1}{|x|} \leftarrow \text{round}$$

$$\text{For } \delta \text{ to be max}^m \text{ possible value}$$

$$\delta \leq \frac{1}{2} \beta^{-m} \beta^e$$

$$\delta = \frac{1}{2} \beta^{-m} \beta^e$$

$$\delta = 10^{-m} = \beta^{-m}$$

$$\delta = \frac{1}{2} \beta^{-m}$$

$$\therefore \delta \leq \frac{1}{2} \beta^{-m}$$

$$\text{Max}^m \text{ Relative error}$$

$$\delta_m = \frac{1}{2} \beta^{-m}$$

$$\text{E}_m (\text{Epsilon } m)$$

$$[\text{Machine Epsilon } m]$$

$$\alpha^2 - 56\alpha + 1 = 0$$

$$\alpha_1 = 28 + \sqrt{783} = 55.78$$

$$\alpha_2 = 28 - \sqrt{783} = 28 - 27.98$$

$$= 0.02$$

$$= 27.982137 \text{ actually}$$

$$= 0.02$$

$$= 27.982137 \text{ actually}$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

$$= 0.02$$

&lt;math display

# Polynomial Interpolation - I

2/10/2025

Friday



# This chapter deals with functions.

# The computer recognizes a function by converting it to a polynomial.

# We'll learn how this is done.

We recognize function as e.g.  $f(x) = \sin x$  and computer as e.g.  $P_n(x)$

# Polynomial: series of terms. e.g.  $2x^2 - 5x + 3$

→ n means degree (highest power of the series of terms)

→ Exponent must be Non-negative & an integer.

$$P_n(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots + a_n x^n$$

where  $\{a_0, a_1, a_2, \dots, a_n\}$  are coefficients/constants

$x \Rightarrow$  Variable

$n \Rightarrow$  degree

# Number of Terms

$$\text{e.g. } 7x^7 - 2x^3$$

max no. of terms with degree 7:  $7x^7 + 0 \cdot x^6 + 0 \cdot x^5 + 0 \cdot x^4 + 0 \cdot x^3 + 0 \cdot x^2 + 0 \cdot x^1$  [8 terms]

∴ for n degree

→ (n+1) terms max

→ (n+1) coefficients max

$$\text{e.g. } f(x) = \sin x$$

$$= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

Truncation error

degree ↑ Accuracy ↑

$$P_6(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} + 0 \cdot x^6 -$$

This computer shows till degree 6

Taylor Series

# As computer has limitations,

$$f(x) \neq P_n(x)$$

$$\text{e.g. } \sin x \neq x - \frac{x^3}{3!} + \frac{x^5}{5!}$$

$$\therefore |f(x) - P_n(x)| = \text{Truncation Error}$$

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x-x_0) + \frac{f''(x_0)}{2!}(x-x_0)^2 \dots$$

$$\text{If there was no error, } n=\infty \quad = \sum_{k=0}^n \frac{f^k(x_0)}{k!}(x-x_0)^k +$$

$$\boxed{\frac{f^{n+1}(x_0)}{(n+1)!}(x-x_0)^{n+1} + \frac{f^{n+2}(x_0)}{(n+2)!}(x-x_0)^{n+2} + \dots}$$

Error starts ↓  
Error Bound

$$\text{e.g. } f(x) = \sin x, n=6$$

$$P_6(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!}$$

$$\therefore \text{Error Bound} = \frac{f^{6+1}(\xi)}{(6+1)!}(x-x_0)^{6+1} \\ = \frac{-6!(\xi)}{7!}(x-x_0)^7$$

$$\begin{aligned} f(x) &= \sin x \\ f'(x) &= \cos x \\ f''(x) &= -\sin x \\ f'''(x) &= -\cos x \\ f''''(x) &= \sin x \\ f''''''(x) &= \cos x \\ f'''''''(x) &= -\sin x \\ f''''''''(x) &= -\cos x \end{aligned}$$

Put n/b/nodes ↑ Accuracy ↑

## Applications:

- Image processing (resizing images)
- Data analysis (filling missing data)
- Engineering (curve fitting)
- Computer graphics (animation and rendering)

Given  $(x_0, x_1, x_2, x_3, \dots, x_n) \Rightarrow n+1$  nodes  $\Rightarrow$  Polynomial will have degree n  
and  $f(x) = \boxed{\quad}$ , find  $P_n(x)$

$$\therefore P_n(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots + a_n x^n$$

$$\begin{aligned} P_n(x_0) &= a_0 + a_1 x_0 + \dots + a_n x_0^n \\ &\vdots \\ P_n(x_n) &= a_0 + a_1 x_n + \dots + a_n x_n^n \end{aligned}$$

$\left. \quad \right\} n+1 \text{ eqn}$

## Matrix Representation

$$\begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} f(x_0) \\ f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$$

$\underbrace{\quad}_{n+1} \quad \underbrace{\quad}_{n+1} \quad \underbrace{\quad}_{n+1}$

Vandermonde matrix

$$\begin{aligned} V \cdot a &= f \\ \therefore a &= V^{-1} f \end{aligned}$$

Q) Using Vandermonde matrix find the interpolating polynomial

(i) given  $f(x) = \frac{1}{x^2}$ ,  $x_0=2$ ,  $x_1=4$  two nodes

$$\# P_n(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots + a_n x^n \quad (\text{general formula of a polynomial})$$

$\#(n+1)$  nodes  $\rightarrow P_n$

$\therefore$  For two nodes degree is 1  $\therefore P_1(x) = a_0 + a_1 x$

$$\therefore P_1(x) = a_0 + a_1 x = a_0 + a_1 x$$

$$P_1(x) = 0.4375 + (-0.09375)x$$

(ii) Find interpolation error when  $x=9$

$$\text{Interpolation Error} = |f(x) - P_1(x)|$$

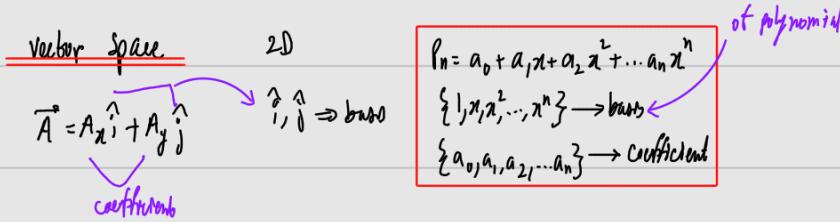
$$= |f(9) - P_1(9)| = \frac{1}{9^2} - (-0.40625) = 0.3939$$

$(n+1)$  nodes  $\rightarrow P_n \Rightarrow$  Existence/Uniqueness Theorem

Only square matrix  
can have inverse matrix

Use calculator

Till now nodes &  $P_n$  were given. Then we use inverse matrix to find the coefficients. But, using inverse matrix takes a lot of resources for computer. Takes a lot of processing time. Mathematician Lagrange gave a theorem to solve this problem without using inverse matrix.



### Lagrange Polynomial

Given  $(n+1)$  nodes,  $f(x)$ : coefficient  
 $P_n(x) = \sum_{k=0}^n f(x_k) \cdot l_k(x)$  Lagrange basis

$\{l_0(x), l_1(x), l_2(x), \dots, l_n(x)\} \rightarrow$  basis

$\{f(x_0), f(x_1), f(x_2), \dots, f(x_n)\} \rightarrow$  coefficient

e.g. Given Nodes = 2,  $f(x)$

$$n=1, P_1(x) = f(x_0) l_0(x) + f(x_1) l_1(x) \quad \square \text{ Product}$$

$\sum$  sum

$$l_k(x) = \prod_{j=0, j \neq k}^n \left( \frac{x-x_j}{x_k-x_j} \right) = \left( \frac{x-x_0}{x_k-x_0} \right) \cdot \left( \frac{x-x_1}{x_k-x_1} \right) \cdot \left( \frac{x-x_2}{x_k-x_2} \right) \cdots \left( \frac{x-x_n}{x_k-x_n} \right)$$

Given  $f(x) = \cos x$

$$\therefore P_2(x) = f(x_0) l_0(x) + f(x_1) l_1(x) + f(x_2) l_2(x) \quad n=2$$

$$x_0 = -\frac{\pi}{4} \rightarrow f(-\frac{\pi}{4}) = \frac{1}{\sqrt{2}} = f(x_0)$$

$$x_1 = 0 \rightarrow f(0) = 1 = f(x_1)$$

$$x_2 = \frac{\pi}{4} \rightarrow f(\frac{\pi}{4}) = \frac{1}{\sqrt{2}} = f(x_2)$$

$$\therefore l_0(x) = \left( \frac{x-x_0}{x_1-x_0} \right) \left( \frac{x-x_1}{x_2-x_1} \right) \left( \frac{x-x_2}{x_3-x_2} \right) = -\frac{1}{\pi^2} (x+\frac{\pi}{4})(x-\frac{\pi}{4})$$

$$\therefore l_1(x) = \left( \frac{x-x_0}{x_1-x_0} \right) \left( \frac{x-x_1}{x_2-x_1} \right) \left( \frac{x-x_2}{x_3-x_2} \right) = \frac{8}{\pi^2} x(x+\frac{\pi}{4})$$

$l_k(x) = \{ \rightarrow \text{Knockel's delta}$   
when  $k=j$

$$\therefore P_2(x) = f(x_0) l_0(x) + f(x_1) l_1(x) + f(x_2) l_2(x) \quad [\text{Put the values and simplify}]$$

# n+1 nodes Langrange coefficient  
 $\# P_n(x) = \sum_{k=0}^n f(x_k) l_k(x)$  Langrange basis

$$\text{eg } f(x) = x^2, P_2 = ?$$

$$(i) \text{ nodes, } x_0=2 \rightarrow f(x_0)=4$$

$$x_1=3 \rightarrow f(x_1)=9$$

Interpolating polynomial

$$\therefore P_1(x) = 4(-(x-3)) + 9(x-2) = -4x + 12 + 9x - 18 = 5x - 6$$

(ii) Find the error for  $x=5$

$$\text{Interpolation error} = |f(x) - P_n(x)| = |(5)^2 - (5(5) - 6)| = |25 - 19| = 6$$

★ (iii) given the function passes through  $(2, 4)$  and  $(3, 9)$

$$\begin{matrix} \downarrow \\ x_0, f(x_0) \end{matrix} \quad \begin{matrix} \downarrow \\ x_1, f(x_1) \end{matrix}$$

(iv) given function passes through  $(1, 2), (2, 4), (3, 9)$  how many nodes are there & what is the degree of the polynomial?

$$= \text{no. of nodes} = 3, \text{ degree} = 3-1 = 2$$

### Langrange

#### Advantages

- Fast Calculation

#### Disadvantages

- Non Dynamic

(Recalculation for adding new node)

Thus we use a different method  $\rightarrow$  Newton form

### Newton / Divided difference form (Dynamic)

$$\text{Node} = (n+1)+1$$

(new) Coll

$$P_{n+1}(x) = P_n(x) + j_{n+1}(x)$$

$$\text{eg } \frac{x^2+3x+5}{\text{new}} = \frac{x+3}{\text{old}} + x^2 + 2x + 2$$

$$\text{eg from 4 nodes } \{-1, 0, 1, 2\}$$

$$f(x) \{5, 1, 1, 12\}$$

$$P_3(x) = \sum_{k=0}^3 a_k n_k(x)$$

$$P_n(x) = \sum_{k=0}^n f(x_k) l_k(x)$$

$$\text{nodes : } P_n(x) = \sum_{k=0}^n a_k n_k(x)$$

$$= a_0 n_0(x) + a_1 n_1(x) + a_2 n_2(x) + \dots + a_n n_n(x)$$

$$= f(x_0) n_0(x) + f[x_0, x_1] n_1(x) + \dots + f[x_0, x_1, x_2, \dots, x_n] n_n(x)$$

$$n_0(x) = 1$$

$$n_1(x) = (x-x_0)$$

$$n_2(x) = (x-x_0)(x-x_1)$$

$$n_3(x) = (x-x_0)(x-x_1)(x-x_2)$$

$$n_k(x) = (x-x_0)(x-x_1) \dots (x-x_{k-1})$$

$$= a_0 n_0(x) + a_1 n_1(x) + a_2 n_2(x) + a_3 n_3(x)$$

$$= f(x_0) \cdot 1 + f[x_0, x_1] (x-x_0) + f[x_0, x_1, x_2] (x-x_0)(x-x_1) + f[x_0, x_1, x_2, x_3] (x-x_0)(x-x_1)(x-x_2) + f[x_0, x_1, x_2, x_3, x_4] (x-x_0)(x-x_1)(x-x_2)(x-x_3)$$

when new node added:

$$\begin{aligned} x_0 = -1 & \quad f[x_0] = 5 \\ x_1 = 0 & \quad f[x_1] = 1 \\ x_2 = 1 & \quad f[x_2] = 1 \\ x_3 = 2 & \quad f[x_3] = 11 \\ x_4 = 3 & \quad f[x_4] = 15 \end{aligned}$$

$$\begin{aligned} f[x_0, x_1] &= \frac{f[x_1] - f[x_0]}{x_1 - x_0} = \frac{1 - 5}{1 - (-1)} = -4 \\ f[x_0, x_1, x_2] &= \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0} = \frac{\frac{f[x_2] - f[x_1]}{x_2 - x_1} - f[x_0, x_1]}{x_2 - x_0} = \frac{\frac{1 - 1}{2 - 1} - 1}{2 - (-1)} = 2 \\ f[x_0, x_1, x_2, x_3] &= \frac{f[x_1, x_2, x_3] - f[x_0, x_1, x_2]}{x_3 - x_0} = \frac{\frac{f[x_2, x_3] - f[x_1, x_2]}{x_3 - x_1} - \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0}}{x_3 - x_0} = \frac{\frac{\frac{11 - 1}{3 - 1} - \frac{1 - 1}{1 - 0}}{3 - 1} - \frac{1 - 1}{1 - 0}}{3 - (-1)} = 5 \\ f[x_0, x_1, x_2, x_3, x_4] &= \frac{f[x_1, x_2, x_3, x_4] - f[x_0, x_1, x_2, x_3]}{x_4 - x_0} = \frac{\frac{f[x_2, x_3, x_4] - f[x_1, x_2, x_3]}{x_4 - x_1} - \frac{f[x_1, x_2, x_3] - f[x_0, x_1, x_2]}{x_3 - x_0}}{x_4 - x_0} = \frac{\frac{\frac{15 - 11}{4 - 1} - \frac{11 - 1}{3 - 1}}{4 - 1} - \frac{1 - 1}{1 - 0}}{4 - (-1)} = 10 \end{aligned}$$

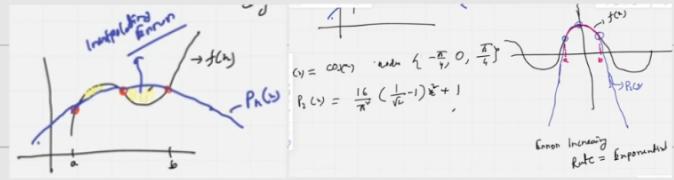
## Interpolation Error

Taylor

$$|f(x) - P_n(x)| = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x-x_0)^{n+1}$$

Cauchy's Theorem Error bound / Lagrange form of the remainder

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



Error Increasing Rate = Exponential

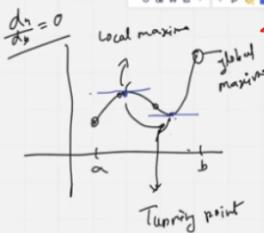
If interval not given, take  $a = -\pi/4, b = \pi/4$

# Max<sup>n</sup> possible error / Upper bound of error

$$|f(x) - P_2(x)| \leq \frac{t^3(\xi)}{6} \left( x^2 - \frac{\pi^2}{16} \right) x$$

$$\text{eg } f(x) = \cos x, \text{ nodes: } \{-\pi/4, 0, \pi/4\}, \text{ Interval } [a, b]$$

$$|f(x) - P_2(x)| = \left| \frac{t^3(\xi)}{3!} (x+\pi/4)(x-0)(x-\pi/4) \right|$$



We're to check for the internal values too in case they have a greater maximum value

$$\begin{aligned} \text{(i) Max: } & \frac{\sin(\xi)}{6} \\ \text{(ii) Max: } & x^2 - \frac{\pi^2}{16} x \rightarrow 3x^2 - \frac{\pi^2}{16} = 0 \quad \omega\left(\frac{\pi}{4}\right) = -0.186 \\ & \omega(x) = x^3 - \frac{\pi^2}{16} x \quad x = \pm \frac{\pi}{4\sqrt{3}} \quad \omega\left(-\frac{\pi}{4\sqrt{3}}\right) = 0.186 \\ & \omega'(x) = 3x^2 - \frac{\pi^2}{16} \quad \omega(-1) = 0.383 \\ & \text{at max/min point, } \omega'(x) = 0 \quad \omega(1) = -0.383 \end{aligned}$$

$$\therefore |f(x) - P_2(x)| \leq \left| \frac{0.186}{6} \times 0.383 \right| \leq 0.0537$$

$$\therefore \text{Max}^n \text{ error} = 0.0537,$$

We'll only deal with the absolute value

## Polynomial Interpolation - 5

02/03/2025

Sunday

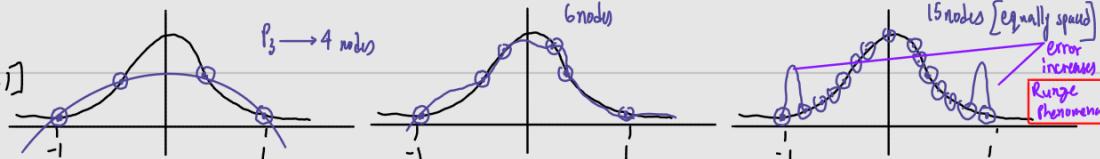
Convergence

nodes ↑ error ↓ → Not always true

$$|f(x) - P_n(x)|$$

∴ nodes ∞ → error = 0

$$\text{eg } f(x) = \frac{1}{1+25x^2}, \text{ Interval } [-1, 1]$$



Increasing the nodes at these endpoints will decrease the error.

Runge Phenomena

- (1) Depends on the function (Mainly symmetric function) Runge function
- (2) Depends on the nodes (Equally spaced)

Solution

(1) Piecewise Interpolation

↳ take smaller intervals and interpolate

↳ add them up.

(2) Non-equidistant nodes (Chebyshev's nodes) (instead of taking randomly)

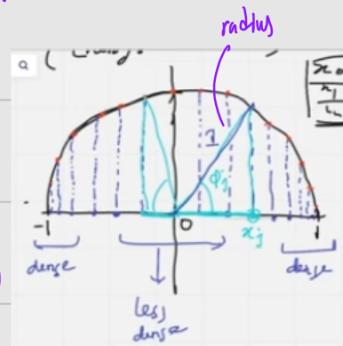
↳ draw a semi-circle

↳ put a no. of equidistant points and identify the corresponding nodes.

↳ Observation: It appears that the nodes are more dense around the interval points

↳ draw triangles with a perpendicular line joining the nodes with their corresponding points.

↳ Chebyshev relates the angles with the corresponding nodes:



$$\rho_j = \frac{(2j+1)\pi}{2(n+1)}$$

$$\cos \theta_j = \frac{x_j}{1}; \text{ radius} = 1$$

$$\Rightarrow x_j = \cos \theta_j$$

$$\therefore x_j = \cos \frac{(2j+1)\pi}{2(n+1)}$$

Runge function

$$\text{eg (i) } f(x) = \frac{1}{1+25x^2}, [-1, 1], n=3 \text{ find the nodes}$$

$$= n=3, \text{ so no. of nodes}=4 \quad \# x_0 = \cos \frac{(2 \times 0 + 1)\pi}{2(n+1)} = \cos \frac{\pi}{8}, x_1 = \cos \frac{3\pi}{8}, x_2 = \cos \frac{5\pi}{8}, x_3 = \cos \frac{7\pi}{8}$$

★ (ii) Using the nodes, define the Lagrange polynomial / Newton's Polynomial

# Polynomial Interpolation - 6

## Hermite Interpolation (aka Derivative Conditions)

Earlier: For  $n+1$  nodes,  
has  $n+1$  functions

$$\begin{bmatrix} x_0, x_1, x_2, \dots, x_n \\ f(x_0), f(x_1), \dots, f(x_n) \end{bmatrix} \quad \begin{array}{l} n \text{ degree} \\ \curvearrowright n+1 \text{ conditions} \end{array}$$

Hermite idea: has  $n+1$  first derivatives  $[f'(x_0), f'(x_1), \dots, f'(x_n)] \rightarrow n+1$  conditions

Total  $2n+2$  Conditions

#  $n+1$  nodes / conditions  $\rightarrow$  degree  $n$

#  $2n+2$  " "  $n$   $\rightarrow$  degree:  $(2n+2)-1 = 2n+1$

Using Lagrange:  $n+1$  nodes give  $n$  degree:  $P_n$

Using Hermite:  $n+1$  nodes give  $2n+1$  degree:  $P_{2n+1} \rightarrow$  higher degree polynomial  $\rightarrow$  Better accuracy

e.g given: 4 nodes

(i)  $n=3$

(ii) Degree of polynomial after using Lagrange/Natural:  $P_n = P_3$

(iii) Using Hermite Interpolation:  $P_{2n+1} = P_{2x_3+1} = P_7$

Hermite Interpolation Theorem: Given nodes  $(x_0, x_1, \dots, x_n)$  and  $f(x_0), f(x_1), \dots, f(x_n)$  and  $f'(x_0), f'(x_1), \dots, f'(x_n)$   $\rightarrow$  Then there exists an unique polynomial  $P_{2n+1}(x)$

Lagrange:  $P_n(x) = \sum_{k=0}^n f(x_k) l_k(x)$

$$l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^n \frac{(x-x_j)}{(x_k-x_j)} = \frac{(x-x_0)}{(x_k-x_0)} \cdot \frac{(x-x_1)}{(x_k-x_1)} \cdot \frac{(x-x_2)}{(x_k-x_2)} \cdots \frac{(x-x_n)}{(x_k-x_n)}$$

Hermite:  $P_{2n+1}(x) = \sum_{k=0}^n [f(x_k) h_k(x) + f'(x_k) \hat{h}_k(x)]$

e.g nodes  $(x_0, x_1)$

$$= n=1 \quad P_{2n+1} = P_{2x_1+1} = P_3 = [f(x_0) h_0(x) + f'(x_0) \hat{h}_0(x)] + [f(x_1) h_1(x) + f'(x_1) \hat{h}_1(x)]$$

$$h_k(x) = [1 - 2(x-x_k) l'_k(x)] \{l_k(x)\}^2 \quad \hat{h}_k(x) = (x-x_k) \{l_k(x)\}^2$$

e.g.  $f(x) = 2\pi x$  nodes  $\{0, \pi/2\}$

$$\begin{array}{lll} n=1 & f(0)=0 & f'(0)=2\pi \\ & f(\pi/2)=1 & f'(1)=0 \end{array}$$

$$\begin{aligned} P_{2x_1+1}(x) &= \left\{ 0 \times h_0(x) + 1 \times \hat{h}_0(x) \right\} + \left\{ 1 \times h_1(x) + 0 \times \hat{h}_1(x) \right\} \\ &= \left\{ (x-0) l_0^2(x) \right\} + \left\{ 1 - 2(x-\pi/2) l_1'(\pi/2) l_1^2(x) \right\} \\ &= x \left( \frac{1-2x}{\pi} \right)^2 + (1-2(x-\pi/2)) \cdot \frac{2}{\pi} \left( \frac{4x^2}{\pi^2} \right) \end{aligned}$$

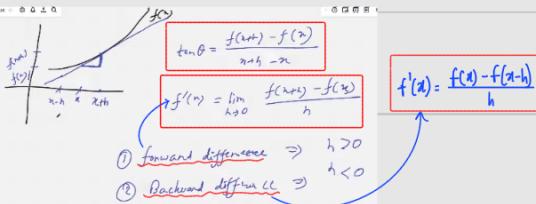
$$l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^n \frac{(x-x_j)}{(x_k-x_j)}$$

$$l_0(x) = \frac{x-\pi/2}{\pi/2-0} = \frac{\pi/2-x}{\pi/2} = 1 - \frac{2x}{\pi}$$

$$l_1(x) = \frac{x-0}{\pi/2-0} = \frac{x-0}{\pi/2} = \frac{2x}{\pi} \quad \left\{ l_0^2(x) = \left( 1 - \frac{2x}{\pi} \right)^2 \right.$$

$$l_1'(x) = \frac{2}{\pi} \quad \downarrow \quad \left. l_1^2(x) = \frac{4x^2}{\pi^2} \right\}$$

## Differentiation - 1



$$f'(x) = P(x) + \text{Error}$$

$$f'(x) = f(x_0) \left( \frac{x_1 - x_0}{x_0 - x_1} \right) + f(x_1) \left( \frac{x_2 - x_0}{x_1 - x_0} \right) + \frac{f''(\xi)}{2!} (x-x_0)(x-x_1)$$

Langrange

Cauchy

$$f'(x) = \frac{f(x_0)}{x_0 - x_1} + \frac{f(x_1)}{x_1 - x_0} + \frac{f''(\xi)}{2} \frac{d}{dx} (\xi) (x-x_0)(x-x_1) + \frac{f''(\xi)}{2} (2x-x_0-x_1)$$

$$f'(x_0) = \frac{f(x_0)}{x_0 - x_1} + \frac{f(x_1)}{x_1 - x_0} + \frac{f''(\xi)}{2} (x_0 - x_1) \quad [x_0 = x_0]$$

$$f'(x_0) = \frac{f(x_0)}{h} + \frac{f(x_0+h)}{h} - \frac{f''(\xi)}{2} (h) \quad [x_1 = x_0 + h]$$

Truncation error

in case of forward differentiation  
Same for backward diff.

$$f'(x_0) = \frac{f(x_0+h) - f(x_0)}{h} - \frac{f''(\xi)}{2} (h)$$

Q) Find forward difference and Truncation error.

$$f(x) = \ln(x), x_0=2$$

$$h=1, 0.1, 0.01, 0.001$$

$$\text{Ans: } f'(x) = \frac{1}{x}$$

$$f'(2) = \frac{1}{2} = 0.5$$

If  $h$  values are not given,  
choose it yourself! small value

$h$	Forward difference	Truncation Error
1	0.49465	$0.09535$
0.1	0.48791	$0.01203$
0.01	- - -	- - -
0.001	- - -	- - -

$$f(x+h) - f(x)$$

$$h$$

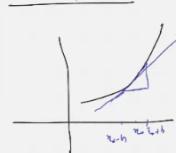
forward

In case of backward difference, curv will be -ve

$$= \frac{\ln(2+1) - \ln(2)}{1} = 0.49465$$

$f'(x) -$

Central difference



$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x-h)}{2h}$$

Q) Find central difference and Truncation error.

$$f(x) = \ln(x), x_0=2$$

$$h=1, 0.1, 0.01, 0.001$$

$$= f'(x) = \frac{1}{x}$$

$$f'(2) = \frac{1}{2} = 0.5$$

$h$	Central difference	Truncation error
1	0.59306	$-0.019306$
0.1	1	↑
0.01	1	↑
0.001	1	↓

$$f(x) = f(x_0) + \frac{(x-x_0)(x-x_1)}{2!} + f(x_0) + \frac{(x-x_0)(x-x_1)(x-x_2)}{3!} + \dots$$

$$f'(x_0) = \frac{f(x_0+h) - f(x_0-h)}{2h} - \frac{f'''(\xi)}{6} (h^3)$$

$$\text{or } f'(x_0) = \frac{f(x_0+h) - f(x_0-h)}{2h} + \frac{f'''(\xi)}{6} (h^3) \quad \text{Error}^{-1} \star$$

$$f'(x_0) = \frac{f(x_0+h) - f(x_0-h)}{2h}$$

error

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

$$= \frac{\ln(2+1) - \ln(2-1)}{2 \cdot 1}$$

## Differentiation - 2

In  $f(x+h), f(x-h)$ ,  $h/h^2$  rounding error occurs

$$\text{Error} = \text{Truncation error} + \text{Rounding error}$$

$$= -\frac{f'''(\xi)}{6} h^2 + \text{Error} \frac{f(x+h) - f(x-h)}{2h}$$

machine epsilon

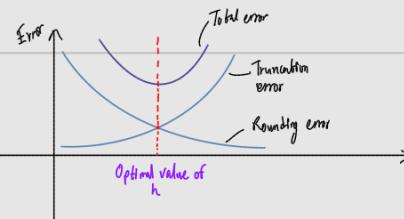
$$\text{Central: } \frac{f'''(\xi)}{6} h^2 \Rightarrow \text{Error is of order } h^2 = \Theta(h^2)$$

$$\text{forward: } \frac{f''(\xi)}{2} h \Rightarrow \Theta(h)$$

$$\text{Error} \leq \left| \frac{f'''(\xi)}{6} h^2 \right| + \text{Error} \frac{|f(x+h) - f(x-h)|}{2h}$$

Error  $\frac{h}{h^2}$

Error  $\frac{1}{h}$



Richardson Extrapolation works for central difference Only

Derivation may come ★

$$\text{forward: } \frac{f(x+h) - f(x)}{h} - \frac{f''(\xi)}{2} h \quad \xrightarrow{\text{Ans}} D_h = \frac{f(x+h) - f(x)}{h} - \frac{f''(\xi)}{2} h$$

further approximation

$$f(x+h) = f(x) + f'(x)h + \frac{f''(x)}{2!} h^2 + \frac{f'''(x)}{3!} h^3 + \dots$$

$$f(x-h) = f(x) - f'(x)h + \frac{f''(x)}{2!} h^2 - \frac{f'''(x)}{3!} h^3 + \dots$$

$$D_h = f'(x) + \frac{f''(x)}{2} h + \frac{f'''(x)}{3!} h^2 + O(h^3) - O(h^4)$$

$$D_{h/2} = f'(x) + \frac{f''(x)}{2} \frac{h}{2} + \frac{f'''(x)}{3!} \frac{h^2}{4} + O(h^4) - O(h^5)$$

$$4D_{h/2} - D_h = 3f'(x) + O(h^4) + \frac{f'''(x)}{5!} h^4 + O(h^4)$$

$$\frac{4D_{h/2} - D_h}{3} = f'(x) - \frac{1}{4} \frac{f'''(x)}{3!} h^4 + O(h^4) = D_h^{(1)}$$

$$D_h^{(1)} = \frac{2^4 D_{h/2} - D_h}{2^4 - 1}$$

$$D_h^{(2)} = f'(x) - \frac{1}{4} \frac{f''(x)}{2!} h^2 + O(h^4) - O(h^5)$$

$$D_{h/2}^{(2)} = f'(x) - \frac{1}{4} \frac{f''(x)}{2!} \frac{h}{2} + O(h^4) - O(h^5)$$

$$16 D_{h/2}^{(2)} - D_h^{(2)} = 15 f'(x) + O(h^4) = D_h^{(2)}$$

We could have applied this Richardson extrapolation procedure without knowing the coefficients of the error series. If we have some general order- $n$  approximation

$$D_h^{(1)} = f'(x) + C(h^n) + O(h^{n+1})$$

then we can always evaluate it with  $\frac{h}{2}$  to get



$$D_{\frac{h}{2}} = f'(x) + C\left(\frac{h}{2}\right)^n + O\left(\left(\frac{h}{2}\right)^{n+1}\right)$$

and then eliminate the  $h^n$  term to get a new approximation

$$D_h^{(2)} = \frac{2^n D_{\frac{h}{2}} - D_h}{2^n - 1} = f'(x) + O(h^{n+1})$$

Solving non-Linear eq<sup>n</sup>

$$\text{eg } x^3 + 2x^2 + 3x + 5 = 0$$

$$\therefore x = \boxed{\quad}, \boxed{\quad}, \boxed{\quad}$$

Finding root

$$\text{eg } f(x) = x^2 - 4$$

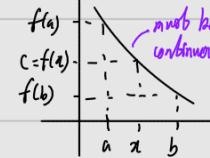
$$f(x) = 0$$

$$\therefore x = \pm 2 \text{ (roots)}$$

$$\text{eg } f(x) = \frac{1}{x} - x$$

$$f(x) = 0$$

$$x = \frac{1}{x}$$



$$c \in [f(a), f(b)]$$

$$x \in [a, b]$$

\*\*\* This is how we find the roots but Computer uses different Algorithms.

1] Internal Bisection/Bisection Method: Based on Intermediate Value Theorem

(4.1: Out of Syllabus but imp for Lab)  $\hookrightarrow$  finding roots of non-linear eq<sup>n</sup>

$$\text{Q1) } f(x) = \frac{1}{x} - 0.5, a_0 = 1, b_0 = 3$$

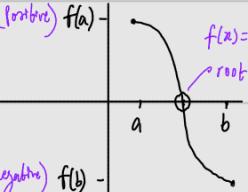
$$f(a_0) = \frac{1}{1} - 0.5 = 0.5, f(b_0) = -0.17$$

$f(a_0) \cdot f(b_0) < 0 \Rightarrow$  At least one root exists  
if  $f(a_0) \cdot f(b_0) > 0 \Rightarrow$  No root exists

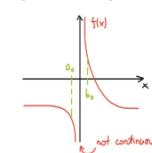
$$\text{midpoint, } m_0 = \frac{a_0 + b_0}{2} = \frac{1+3}{2} = 2 \text{ (positive) } f(a)$$

$$f(m_0) = f(2) = \frac{1}{2} - 0.5 = 0$$

$$m_0 = x_* = 2 \Rightarrow \text{root}$$



Example  $\rightarrow$  Same example with initial interval [-0.5, 0.5].



In this case  $f(a_0) \cdot f(b_0) < 0$ , but there is no root in the interval.

1st Iteration  $\rightarrow f(x) = \frac{1}{x} - 0.5, a_0 = 1, b_0 = 3$

$$f(a_0) = 0.5, f(b_0) = -0.17$$

$\therefore f(a_0) \cdot f(b_0) < 0 \Rightarrow$  root b exists

$$\# m_0 = \frac{1.5+3}{2} = 2.25 \# f(m_0) = \frac{1}{2.25} - 0.5 = -0.056$$

$$= 2.25$$

$\Rightarrow$  b  $\rightarrow$  re, update b

2nd Iteration  $\rightarrow a_1 = 1.5, b_1 = 2.25, f(a_1) \cdot f(b_1) < 0$

$$f(a_1) = 0.17$$

$$f(b_1) = -0.056$$

$$m_1 = \frac{a_1 + b_1}{2} = 1.875$$

$$\rightarrow f(m_1) = 0.033 \rightarrow \text{re, update } a$$

no. of iterations will be given in the question

3rd Iteration  $\rightarrow a_2 = m_1 = 1.875, b_2 = 2.25, f(m_2) = -0.01525 \rightarrow \text{re, update } b_2$

$$m_2 = \frac{a_2 + b_2}{2} = 2.0625$$

must check  $f(a_2) \cdot f(b_2) < 0$

Taking  $m_0$  close to zero

Convergence  
↓  
Slow

The rate of convergence is steady, so we can pre-determine how many iterations will be needed to converge to a given accuracy. After K iterations, the interval has length

$$|b_k - a_k| = \frac{|b_0 - a_0|}{2^k}$$

so the error in the mid-point satisfies

$$|m_k - x_*| \leq \frac{|b_0 - a_0|}{2^{k+1}}$$

In order for  $|m_k - x_*| \leq \delta$  we need n iterations, where

$$\frac{|b_0 - a_0|}{2^{k+1}} \leq \delta \implies \log(\frac{|b_0 - a_0|}{\delta}) - (n+1) \log(2) \leq \log(\delta) \implies n \geq \frac{\log(|b_0 - a_0|) - \log(\delta)}{\log(2)} - 1$$

Example  $\rightarrow$  With  $a_0 = 1.5, b_0 = 3$ , as in the above example, then for  $\delta = \epsilon_M = 1.1 \times 10^{-10}$  we would need

$$n \geq \frac{\log(1.5) - \log(1.1 \times 10^{-10})}{\log(2)} - 1 \implies n \geq 53 \text{ iterations.}$$

Advantage: Method is robust

Disadvantage: The convergence is pretty slow.  
Only works in One dimension

a	b	m	f(a)	f(b)	f(m)
K=0					
K=1					
K=2					
:					

$$a_0 = 1.5, b_0 = 3, b_1 = 2.25, a_1 = 1.5, \therefore |b_1 - a_1| = \frac{|b_0 - a_0|}{2} = \frac{1.5}{2} = 0.75$$

$$|b_0 - a_0| = 1.5$$

$$|b_1 - a_1| = 0.75$$

$$\# |b_2 - a_2| = \frac{|b_1 - a_1|}{2} = \frac{0.75}{2} = 0.375 \therefore |b_k - a_k| = \frac{|b_0 - a_0|}{2^k}, k = \text{no. of iteration} + 1$$

# As  $m_k$  is never equal to  $x_*$ ,

$|m_k - x_*| = \text{Error}$   
actual error

$$a_k \quad m_k \quad x_* \quad b_k$$

$$\begin{aligned} \therefore |m_k - x_*| &\leq \frac{|b_k - a_k|}{2} \\ \Rightarrow |m_k - x_*| &\leq \frac{|b_0 - a_0|}{2^k} \\ \therefore |m_k - x_*| &\leq \frac{|b_0 - a_0|}{2^{k+1}} \end{aligned}$$

Actual error      max<sup>m</sup> possible error

$$\begin{aligned} \frac{|b_0 - a_0|}{2^{k+1}} &\leq \delta \\ \Rightarrow \ln \left| \frac{|b_0 - a_0|}{2^{k+1}} \right| &\leq \ln \delta \\ \Rightarrow \ln |b_0 - a_0| - \ln 2^{k+1} &\leq \ln \delta \\ \Rightarrow -\ln 2^{k+1} &\leq \ln \delta - \ln |b_0 - a_0| \\ \Rightarrow \ln 2^{k+1} &\geq \ln |b_0 - a_0| - \ln \delta \end{aligned}$$

$$(k+1) \ln 2 \geq \ln |b_0 - a_0| - \ln \delta$$

$$\therefore k \geq \frac{\ln |b_0 - a_0| - \ln \delta}{\ln 2} - 1$$

$$\text{eg } \int = 0.00001, a_0 = 1.5, b_0 = 3$$

$$K \geq 16.19$$

$$K \geq \frac{\ln |3-1.5| - \ln (0.00001)}{\ln 2} - 1 \therefore K = 17 \text{ (18 iterations)}$$

$$K \geq 52.6$$

$$\therefore K = 53 //$$

# Bisection Method Disadvantages

- Convergence is slow
- Multiple roots can't be found

Hence

Fixed point iteration

- Convergence is faster
- Can find multiple roots

$f(x) = 0$   
 $g(x) = x \Rightarrow$  fixed point ( $x_g$ )

e.g.  $f(x) = x^2 - 2x - 3$   
using middle term,  
 $x_* = 3, x_* = -1$

$f(x) = x$  will be given

$$\begin{aligned} x^2 - 2x - 3 &= 0 \\ x^2 &= 2x + 3 \end{aligned}$$

$x = \sqrt{2x+3} = g(x)$ , given  $x_0 = 0$

$$\begin{aligned} g(0) &= \sqrt{3} = 1.7321 & x_0 = 0 \\ g(\sqrt{3}) &= 2.5425 & x_1 = \sqrt{3} \\ g(2.5425) &= 2.8939 & x_2 = 2.8939 \\ g(2.8939) &= 2.94 & x_3 = 2.94 \end{aligned}$$

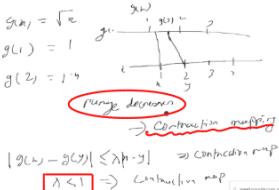
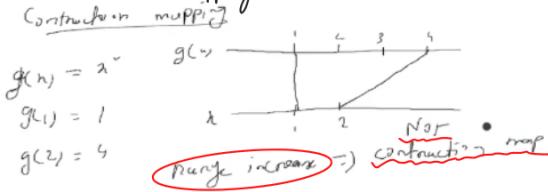
initial value

4 iterations

$$\begin{aligned} g(x) &= \frac{3}{x-2} \Rightarrow 33 \text{ iterations} \\ g(x) &= \frac{x^2-3}{2x} \Rightarrow 5 \text{ iterations} \\ g(x) &= \frac{x^2-3}{2x-2} \Rightarrow 20 \text{ iterations} \end{aligned}$$

Computer chooses the best one; using Contraction Method

### Contraction Mapping



$\lambda \propto \frac{1}{\text{Convergence Rate}}$  ★

Contraction Mapping is when  $\lambda = 1$  ★

Bisection Method:  $\lambda = \frac{1}{2}$   
 $|g(x_2) - g(x_1)| \leq \lambda |x_2 - x_1|$   
 $|g(x) - g(y)| \leq |g'(x)| |x - y|$

eg.  $g(x) = \sqrt{2x+3}$   
 $g'(x) = \frac{1}{2}(2x+3)^{-1/2} \cdot 2$   
 $= \frac{1}{\sqrt{2x+3}} \quad \lambda < 1$   
 $\therefore$  Contraction mapping of function  
 $g(M) \text{ is possible if } x > -1$

eg.  $f(x) = \sqrt{2x+3}, x = \frac{1}{2}, y = 4$   
 $g(-\frac{1}{2}) = \sqrt{2}$   
 $g(4) = g(4) = \sqrt{11}$   
 $\therefore |g(x) - g(y)| \leq \lambda |x - y|$   
 $\therefore 3.32 < 9.5$   
 $\therefore$  Contraction mapping possible

★ if  $x, y$  given but not  $\lambda$  value, then contraction mapping is possible if  
 $|g(x) - g(y)| < |x - y|$

Q)  $g(x) = \sqrt{2x+3}, x_0 = 0, x_* \rightarrow 3$

$x_k$	$ x_* - x_k $	$ x_k - x_{k-1}  /  x_k - x_{k-1} $
$x_0 \rightarrow 0$	3	Can't calculate
$x_1 \rightarrow \sqrt{3}$	1.2675 ( $3 - \sqrt{3}$ )	0.4226
$x_2 g(\sqrt{3})$	0.4575	0.3608
$x_3$		0.3423
$x_4$		0.3365
		0.3339

going towards  $0.333 = \frac{1}{3}$   
 $\therefore \lambda = \frac{1}{3}$

Q)  $g(x) = \frac{x+5}{2x-2}, x_0 = 0$

0.5  
0.1  
0.012  
0.00012  
1  
J = 0

high rate of convergence

★ if  $\lambda = 0 \Rightarrow$  Super Linear Convergence  
 $\lambda = \frac{1}{3} / \frac{1}{2} / \dots \Rightarrow$  Linear Convergence

★ if  $\lambda > 1 \Rightarrow$  Divergence. Hence, we can't find root

### Newton's Method & Aitken Acceleration 1

So far, (1)  $f(x) = \sqrt{2x+3}, x_* = 3$

$$g'(x) = \frac{1}{\sqrt{2x+3}}$$

$$g'(x_*) = g'(3) = \frac{1}{\sqrt{2x+3}} = \frac{1}{\sqrt{2 \cdot 3 + 3}} = \frac{1}{3} = \lambda \quad g'(-) = 0 = \lambda \quad (\text{Super Linear Convergence})$$

(2)  $f(x) = \frac{x^2+3}{2x-2}, x_* = -1$

$$g'(x) = \frac{x^2-2x-3}{2(x-1)^2}$$

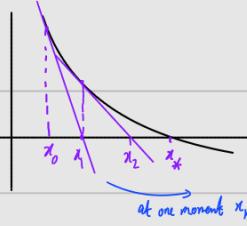
### Disadvantage of Bisection Method:

- Convergence rate is  $\frac{1}{2}$  which is very low

Thus we used fixed point iteration to reduce value of  $\lambda$  (less than  $\frac{1}{2}$ )  
but not always zero.

\* \* \* Newton's Method convergence rate is always zero

Newton's Method (4.4) : This is a particular fixed point iteration that is very widely used because it usually converges Superlinearly.



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

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

$$x_3 = x_2 - \frac{f(x_2)}{f'(x_2)}$$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

$$f'(x) = \frac{f(x) \times f''(x)}{(f'(x))^2}, f(x_*) = 0$$

$\lambda$  tends to zero

$$f(x) = x - \frac{f(x)}{f'(x)}$$

$$f'(x) = \frac{f(x) \times f''(x)}{(f'(x))^2}, f(x_*) = 0$$

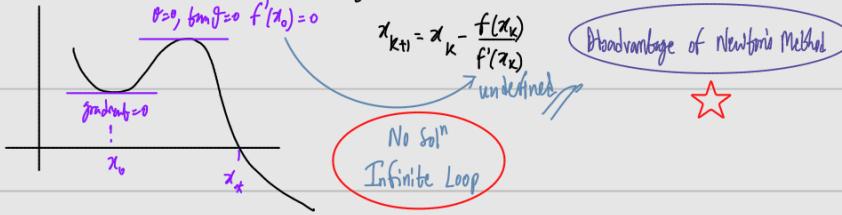
$$f'(x_*) = \frac{f(x_*) \times f''(x_*)}{(f'(x_*)^2)} = 0$$

$$f'(x_*) = 0 \rightarrow \lambda //$$

$$\text{Derivation: } m = f'(x_k) = \frac{0 - f(x_k)}{x_{k+1} - x_k}$$

$$\therefore x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

# Sometimes, Newton's Method doesn't give a soln



In fact, Newton only applied the method to polynomial equations, and without using calculus. The general form using derivatives ("fluxions") was first published by Thomas Simpson in 1740. [See "Historical Development of the Newton-Raphson Method" by T.J. Ypma, SIAM Review 37, 531 (1995).]

$$\text{eg } f(x) = \frac{1}{x} - a, x_* = \frac{1}{a}, a > 0$$

$$\begin{aligned} x_{k+1} &= x_k - \frac{f(x_k)}{f'(x_k)} \\ &= x_k - \frac{\left(\frac{1}{x_k} - a\right)}{-\frac{1}{x_k^2}} \\ &\vdots \\ &= 2x_k - ax_k^2 \end{aligned}$$

$x_k$	$ x_* - x_k $	$ x_* - x_k  /  x_k - x_{k-1}  = \lambda$
$x_0 = 1$	1	—
$x_1 = 0.5$	0.5	0.5
$x_2 = 0.125$	0.125	0.125
$x_3 = 0.09765625$	0.00781	0.0625
$x_4 = 0.099997$	0.00003	0.0003
$x_5 = 0$	0	0 $\rightarrow \lambda$

$$\text{eg } f(x) = x^3 - 2x + 2, x_* \approx -1.7693, x_0 = 0$$

$$f'(x) = 3x^2 - 2$$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} = \frac{x_k^3 - 2x_k + 2}{3(x_k)^2 - 2}$$

(4.5: Out of syllabus)

Solving two Equations

$$f_1(x, y) = xy - y^3 - 1 = 0 \quad x_0 = 2,$$

$$f_2(x, y) = x^2y + y - 5 = 0 \quad y_0 = 3$$

$$\frac{\partial f_1}{\partial x} = y = 3, \quad \frac{\partial f_1}{\partial y} = x - 3y^2 = 2 - 3(3)^2 = -25$$

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} = \begin{pmatrix} x_k \\ y_k \end{pmatrix} - \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{pmatrix}^{-1} \begin{pmatrix} f_1(x_k, y_k) \\ f_2(x_k, y_k) \end{pmatrix}$$

↓ Jacobian Matrix

$$\frac{\partial f_2}{\partial x} = 2y = 2 \cdot 3 = 6, \quad \frac{\partial f_2}{\partial y} = x^2 + 1 = 2^2 + 1 = 5$$

$$\therefore \begin{pmatrix} 3 & -25 \\ 12 & 5 \end{pmatrix}^{-1} = \begin{pmatrix} 1/3 & 5/63 \\ -4/105 & 1/105 \end{pmatrix}$$

$$\therefore \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 1/3 & 5/63 \\ -4/105 & 1/105 \end{pmatrix} \begin{pmatrix} -22 \\ 10 \end{pmatrix} \rightarrow \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1.54 \\ 1.47 \end{pmatrix}$$

$$= \begin{pmatrix} 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 1/3 & 5/63 \\ -4/105 & 1/105 \end{pmatrix} \begin{pmatrix} 14/3 \\ 14/15 \end{pmatrix} \rightarrow \begin{pmatrix} x_3 \\ y_3 \end{pmatrix} = \begin{pmatrix} 1.48 \\ 1.47 \end{pmatrix}$$

$$\therefore \begin{pmatrix} x_4 \\ y_4 \end{pmatrix} = \begin{pmatrix} 1.47 \\ 1.47 \end{pmatrix} \rightarrow \begin{pmatrix} x_5 \\ y_5 \end{pmatrix} = \begin{pmatrix} 1.499 \\ 1.499 \end{pmatrix} \rightarrow \text{Approaches } \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

Q)  $f(x) = x^2 - 2x e^{-x} + e^{-2x}$ ,  $x_0 = 1$  with in  $10^{-5}$   $|f(x_k)| < 10^{-5}$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

$$x_{k+1} = x_k - \frac{x_k^2 - 2x_k e^{-x_k} + e^{-2x_k}}{2x_k - 2e^{-x_k} + 2x_k e^{-x_k} + 2e^{-2x_k}}$$

$x_k$	$f(x_k)$	is $ f(x_k)  < 10^{-5}$ ?
$x_0 = 1$	0.399576	No
$x_1 = 0.768591$	0.073222	No
$x_2 = 0.664990$	0.022532	No
$x_3 = 0.56625$	$0.5 \times 10^{-5}$	Yes ✓

Use Calculator Tricks

\* \* \* \* \* Sometimes convergence rate is very low, hence it needs to be increased

### Aitken Acceleration → Increases convergence Rate

(4.6: Out of syllabus)

e.g.  $f(x) = \frac{1}{x} - 0.5$ ,  $x_0 = 2$

$$g(x) = x + \frac{1}{16} \left( \frac{1}{x} - 0.5 \right)$$

f(x)

$$g'(x) = 1 + \frac{1}{16} \left( -\frac{1}{x^2} \right)$$

$$\lambda = g'(x_0) = g'(2) = 0.989375 \Rightarrow$$

e.g.  $f(x) = \frac{1}{x} - 0.5$ ,  $x_0 = 1.5$

$$g(x) = x + \frac{1}{16} \left( \frac{1}{x} - 0.5 \right)$$

$$\hat{x}_{k+2} = x_k - \frac{(x_{k+1} - x_k)^2}{x_{k+2} - 2x_{k+1} + x_k}$$

After two iterations  
needs to find the

$$x_1 = g(x_0) = 1.510417$$

$$x_2 = g(x_1) = 1.521596$$

$$\hat{x}_2 = x_0 - \frac{(x_1 - x_0)^2}{x_2 - 2x_1 + x_0} = 1.877604$$

$$x_3 = g(\hat{x}_2) = 1.879691$$

$$x_4 = g(x_3) = 1.881642$$

$$\hat{x}_4 = x_2 - \frac{(x_3 - x_2)^2}{x_4 - 2x_3 + \hat{x}_2} = 1.992634$$

Using fixed point iteration, it would have taken 819 iterations,  
here it only took eight

Advantage of Newton's Method

\* \* \* \* \* Newton's Method has high computation cost. Thus to reduce cost,  $f'(x)$  is replaced by an approximate function  $g(x)$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

$f'(x) \approx g(x)$  : Quasi-Newton's Method (4.7)

↳ Secant method

↳ D'Alembert's method

(out of syllabus)

#  $x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$ , Backward Difference =  $\frac{f(x_k) - f(x_{k-h})}{h}$

$$g_k = \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$$

$$\Rightarrow x_{k+1} = x_k - \frac{f(x_k)}{\partial_k} = x_k - \frac{f(x_k) [x_k - x_{k-1}]}{f(x_k) - f(x_{k-1})}$$

No  $f'(x_k)$  hence less cost

★ Secant Method requires two starting points ★

e.g.  $f(x) = \frac{1}{x} - 0.5$ ,  $x_k = 2$ ,  $x_0 = 0.125$ ,  $x_1 = 0.5$

$$x_2 = x_1 - \frac{f(x_1)(x_1 - x_0)}{f(x_1) - f(x_0)} = 0.6875$$

$$x_3 = x_2 - \frac{f(x_2)(x_2 - x_1)}{f(x_2) - f(x_1)}$$

$$x_4 = 1.99916$$

$$x_5 = 2$$

$$\frac{|x_k - x_{k-1}|}{|x_k - x_{k-1}|} = 0.75$$

$$= 0.65625$$

$$= 0.15194$$

$$\rightarrow 0 \quad \text{Super Linear Convergence}$$

Use Calculator Tricks

Using Newton's Method,  $\lambda$  is always zero  
As Secant Method is part of / improved version  
of Newton's Method, hence here  $\lambda$  is also zero  
which is Super Linear Convergence

# Gaussian Elimination method and LU decomposition - 1

23/04/2025

Wednesday

$$\begin{aligned} x_1 - 2x_2 + x_3 &= 0 \\ x_1 - 2x_2 + 2x_3 &= 4 \\ 2x_1 + 2x_2 - 2x_3 &= 1 \end{aligned}$$

$$\begin{pmatrix} 1 & 2 & 1 \\ 1 & -2 & 2 \\ 2 & 2 & -2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 4 \\ 1 \end{pmatrix}$$

$$A \cdot x = b$$

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

Don't mix up with Vandermonde Matrix  
That solves Polynomial and now we're solving Linear system. ★

# Unique soln exists if:

- (1) A square matrix
- (2)  $\det(A) \neq 0$  [non singular matrix]

finding inverse matrix is of computationally high cost  
hence computer uses a different approach to solve the linear system by using:

In general, any lower triangular system  $Lx = b$  can be solved by forward substitution

$$x_j = \frac{b_j - \sum_{k=1}^{j-1} l_{jk} x_k}{l_{jj}}, \quad j = 1, \dots, n. \quad (5.5)$$

Similarly, an upper triangular matrix  $U$  has the form

$$U = \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & u_{nn} \end{pmatrix}, \quad (5.6)$$

and an upper-triangular system  $Ux = b$  may be solved by backward substitution

$$x_j = \frac{b_j - \sum_{k=j+1}^n u_{kj} x_k}{u_{jj}}, \quad j = n, \dots, 1. \quad (5.7)$$

To estimate the computational cost of forward substitution, we can count the number of floating-point operations (+, -, ×, ÷).

Example → Number of operations required for forward substitution.

Consider each  $x_j$ . We have

$$\begin{aligned} j = 1 &\rightarrow 1 \text{ division} \\ j = 2 &\rightarrow 1 \text{ division} + [1 \text{ subtraction} + 1 \text{ multiplication}] \\ j = 3 &\rightarrow 1 \text{ division} + 2 \times [1 \text{ subtraction} + 1 \text{ multiplication}] \\ &\vdots \\ j = n &\rightarrow 1 \text{ division} + (n-1) \times [1 \text{ subtraction} + 1 \text{ multiplication}] \end{aligned}$$

So the total number of operations required is

$$\sum_{j=1}^n (1 + 2(j-1)) = 2 \sum_{j=1}^n j - \sum_{j=1}^n 1 = n(n+1) - n = n^2.$$

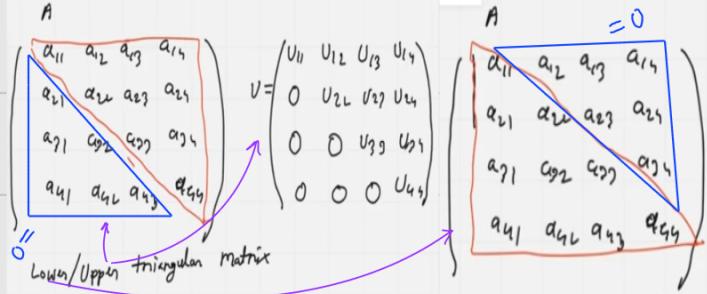
Upper triangular matrix → Bottom up approach

Lower triangular matrix → Top down approach

Computational Complexity of the algorithm

$A = n \times n$  Lower/Upper triangular matrix transformation =  $O(n^2)$

## (1) Gaussian Elimination: (5.1, 5.2)



### Augmented Matrix

$$\left( \begin{array}{ccc|c} 1 & 2 & 1 & 0 \\ 0 & -2 & 2 & 4 \\ 0 & 1 & -2 & 4 \\ 0 & 0 & 1 & 4 \end{array} \right)$$

Row multipliers:  $m_{21} = \frac{a_{21}}{a_{11}} = \frac{-2}{1} = -2$ ,  $m_{31} = \frac{a_{31}}{a_{11}} = \frac{1}{1} = 1$ ,  $m_{41} = \frac{a_{41}}{a_{11}} = \frac{0}{1} = 0$

$R_2 \rightarrow R_2 - R_1 \times m_{21}$ ,  $R_3 \rightarrow R_3 - R_1 \times m_{31}$ ,  $R_4 \rightarrow R_4 - R_1 \times m_{41}$

Upper triangular matrix:  $\left( \begin{array}{ccc|c} 1 & 2 & 1 & 0 \\ 0 & -4 & 1 & 4 \\ 0 & 0 & -1 & 4 \\ 0 & 0 & 0 & 4 \end{array} \right)$

Bottom up approach:

$$\begin{pmatrix} 1 & 2 & 1 \\ 0 & -4 & 1 \\ 0 & 0 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 4 \\ 12 \end{pmatrix}$$

$x_3 = 1$ ,  $x_2 = -4$ ,  $x_1 = 0$

## LU Decomposition

$$A = LUV$$

$$Ax = b$$

$$L\bar{U}x = b, \quad Ux = y$$

$$j=?$$

Unit b: Lower triangular matrix

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$\det(L) = \det(V)$

$L$ : Lower,  $V$ : Upper of  $A$

Both  $L$  &  $V$  of  $A$  give same soln

$$A^{(3)} = V$$

## Gaussian Elimination method and LU decomposition - 2

$$(5.3) \quad \begin{pmatrix} 1 & 2 & 1 \\ 1 & -2 & 2 \\ 2 & 2 & -2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 4 \\ 1 \end{pmatrix}$$

$$A^{(1)} \quad x \quad b$$

$$M_{21} = \frac{a_{21}}{a_{11}} = \frac{-2}{1} = -2$$

$$M_{31} = \frac{a_{31}}{a_{11}} = \frac{2}{1} = 2$$

$$\text{Frobenius Matrix, } F^{(1)} = \begin{pmatrix} 1 & 0 & 0 \\ -m_{21} & 1 & 0 \\ -m_{31} & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -2 & 0 & 1 \end{pmatrix}$$

$$A^{(2)} = F^{(1)} \times A^{(1)} = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -2 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 1 \\ 1 & -2 & 2 \\ 2 & 2 & -2 \end{pmatrix} = \begin{pmatrix} 1 & 2 & 1 \\ 0 & -4 & 1 \\ 0 & 0 & -3 \end{pmatrix}$$

$$A^{(2)} = \begin{pmatrix} 1 & 2 & 1 \\ 0 & -4 & 1 \\ 0 & 0 & -3 \end{pmatrix} \quad F^{(2)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad M_{32} = \frac{a_{32}}{a_{22}} = \frac{1}{-4} = -\frac{1}{4}$$

$$M_{32} = \frac{a_{32}}{a_{22}} = \frac{1}{-4} = -\frac{1}{4}$$

$$\therefore F^{(2)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$V = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 1 \\ 0 & -4 & 1 \\ 0 & 0 & -3 \end{pmatrix} = \begin{pmatrix} 1 & 2 & 1 \\ 0 & -4 & 1 \\ 0 & 0 & -3 \end{pmatrix} = V$$

$$L = \begin{pmatrix} 1 & 0 & 0 \\ m_{21} & 1 & 0 \\ m_{31} & m_{32} & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 0 & -\frac{1}{4} & 1 \end{pmatrix}$$

$$\text{eg } \begin{pmatrix} 0 & 3 & 0 \\ 2 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix} \rightarrow \text{Home swap the row}$$

$$M_{21} = \frac{a_{21}}{a_{11}} = \frac{2}{0} \text{ (undefined)}$$

$$A^{(1)} = L \times V \quad Ax = b$$

$$L \bar{U}x = b, \quad Vx = y$$

$$Ly = b \quad Ly = b$$

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & -2 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} \quad \# x_1 = 0, \quad x_2 = 1, \quad x_3 = 1$$

$$\therefore \begin{pmatrix} 1 & 2 & 1 \\ 0 & -4 & 1 \\ 0 & 0 & -3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

$$x_1 = 1, \quad x_2 = -\frac{1}{4}, \quad x_3 = -\frac{1}{3}$$

Substitution Process  
 $\rightarrow O(n^2)$

### LU decomposition

LU decomposition offers several advantages over Gaussian elimination, particularly when solving multiple systems of linear equations with the same coefficient matrix.

If you're solving the same system for different  $b$  vectors, you don't need to redo the decomposition each time! **I just reuse L and U!** In contrast, in the Gaussian elimination method, if  $b$  changes, we need to restart row operations from the beginning.

→ So far we used 3 unknowns, 3 equations.  $3 \times 3$  Matrix

→ But Gram-Schmidt elimination won't work when no. of unknowns  $\neq$  no. of equations. This is called Overdetermined System

$$\vec{x} = (2, 3, 5) \quad \vec{a} = \begin{bmatrix} 2 \\ 3 \\ 5 \end{bmatrix}, \quad \vec{a}^T = [2 \ 3 \ 5] \quad \text{magnitude/length/norm, } |\vec{a}| = \sqrt{2^2 + 3^2 + 5^2} = \sqrt{38}$$

$$|\vec{x}| = \sqrt{\vec{x}^T \vec{x}} = \sqrt{38} \Rightarrow L_2\text{-norm}$$

Transpose

Same

$$|\vec{x}| = \sqrt{\vec{x} \cdot \vec{x}} = \sqrt{(2, 3, 5) \cdot (2, 3, 5)} = \sqrt{4+9+25} = \sqrt{38}$$

Dot product/inner product

$$\vec{x} \cdot \vec{j} = |\vec{x}| \cdot |\vec{j}| \cos \theta = \vec{x}^T \vec{j}$$

$$\vec{x} = (2, 3, 5)^T \Rightarrow L_1\text{-norm}$$

$$= |2| + |3| + |5|$$

$$= 10$$

Orthogonality (6.1)

$$\vec{x}_1 \cdot \vec{x}_2 = \vec{x}_1 \cdot \vec{j} = 0, \quad \vec{x}_1^T \vec{j} = 0$$

Orthonormal

$$\textcircled{1} \quad \vec{x}_1^T \vec{j} = 0 \Rightarrow \vec{x}_1 \cdot \vec{x}_2 = 0$$

$$\textcircled{2} \quad \vec{x}_1^T \vec{x}_1 = 1, \quad \vec{j}^T \vec{j} = 1 \Rightarrow \text{Norm Unity}$$

e.g.  $S = \left\{ \frac{1}{\sqrt{5}}(2, 1)^T, \frac{1}{\sqrt{5}}(1, -2)^T \right\}$

$$\vec{u} = \frac{1}{\sqrt{5}}(2, 1)^T \quad \vec{v} = \frac{1}{\sqrt{5}}(1, -2)^T$$

$$\vec{u}^T \vec{u} = \vec{u} \cdot \vec{u} = \frac{1}{5}(4+1) = 1$$

$$\vec{v}^T \vec{v} = \vec{v} \cdot \vec{v} = \frac{1}{5}(1+(-2)^2) = 1$$

$$\vec{u}^T \vec{v} = \vec{u} \cdot \vec{v} = \frac{1}{5}(2 \cdot 1 - 1 \cdot 2) = 0$$

$$S = \left\{ \vec{x}_i \mid \vec{x}_i \in \mathbb{R}^n, \vec{x}_i^T \vec{x}_j = \delta_{ij} \right\}$$

$$\delta_{ij} = \begin{cases} 1, & i=j \rightarrow \text{normality} \\ 0, & i \neq j \rightarrow \text{orthogonality} \end{cases}$$

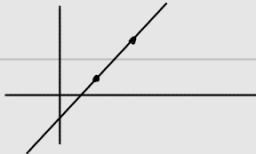
Kronecker Delta

Theorem:  $Q^{-1} = Q^T$ 

$$QQ^T = Q^T Q = I$$

Kronecker delta

#  $P_1(x)$  Nodes = 2  
 $P_1(x) = a_0 + a_1 x$



# Now for nodes = 3  
And we want  $P_1(x)$

This is not possible, as  
there are two straight lines

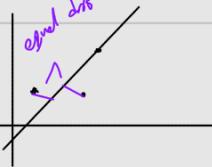
can't draw 1 degree  
st line passing through  
all these points.

This is called Overdetermined System

(more nodes available than needed)

relates to Machine Learning

Solution: A straight line such that all points are equal dist from the line.



Q)  $n=1, \quad x_0 = -3, \quad x_1 = 0, \quad x_2 = 6$   
 $f(x_0) = 0, \quad f(x_1) = 0, \quad f(x_2) = 2$

produces a square matrix  
#  $Ax = b$  (6.2)  
 $(A^T A)x = A^T b$

Ans  $P_1(x) = a_0 + a_1 x$   
 $P_1(x_0) = a_0 + a_1 x_0 = f(x_0)$   
 $P_1(x_1) = a_0 + a_1 x_1 = f(x_1)$   
 $P_1(x_2) = a_0 + a_1 x_2 = f(x_2)$

No square matrix  
hence No inverse

$$\begin{pmatrix} 1 & -3 \\ 1 & 0 \\ 1 & 6 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 2 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} 1 & 1 & 1 \\ -3 & 0 & 6 \end{pmatrix} \begin{pmatrix} 1 & -3 \\ 1 & 0 \\ 1 & 6 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ -3 & 0 & 6 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 2 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} 3 & 3 \\ 3 & 45 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 2 \\ 12 \end{pmatrix}$$



Theorem 6.2. Let  $Q$  be an  $m \times n$  matrix. The columns of  $Q$  form an orthonormal set iff  $Q^T Q = I_n$ . If  $m = n$ , then such a  $Q$  is called an orthogonal matrix. For  $m \neq n$ , it is just called a matrix with orthonormal columns.

Proof. Let  $q_1, q_2, \dots, q_n$  be the columns of  $Q$ . Then

$$Q^T Q = \begin{pmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_n^T \end{pmatrix} \begin{pmatrix} q_1 & q_2 & \cdots & q_n \end{pmatrix} = \begin{pmatrix} q_1^T q_1 & q_1^T q_2 & \cdots & q_1^T q_n \\ q_2^T q_1 & q_2^T q_2 & \cdots & q_2^T q_n \\ \vdots & \vdots & \ddots & \vdots \\ q_n^T q_1 & q_n^T q_2 & \cdots & q_n^T q_n \end{pmatrix}. \quad (6.5)$$

So orthonormality  $q_i^T q_j = \delta_{ij}$  is equivalent to  $Q^T Q = I_n$ , where  $I_n$  is the  $n \times n$  identity matrix.

► Note that the columns of  $Q$  are a basis for  $\text{range}(Q) = \{Qx : x \in \mathbb{R}^n\}$ .

Example → The set  $S = \left\{ \frac{1}{\sqrt{5}}(2, 1)^T, \frac{1}{\sqrt{5}}(1, -2)^T \right\}$ . The two vectors in  $S$  are orthonormal, since

$$\frac{1}{\sqrt{5}}(2, 1) \cdot \frac{1}{\sqrt{5}}(2, 1) = 1, \quad \frac{1}{\sqrt{5}}(1, -2) \cdot \frac{1}{\sqrt{5}}(1, -2) = 1, \quad \frac{1}{\sqrt{5}}(2, 1) \cdot \frac{1}{\sqrt{5}}(1, -2) = 0.$$

Therefore  $S$  forms a basis for  $\mathbb{R}^2$ . If  $x$  is a vector with components  $x_1, x_2$  in the standard basis  $\{(1, 0)^T, (0, 1)^T\}$ , then the components of  $x$  in the basis given by  $S$  are

$$\begin{pmatrix} \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} & -\frac{2}{\sqrt{5}} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

► Inner products are preserved under multiplication by orthogonal matrices, since  $(Qx)^T Qy = x^T (Q^T Q)y = x^T y$ . This means that angles between vectors and the lengths of vectors are preserved. Multiplication by an orthogonal matrix corresponds to a rigid rotation (if  $\det(Q) = 1$ ) or a reflection (if  $\det(Q) = -1$ ).

Theorem 6.3. The matrix  $A^T A$  is invertible iff the columns of  $A$  are linearly independent, in which case  $Ax = b$  has a unique least-squares solution  $\hat{x} = (A^T A)^{-1} A^T b$ .

Proof. If  $A^T A$  is singular (non-invertible), then  $A^T A \hat{x} = 0$  for some non-zero vector  $\hat{x}$ , implying that

$$x^T A^T A \hat{x} = 0 \implies \|A\hat{x}\|_2^2 = 0 \implies A\hat{x} = 0. \quad (6.8)$$

This implies that  $A$  is rank-deficient (i.e. its columns are linearly dependent).

Conversely, if  $A$  is rank-deficient, then  $A\hat{x} = 0$  for some  $\hat{x} \neq 0$ , implying  $A^T A \hat{x} = 0$  and hence that  $A^T A$  is singular. □

► The  $n \times m$  matrix  $(A^T A)^{-1} A^T$  is called the Moore-Penrose pseudoinverse of  $A$ . In practice, we would solve the normal equations (6.7) directly, rather than calculating the pseudoinverse itself.

Matrix  $A$

$$\begin{pmatrix} 1 & x_0 \\ 1 & x_1 \\ 1 & x_2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} f(x_0) \\ f(x_1) \\ f(x_2) \end{pmatrix} \quad \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 3 & 3 \\ 3 & 45 \end{pmatrix}^{-1} \begin{pmatrix} 2 \\ 12 \end{pmatrix}$$

$$\det(A^T A) = 126 \neq 0 \text{ hence has inverse}$$

$$\therefore \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \frac{1}{126} \begin{pmatrix} 45 & -3 \\ -3 & 3 \end{pmatrix} \begin{pmatrix} 2 \\ 12 \end{pmatrix} = \begin{pmatrix} 3/14 \\ 5/21 \end{pmatrix}$$

$$\therefore P_1(x) = a_0 + a_1 x$$

$$= \frac{3}{14} + \frac{5}{21}x$$

$$\# Ax = b$$

$$\# Ax = b$$

$$A^T A x = A^T b$$

$$\# A = QR$$

$$x = (A^T A)^{-1} A^T b$$

$$\# R = \text{Upper triangular matrix}$$

inverting is a hassle  $\# Q = \text{All columns form an orthonormal set}$

In practice the normal matrix  $A^T A$  can often be ill conditioned. Thus we do QR decomposition.

$$\boxed{A = QR}$$

Theorem

$$Q = \begin{bmatrix} 2 & 4 \\ 3 & 5 \\ 1 & 6 \\ q_3 \\ q_4 \end{bmatrix} \quad \begin{aligned} q_1 &= (2, 3, 1)^T \\ q_2 &= (4, 5, 6)^T \end{aligned} \quad 2 \times 2 \quad \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$Q^T Q = I \rightarrow \text{Identity matrix}$$

$$3 \times 3 \quad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

### Gram-Schmidt Process

Q) From three given two vectors, find the orthonormal vectors using Gram-Schmidt process.

$$\text{Given } A = \begin{bmatrix} 3 & 1 \\ 6 & 2 \\ 0 & 2 \end{bmatrix} \text{ or } u_1 = (3, 6, 0)^T, u_2 = (1, 2, 2)^T$$

$$P_k = u_k - \sum_{i=1}^{k-1} (u_i^T q_i) q_i, \quad q_k = \frac{P_k}{\|P_k\|}$$

$$\text{Ans} \quad p_1 = u_1 = \begin{bmatrix} 3 \\ 6 \\ 0 \end{bmatrix}, \quad q_1 = \frac{p_1}{\|p_1\|} = \frac{\begin{pmatrix} 3 \\ 6 \\ 0 \end{pmatrix}}{\sqrt{3^2 + 6^2 + 0^2}}$$

$$p_2 = u_2 - (u_1^T q_1) q_1, \quad q_2 = \frac{p_2}{\|p_2\|}$$

$$= \left( \frac{1}{2} \right) - \left\{ \left( 1, 2, 2 \right) \cdot \frac{1}{\sqrt{45}} \begin{pmatrix} 3 \\ 6 \\ 0 \end{pmatrix} \right\} \cdot \frac{1}{\sqrt{45}} \begin{pmatrix} 3 \\ 6 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 2 \end{pmatrix} \quad \therefore Q = [q_1 \ q_2]$$

$$Ax = b, A = QR$$

$$\Rightarrow A^T A x = A^T b$$

$$\Rightarrow (QR)^T QR x = (QR)^T b$$

$$\Rightarrow Q^T R^T QR x = Q^T R^T b$$

$$\Rightarrow (Q^T Q) R^T R x = R^T Q^T b$$

$$\Rightarrow R x = Q^T b$$

$$\therefore x = R^{-1} Q^T b$$

$$\text{Q) Given } (-3, 0), (0, 0), (6, 2) \text{ find } P_1(x)$$

3 nodes but degree 1: Overdetermined system

$$\text{Ans: } \begin{pmatrix} 1 & -3 \\ 1 & 0 \\ 1 & 6 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 2 \end{pmatrix}$$

$$A \quad x \quad b$$

$$u = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad u_2 = \begin{pmatrix} -3 \\ 0 \\ 6 \end{pmatrix}$$

$\because R = 2 \times 2$  upper triangular matrix

$$i.e. Q = \begin{pmatrix} 1/\sqrt{3} & -4/\sqrt{42} \\ 1/\sqrt{3} & -1/\sqrt{42} \\ 1/\sqrt{3} & 5/\sqrt{42} \end{pmatrix}$$

$$R = \begin{pmatrix} u_1^T q_1 & u_2^T q_1 \\ 0 & u_2^T q_2 \end{pmatrix} = \begin{pmatrix} \sqrt{3} & \sqrt{3} \\ 0 & \sqrt{42} \end{pmatrix}$$

$$\# R x = Q^T b$$

$$\begin{pmatrix} \sqrt{3} & \sqrt{3} \\ 0 & \sqrt{42} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \sqrt{3} & \sqrt{3} & 1/\sqrt{3} \\ -4/\sqrt{42} & -1/\sqrt{42} & 5/\sqrt{42} \end{pmatrix} \begin{pmatrix} 0 \\ 2 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \sqrt{3} & \sqrt{3} \\ 0 & \sqrt{42} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} 2/\sqrt{3} \\ 10/\sqrt{42} \end{pmatrix}$$

$$\therefore \sqrt{42} a_1 = 10/\sqrt{42} \quad \left| \begin{array}{l} \sqrt{3} a_0 + \sqrt{3} a_1 = 2/\sqrt{3} \\ a_0 = 3/7 \end{array} \right.$$

$$a_1 = 5/21$$

$$\therefore P_1(x) = a_0 + a_1 x$$

$$= 3/7 + 5/21 x$$

OR

$$A = QR$$

$$Q^T A = Q^T Q R$$

$$\therefore Q^T A = R$$

$$\therefore R = \begin{pmatrix} \sqrt{3} & \sqrt{3} & 1/\sqrt{3} \\ -4/\sqrt{42} & -1/\sqrt{42} & 5/\sqrt{42} \end{pmatrix} \begin{pmatrix} 1 & -3 \\ 1 & 0 \\ 1 & 6 \end{pmatrix}$$

$$= \begin{pmatrix} \sqrt{3} & \sqrt{3} \\ 0 & \sqrt{42} \end{pmatrix}$$

same result

## Numerical Integration - 1

$$I(f) = \int_a^b f(x) dx, \quad I_n(f) = \int_a^b P_n(x) dx$$

To computer, it's polynomial  
Lagrange basis

$$P_n(x) = \sum_{k=0}^n l_k(x) f(x_k)$$

$$\therefore I_n(f) = \int_a^b \sum_{k=0}^n l_k(x) f(x_k) dx$$

constraint

$$= \sum_{k=0}^n f(x_k) \int_a^b l_k(x) dx$$

Weight factor,  $\delta_k = \int_a^b l_k(x) dx$

$$I_n(f) = \sum_{k=0}^n \delta_k f(x_k)$$

Newton-Cotes formula  
(7.1)

The nodes given, will all be equidistant to each other.  $a = x_0 < x_1 < x_2 \dots < x_n = b$

# The formula for these kind of nodes is called Closed Newton-Cotes Formula.  $h = \frac{b-a}{n}$

# When  $a < x_0 < x_1 < x_2 < \dots < x_n < b$ , it's called Open Newton-Cotes Formula.  $h = \frac{b-a}{n+2}$

# When in Newton-Cotes formula  $n=1$  (1 degree polynomial)  $P_1(x)$  then  $h = b-a$ . This is called Trapezium Rule.  $x_0 = a, x_1 = b$

$$\# P_1(x) = l_0(x) \cdot f(x_0) + l_1(x) \cdot f(x_1)$$

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

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

$$\# \delta_0 = \int_a^b l_0(x) dx = \int_a^b \frac{x-b}{a-b} dx = \frac{1}{a-b} \int_a^b x-b dx = \frac{1}{a-b} \left[ \frac{x^2}{2} - bx \right]_a^b = \frac{1}{a-b} \left( \frac{b^2}{2} - b^2 - \frac{a^2}{2} + ab \right) = \frac{b-a}{2}$$

$$\# \delta_1 = \int_a^b l_1(x) dx = \int_a^b \frac{x-a}{b-a} dx = \frac{b-a}{2}$$

$$\# I_n(f) = \sum_{k=0}^n \delta_k f(x_k)$$

$$I_n(f) = \delta_0 f(x_0) + \delta_1 f(x_1) = \frac{b-a}{2} f(a) + \frac{b-a}{2} f(b)$$

Q)  $a=0, b=2, f(x)=e^x$

$$\text{Ans} \quad I = \int_a^b f(x) dx = \int_0^2 e^x dx = [e^x]_0^2 = e^2 - e^0 = 6.389$$

$$I_1 = \frac{b-a}{2} [f(a) + f(b)] = \frac{2-0}{2} [e^0 + e^1] = 3.389$$

$$\text{Relative error} = \frac{|I - I_1|}{|I|} \times 100 = \frac{6.389 - 3.389}{6.389} \times 100 \approx 31.3\%$$

Upper Bound =  $\frac{f^{(n+1)}(\xi)}{(n+1)!} (x-x_0)(x-x_1)\dots(x-x_n), \quad \xi \in [a, b]$

As 1 degree polynomial, this high relative error  
Higher degree will have higher Accuracy

for  
integration

Q)  $a=0, b=2, f(x)=e^x, n=1$

$$\text{Ans} \quad \left| \frac{f^{(n+1)}(\xi)}{(n+1)!} \right|_{\max} \rightarrow \left| \frac{f''(\xi)}{2!} \right|_{\max} = \frac{1}{2} e^\xi \Big|_{\max \xi \in [0, 2]} = \frac{e^2}{2}$$

$$\int_0^2 (x-0)(x-2) dx = \int_0^2 (x^2 - 2x) dx = \left[ \frac{x^3}{3} - 2\frac{x^2}{2} \right]_0^2 = \frac{4}{3}$$

∴ Upper bound  $\leq \frac{e^2}{2} \times \frac{4}{3} = 4.926$

► Theorem 7.1 suggests that the accuracy of  $I_n$  is limited both by the smoothness of  $f$  (outside our control) and by the location of the nodes  $x_k$ . If the nodes are free to be chosen, then we can use Gaussian integration (see later).

► As with interpolation, taking a high  $n$  is not usually a good idea. One can prove for the closed Newton-Cotes formula that

$$\sum_{k=0}^n |\sigma_k| \rightarrow \infty \quad \text{as } n \rightarrow \infty.$$

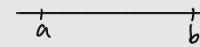
This makes the quadrature vulnerable to rounding errors for large  $n$ .

# Composite Newton-Cotes Formulae (7.2)

## Numerical Integration-2

[a, b] Here make some extra points between a & b  
 $n=1$  to make the accuracy better.  
 2 nodes

m ↑ Accuracy ↑



# m = no. of trapezoids if 1 is divided  $\rightarrow m=1$  means no iteration

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

m+1 nodes

★ Only 1 degree polynomials

Q)  $a=0, b=2, f(x)=e^x, m=2$ . Use composite Newton-Cotes Formula

$$I = \int_0^2 e^x dx = 6.389 \quad h = \frac{2-0}{2} = 1, \quad x_0=0, \quad x_1=0+1=1, \quad x_2=1+1=2$$

$$\text{Q) } m=3, \quad h = \frac{2-0}{3} = \frac{2}{3}$$

$$x_0=0, \quad x_1=0+\frac{2}{3}=2/3, \quad x_2=\frac{2}{3}+\frac{2}{3}=\frac{4}{3}, \quad x_3=\frac{4}{3}+\frac{2}{3}=2$$

$$\text{Q) } m=4, \quad h = \frac{2-0}{4} = \frac{1}{2}$$

$$x_0=0, \quad x_1=0+1/2=1/2, \quad x_2=1, \quad x_3=3/2, \quad x_4=2$$

$$C_{1,2} = \frac{h}{2} [f(x_0) + 2f(x_1) + f(x_2)]$$

$$= \frac{h}{2} [e^0 + 2e^1 + e^2] = 6.9126$$

h halving  $\rightarrow$  error decreases by 4

$\therefore O(h^2)$  Quadratic Convergence

$$\therefore C_{1,3} = \frac{h}{2} [f(x_0) + 2f(x_1) + 2f(x_2) + f(x_3)]$$

$$= \frac{2}{3/2} [e^0 + 2e^{2/3} + 2e^{4/3} + e^2] = 6.623$$

Closer to the actual value

$$\therefore C_{1,4} = \frac{0.5}{2} [f(x_0) + 2(f(x_1) + f(x_2) + f(x_3)) + f(x_4)] = 6.5216$$

$$C_{1,m} = \frac{h}{2} [f(x_0) + 2f(x_1) + \dots + 2f(x_{m-1}) + f(x_m)]$$

1 degree

#  $n=1 \rightarrow$  Trapezoidal Rule

#  $n>1 \rightarrow$  Simpson's Rule

\* Using Simpson's Rule, till  $n=2$  &  $3$  we can get the exact value.

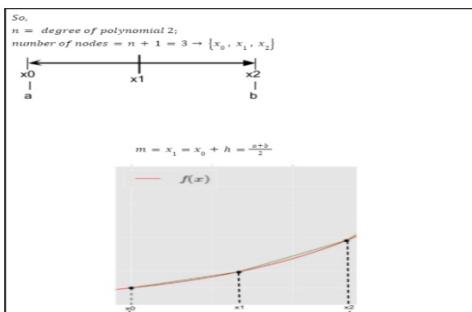
## Simpson's Rule (Derivation is Important)

(This is the  $n=2$  Newton-Cotes formula)

### Simpson's rule

From, Trapezium / Trapezoidal rule  $\rightarrow I_1(f) = \int_a^b P_1(x) dx$

Simpson's rule  $\rightarrow I_2(f) = \int_a^b P_2(x) dx$



$$\begin{aligned} I_2(f) &= \int_a^b P_2(x) dx \\ I_2(f) &= \int_a^b [l_2(x)(x-x_0)(x-x_2)] dx \\ &= I_1(f) + l_2(x_0)dx f(x_0) + l_2(x_2)dx f(x_2) \\ &= I_1(f) + \sigma_0 f(x_0) + \sigma_1 f(x_1) + \sigma_2 f(x_2) \\ &= I_1(f) + \sigma_0 f(a) + \sigma_1 f(m) + \sigma_2 f(b) \\ \text{Here,} \\ \sigma_0 &= \int_a^b l_2(x) dx \\ &= \int_a^b \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} dx \\ &= \sigma_0 = \int_a^b \frac{(x-m)(x-b)}{(m-a)(m-b)} dx \\ &= \sigma_0 = \frac{1}{(m-a)(m-b)} \int_a^b (x-m)(x-b) dx \\ &= \sigma_0 = \frac{1}{(m-a)(m-b)} \times \frac{(b-a)(m-b)(m-a)}{6} \\ &= \sigma_0 = \frac{1}{6}(b-a) \end{aligned}$$



$$\begin{aligned} \sigma_1 &= \int_a^b l_1(x) dx \\ &= \int_a^b \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} dx \\ &= \sigma_1 = \int_a^b \frac{(x-a)(x-m)}{(m-a)(m-b)} dx \\ &= \sigma_1 = \frac{1}{(m-a)(m-b)} \int_a^b (x-a)(x-m) dx \\ &= \sigma_1 = \frac{1}{(m-a)(m-b)} \times \frac{(b-a)(m-b)(m-a)}{6} \\ &= \sigma_1 = \frac{1}{6}(b-m) \end{aligned}$$

$$\begin{aligned} \sigma_2 &= \int_a^b l_2(x) dx \\ &= \sigma_2 = \int_a^b \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} dx \\ &= \sigma_2 = \int_a^b \frac{(x-b)(x-m)}{(m-a)(m-b)} dx \\ &= \sigma_2 = \frac{1}{(b-a)(m-b)} \int_a^b (x-b)(x-m) dx \\ &= \sigma_2 = \frac{1}{(b-a)(m-b)} \times \frac{(b-a)(m-b)(m-a)}{6} \\ &= \sigma_2 = \frac{1}{6}(b-a) \end{aligned}$$

$$\begin{aligned} \text{So,} \\ I_2(f) &= \sigma_0 f(x_0) + \sigma_1 f(x_1) + \sigma_2 f(x_2) \\ I_2(f) &= \sigma_0 f(a) + \sigma_1 f(m) + \sigma_2 f(b) \\ I_2(f) &= \frac{b-a}{6} f(a) + \frac{2(b-a)}{3} f\left(\frac{m+a}{2}\right) + \frac{b-a}{6} f(b) \\ I_2(f) &= \frac{b-a}{6} f(a) + \frac{4(b-a)}{6} f\left(\frac{a+b}{2}\right) + \frac{b-a}{6} f(b) \end{aligned}$$

$$I_2(f) = \frac{b-a}{6} [f(a) + 4f\left(\frac{a+b}{2}\right) + f(b)] \rightarrow \text{Simpson's formula}$$

### 1) Example:

Find  $I_1(f)$  and  $I_2(f)$  of the function  $e^x$  on interval [0, 2]. Show relative error.

Solution  $\rightarrow$

$$I_1(f) = \int_0^2 e^x dx = [e^x]_0^2 = e^2 - e^0 = 6.389 \quad [\text{Actual Integration}]$$

$$I_1(f) = \frac{3-a}{6} [f(a) + 4f\left(\frac{a+b}{2}\right) + f(b)] \quad [\text{Approximate Integration}]$$

$$= I_1(f) = \frac{3-0}{6} [e^0 + 4e^1 + e^2]$$

$$= I_1(f) = \frac{1}{3} [e^0 + 4e^1 + e^2] = 6.421$$

$$\text{relative error (in percentage)} = \frac{|I_1 - I_2|}{|I_1|} \times 100 = \frac{|6.421 - 6.389|}{|6.389|} \times 100 = 0.501\%$$

### 2) Example:

Compute the upper bound error for  $I_2$  of the function  $e^x$  on interval [0, 2].

Solution  $\rightarrow$

Here,  $n = 2$

$$f(x) = e^x$$

$$[a, b] = [0, 2]$$

For integration, upper bound error  $\rightarrow$

$$|I - I_n| \leq \left| \frac{f^{(n+1)}(t)}{(n+1)!} \right| \int_a^b |(x - x_0)(x - x_1) \dots (x - x_n)| dx$$

$$\Rightarrow |I - I_2| \leq \left| \frac{f^{(3)}(t)}{3!} \right| \int_0^2 |(x - x_0)(x - x_1)(x - x_2)| dx$$

$$\Rightarrow |I - I_2| \leq \left| \frac{f^{(3)}(t)}{3!} \right| \int_0^2 |(x - a)(x - \frac{a+b}{2})(x - b)| dx$$

Now, finding the max of  $\left| \frac{f^{(3)}(t)}{3!} \right|$  within [0, 2]

$$= \left| \frac{f'''(x)}{3!} \right| = \left| \frac{e'''}{6} \right| = \left| \frac{e''}{2} \right| = \left| \frac{e'}{6} \right| = \left| \frac{e^x}{6} \right| \quad [\text{where } |f''(x)| = |e^x|]$$

Now, finding the max of  $\left| \int_0^2 |(x - a)(x - \frac{a+b}{2})(x - b)| dx \right|$  within [0, 2]

$$= \int_0^2 |(x - a)(x - \frac{a+b}{2})(x - b)| dx$$

$$= \int_0^2 |(x - 0)(x - \frac{0+2}{2})(x - 2)| dx$$

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

$$= \int_0^2 |(x^3 - 3x^2 + 2x)| dx$$

$$= \left| \frac{x^4}{4} - \frac{3x^3}{3} + \frac{2x^2}{2} \right|_0^2 = 2$$

So, upper bound of error  $\leq \frac{e^2}{6} \times 2 \approx 2.463$

A coordinate system is shown on a background of a 20x20 grid of small grey dots. A vertical black arrow points upwards from the origin, and a horizontal black arrow points to the right from the origin. The word "THE" is written in black lines above the origin, and the word "END" is written in black lines below the origin.

THE

END