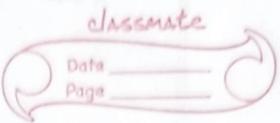


# SLIDES - LEC 2



## ERRORS

① Truncation error - Replacing infinite process with a finite one.

$\pi$  value with 3.14

② Rounding error - Due to limitations in representing numbers with finite precision.

A small error in the input can be amplified.

Reason -

- ① by a bad problem - example  $y = 1/x$
- ② by bad algorithm

Two types of error -

① Absolute error =  $||\text{True Value} - \text{Approx Value}||$

② Relative error =  $\frac{\text{Absolute Error}}{||\text{True Value}||}$

For Rounding errors - Relative error is used.

Normalized floating point representation -

$x \in \mathbb{R}$  and  $B > 1$

$$x = \pm (d_1 d_2 \dots d_t \dots)_{\beta} \times \beta^e$$

$$0 \leq d_i \leq \beta \quad d_1 \neq 0, e \in \mathbb{Z}$$

$\beta \rightarrow \text{base}$

$e \rightarrow \text{exponent}$

$f = (d_1 d_2 \dots d_t \dots)$  → mantissa / fraction

$$f = \pm d_1 (\cdot d_1 d_2 \dots d_t \dots)$$

such a way that it is  $\infty$

$$= \sum_{i=1}^{\infty} \frac{d_i}{\beta^i}$$

$$x = \pm f \times \beta^e$$

The number of digits allowed in  $f$   
called precision of floating point representation.  
which in above case is  $\infty$ .

Computer use  $\beta=2$ .  $\Rightarrow$  Floating point system

$$F(\beta, t, e_{\min}, e_{\max})$$

$$\hookrightarrow f = \pm (d_1 d_2 \dots d_t)_{\beta} \times \beta^e, \quad 0 \leq d_i < \beta, \quad d_i \neq 0, \\ e_{\min} \leq e \leq e_{\max} \quad \exists \cup \{0\}$$

$$f = (d_1 d_2 \dots d_t)_{\beta} = \sum_{j=1}^{t+1} \frac{d_j}{\beta^j} \leq 1 - \beta^{-t} < 1$$

$$\text{Example} - +0.101 \times 2^2 = (1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3}) \times 2^2$$

$$= \frac{5}{2} \times$$

$$\text{no. of } F(\beta, t, e_{\min}, e_{\max}) \rightarrow 2(\beta-1)(\beta)^{t-1}(e_{\max} - e_{\min}) + 1$$

0.9999

0.00001

Underflow Threshold -

$$= (0.1000)_\beta \times \beta^{e_{\min}} = \beta^{e_{\min}-1}$$

Overflow method -

$$= (\cdot (B-1) (B-1) \dots) \times \beta^{e_{\max}}$$

$$= (1 - \beta^{-t}) \times \beta^{e_{\max}}$$

why?  $\rightarrow$  $(\cdot (B-1) (B-1) (B-1) \dots) \rightarrow$  largest mantissa

$$= \sum_{j=1}^t \frac{\beta-1}{\beta^j}$$

$$= \beta - \sum_{j=1}^t \beta^{-j} = (\beta-1) \frac{(\beta^{-t} - 1)}{\beta^{-1} - 1}$$

$$= (\beta-1) \frac{(-1 + \beta^{-t})}{1 - \beta} = 1 - \beta^{-t}$$

Machine Epsilon / Precision  $\Rightarrow \text{eps} = \beta^{1-t}$ 

$$\text{eps} = (0.10 \dots 01)_\beta \times \beta^1 - (-100 \dots 00)_\beta \times \beta^1$$

$$= \beta (\beta^{-t}) = \beta^{1-t}$$

Machine epsilon is the smallest the no mat when added to 1, result in a value different from 1 due to the bin

# Lec-3 SLIDES

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

IEEE single precision

(... floating point representation)

IEEE specifies  $f(2, t, e_{\min}, e_{\max})$  -

sign (1 bit)	exponent (8 bits)	mantissa (23 bits)
--------------	-------------------	--------------------

# Exponents =  $2^8 = 256$

Largest Exponent =  $255$  (i.e.,  $b=127$ )

Exponent Range :  $0 \leq b \leq 255$  and  
 $-127 \leq e \leq 128$

$$x = (-1)^s \times (1 \cdot d_1 \dots d_{23})_2 \times 2^{b-127}$$

$$x = (-1)^s \times (1 \cdot d_1 \dots d_{23})_2 \times 2^{b-127}$$

$$\therefore x = ((\pm 1) \dots (\pm 1)) \dots (\pm 1) = \text{where } e = b - 127$$

Effective Range  $\rightarrow -126 \leq e \leq 127$

$e = -127$  is reserved for  $\pm 0$  and

$e = 128$  reserved for  $\pm \infty$  and NaN

## IEEE Double Precision Floating Point Representation

IEEE specifies  $f(2, t, e_{\min}, e_{\max}) -$

$(2^{e-b}) \times (1.d_1 \dots d_{52})_2$

$(\text{sign} (1 \text{ bit})) \mid \text{exponent } (8 \text{ bits}) \mid \text{mantissa } (52 \text{ bits})$

# exponents =  $2^8 = 2048$

Largest Exponent =  $2^8 - 1 = 2047$

Exponent range :  $0 \leq b \leq 2047$  and  $-1023 \leq e \leq 1024$

$$x = (-1)^s \times (1.d_1 \dots d_{52})_2 \times 2^{b-1023}$$

where  $e = b - 1023$

Effective range :  $-1022 \leq e \leq 1023$

$e = -1023$  is reserved for  $\pm 0$  and

$e = 1025$  is reserved for  $\pm \infty$  and NaN

## Gap b/w floating point numbers

Let  $x \in f(\beta, t, l_{\min}, l_{\max})$

$$x = (d_1 d_2 \dots d_t) \times \beta^e$$

$$\text{next}(x) = x + \text{ulp}(x)$$

$x$  and  $\text{next}(x)$  are consecutive floating point numbers.

$\text{next}(x)$  is the smallest floating point number larger than  $x$ .

Why?

$$\text{ulp} = \beta^{e-t}$$

If we increase the  $d_t$  by one then  $\beta^{-t}$  is added in the significant part.

so if the mantissa is increased by  $\beta^{-t}$  then  $x$  is increased by  $\beta^{-t} \cdot \beta^e$

$$= \beta^{e-t} = \text{ulp}(x)$$

Relative gap b/w floating point numbers

$$\text{Relative gap} = \frac{\text{next}(x) - x}{x}$$

$$= \frac{\beta^{e-t}}{(d_1 d_2 \dots d_t)_\beta \times \beta^e}$$

$$\leq \frac{\beta^{e-t}}{(0.10000)_\beta \times \beta^e}$$

$$\leq \frac{\beta^{-t}}{\beta^{-1}} \leq \beta^{1-t} = \text{eps}$$

$$\text{Relative gap} \leq \text{eps}$$

## Rounding

A number  $x \in R$  can be represented by a number in  $F(\beta, t, e_{\min}, e_{\max})$  provided  $\text{realmin} \leq |x| \leq \text{realmax}$ .

Let  $f(x)$  denotes the floating point representation of  $x$  that is -

$$f(x) : R \rightarrow F(\beta, t, e_{\min}, e_{\max})$$

$$x \rightarrow f(x)$$

$$\text{Let } x = ( \cdot d_1 d_2 \dots d_t \dots )_{\beta} \times \beta^e$$

$$x_L = ( \cdot d_1 d_2 \dots d_t )_{\beta} \times \beta^e$$

$$x_R = x_L + \beta^{e-t}$$

$$x_L \leq x \leq x_R$$

$$f(x) = \begin{cases} x_L & \text{round down} \\ x_R & \text{round up} \\ x_L \text{ or } x_R & \text{round to the nearest} \end{cases}$$

# Rounding Error

Theorem -

Let  $f: \mathbb{R} \rightarrow F(\beta, t, L, U)$ ,  $x \mapsto f(x)$

Then,

$$f(x) = x(1 + \delta) \text{ with } |\delta| \leq u$$

$$u = \begin{cases} \beta^{1-t} & \text{round down} \\ \frac{\beta^{1-t}}{2} & \text{round to nearest} \end{cases}$$

Proof -

$$x = d_0(d_1 d_2 \dots d_t \dots) \times \beta^e$$

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

$$f(x) = x(1 + \delta)$$

Now -

$$|f(x) - x| \leq \beta^{e-t} \quad \left\{ \begin{array}{l} \text{for round to down} \\ \leq \frac{1}{2} \cdot \beta^{e-t} \quad \left\{ \begin{array}{l} \text{for round to nearest} \\ t^m \end{array} \right. \end{array} \right.$$

why?

$$(f(x) - x) = (0.0000 \dots d_{t+1} d_{t+2} \dots) \times \beta^e$$

$$\leq (0.0000 \dots 1) \times \beta^e = \beta^{e-t}$$

$$|x - f(x)| \leq \beta^{e-t}$$

$$\frac{|x - f(x)|}{|x|} \leq \frac{\beta^{e-t}}{(0. d_1 d_2 \dots d_t) \times \beta^e}$$

$$\leq \frac{\beta^{e-t}}{(0.100 \dots 1) \beta^e} = \beta^{1-t}$$

## Floating Point Arithmetic

Case 1

If  $x, y \in F(\beta, t, e_{\min}, e_{\max})$

$$f_1(x \oplus y) = (x \oplus y)(1 + \delta)$$

where  $|\delta| \leq u$

Case 2

If  $x$  and  $y$  are not in floating point representation.

$$\hat{x} = f_1(x) = x(1 + \delta_1)$$

$$\hat{y} = f_1(y) = y(1 + \delta_2)$$

$$(x \pm y) = (x \pm y)(1 + \frac{x\delta_1 + y\delta_2}{x \pm y})$$

Relative Error -

$$\left| \frac{x\delta_1 + y\delta_2}{x + y} \right| \leq \delta_1 + \delta_2 \leq 2u \quad \{ \text{small} \}$$

For +ve

Relative Error -

$$\left| \frac{x\delta_1 - y\delta_2}{x - y} \right| \quad \text{can be very large}$$

when  $x - y$  is very small.

(THIS IS LOSS OF SIGNIFICANCE  
OR CATASTROPHIC CANCELLATION)

## Loss of significance / catastrophic cancellation -

- ① When two close float points are subtracted.
- ② Their leading digits cancels out.
- ③ The result has fewer significant than the input floating point.

Example -

$$\text{for } ax^2 + bx + c = 0$$

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

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

Here to avoid this  
x<sub>2</sub> calculation  
to avoid subtraction  
and hence to avoid  
loss of significance.

Correct way -

$$\text{calculate } x_1 \text{ and then } x_1 x_2 = \frac{c}{a}.$$

Remedy -

Avoid subtraction by changing the algorithm method.

## Rounding Error of Differentiation -

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

For small  $h$ ,  $D_h f(x)$  exhibits cancellation error due to rounding, which is magnified by the division by  $h$ .

$$\text{let } \hat{f}(x) = f(x + \delta(x))$$

$$\hat{f}(x) = f(x)(1 + \delta(x))$$

where  $\delta(x) \leq \text{eps}$

$$\overset{\wedge}{D}_h f(x) = \frac{\hat{f}(x+h) - \hat{f}(x)}{h}$$

$$= \frac{f(x+h)(1 + \delta(x+h)) - f(x)(1 + \delta(x))}{h}$$

$$= D_h f(x) + \frac{\delta(x+h)f(x+h) - \delta(x)f(x)}{h}$$

Rounding Error =

$$= e_r(x, h) = \left| \frac{\delta(x+h)f(x+h) - \delta(x)f(x)}{h} \right|$$

$$\leq \left| \frac{\delta(x+h)f(x+h) + \delta(x)f(x)}{h} \right|$$

$$\leq \left| \frac{f(x+h) \cdot \text{eps} + \delta(x)f(x) \cdot \text{eps}}{h} \right|$$

$$\leq \text{eps} \cdot \left| \frac{f(x+h) + f(x)}{h} \right|$$

continued ↗

Now By Taylor's Thm -

Truncation error made when approximating a sum by finite sum

$$f(x+h) = f(x) + f'(x) \cdot h + \frac{f''(\xi)}{2} \cdot h^2$$

$$D_h(f(x)) = f'(x) + \frac{f''(\xi)}{2} \cdot h^{\frac{1}{2}}$$

Truncation Error ←  $E(x, h) = \left| \frac{f''(\xi)}{2} \cdot h \right|$

This error should decrease as  $h \rightarrow 0$ .

Rounding Error =  $\left| f(x+h) + f(x) \right| \text{eps}_h$

This error increases as  $h \rightarrow 0$

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

$$\approx \frac{h}{2} |f''(\xi)| + \frac{\epsilon}{h} (|f(x+h) + f(x)|)$$

Ignoring constants + magnitudes.

for balancing the two errors to minimize total error -

$$h^2 \approx \epsilon$$

$$h \approx \sqrt{\epsilon}$$

we get better differentiation results.

## POLYNOMIAL

## INTERPOLATION - I

Polynomial Interpolation ?

Problem - Given a dataset  $(x_0, f_0), \dots, (x_n, f_n)$   
consisting of

distinct nodes :  $[x_0, x_1, \dots, x_n]$

and values :  $[f_0, f_1, \dots, f_n]$

construct a polynomial  $p(x)$  of lowest degree  
such that  $p(x_j) = f_j$  for  $j = 0 : n$ .

Vandermonde Interpolating Polynomial.

Thm - Consider nodes  $[x_0, \dots, x_n]$   
values  $[f_0, \dots, f_n]$

There exist a polynomial of degree  
at most  $n$  such that  $p(x_j) = f_j$   
for  $j = 0 : n$

Proof - Consider the polynomial  $p(x) = a_0 x_0^0 + a_1 x_1^1 + \dots + a_n x_n^n$

$p(x_j) = f_j$  for  $j = 0 : n$

$$\begin{bmatrix} 1 & x_0 & \dots & x_0^n \\ 1 & x_1 & \dots & x_1^n \\ 1 & \vdots & \vdots & \vdots \\ 1 & x_n & \dots & x_n^n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_n \end{bmatrix}$$

Vandermonde Matrix

$$VA = Y$$

det

$$\det(V) \neq 0 \quad \text{so unique solution.}$$

- The computation is numerically unstable.
- Solution of system require  $O(n^3)$  operations.
- Any additional new data  $(x_{n+1}, f_{n+1})$  requires recompute.
- Evaluation of  $p(x)$  at given require  $O(n^2)$  operations.

### Interpolating polynomial in General basis

Let  $P_n$  denote the vector space of polynomials of degree at most  $n$ . Let  $\phi_0(x), \phi_1(x), \dots, \phi_n(x)$  be a basis of  $P_n$ . Let  $p(x) = a_0 \phi_0(x) + a_1 \phi_1(x) + \dots + a_n \phi_n(x)$ .

$\phi_0(x_0)$	$\phi_1(x_0)$	$\dots$	$\phi_n(x_0)$	$a_0$	$f_0$
$\phi_0(x_1)$	$\phi_1(x_1)$	$\vdots$	$\phi_n(x_1)$	$a_1$	$f_1$
$\phi_0(x_2)$	$\phi_1(x_2)$	$\vdots$	$\phi_n(x_2)$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\phi_0(x_n)$	$\phi_1(x_n)$	$\dots$	$\phi_n(x_n)$	$a_n$	$f_n$

Coefficient matrix is non singular.

The linear system has unique solution.

LAGRANGE INTERPOLATING POLYNOMIAL

Example - Let say for three data points -

$$(x_0, f_0), (x_1, f_1), \dots (x_2, f_2)$$

define,

$$\begin{aligned} P(x) = & \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} f_0 + \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} f_1 \\ & + \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} f_2 \end{aligned}$$

$$\begin{aligned} \text{Then } P(x_0) &= f_0 \\ P(x_1) &= f_1 \\ P(x_2) &= f_2 \end{aligned}$$

NOW IN GENERAL -

for data set  $(x_0, f_0), (x_1, f_1), \dots (x_n, f_n)$

$$\text{define } w(x) = (x-x_0)(x-x_1) \dots (x-x_n) \in P_{n+1}$$

$$L(x) = \prod_{i \neq j} \left( \frac{x-x_i}{x_j-x_i} \right) \# = \prod_{i \neq j} \frac{(x-x_i)}{(x_j-x_i)}$$

$$w(x) = \prod_{m=0}^n (x-x_m)$$

$$w(x) = (x - x_j) \cdot \prod_{\substack{i=0 \\ i \neq j}}^n (x - x_i) \quad \text{--- (1)}$$

$$\begin{aligned} w'(x_j) &= \frac{d}{dx} \left[ (x - x_j) \cdot \prod_{\substack{i=0 \\ i \neq j}}^n (x - x_i) \right] \\ &= \prod_{\substack{i=0 \\ i \neq j}}^n (x_j - x_i) + (x_j - x_j) \cdot \cancel{\frac{d}{dx} \left( \prod_{\substack{i=0 \\ i \neq j}}^n (x - x_i) \right)} \end{aligned}$$

$$w'(x_j) = \prod_{\substack{i=0 \\ i \neq j}}^n (x_j - x_i) \quad \text{--- (2)}$$

from (1) and (2)

$$l_j(x) = \frac{w(x)}{(x - x_j) \cdot w'(x_j)}, \text{ for } j=0 \dots n$$

$$l_j(x) = \delta_{ij} \quad \text{where } \delta_{ij} \text{ is Dirac Delta Function.}$$

$$\text{Hence, } p(x) = f_0 l_0(x) + \dots + f_n l_n(x)$$

interpolates the dataset  $(x_0, f_0), \dots, (x_n, f_n)$ .

The basis  $l_0(x), l_1(x), \dots, l_n(x)$   
 is called Lagrange basis of  $P_n$  and  $p(x)$   
 is called Lagrange interpolating polynomial.

IMP

$$\sum_{j=0}^n l_j(x) = 1$$

at given  $x$

→ Computation of  $p(x)$  require  $O(n^2)$  operation.

Why?

Calculating single  $l_j(x)$  takes  
multiplication terms so  
takes  $O(n)$ .

$$\text{Total time} = O(n^2)$$

→ Cannot accommodate new  $(x_{n+1}, b_{n+1})$ . Any  
additional data requires recomputation.

A

Doubt

Computation of  $p(x)$  requires  $O(n^2)$  operations?  
G. In unexpanded form

## BARYCENTRIC LAGRANGE INTERPOLATION

$$l_j = \frac{w(x)}{(x - x_j) \cdot w'(x_j)} = \frac{w(x) \cdot w_j}{(x - x_j)}$$

where  $w_j = \frac{1}{w'(x_j)}$

$$\begin{aligned}
 p(x) &= \sum_{j=0}^n l_j(x) \cdot y_j(x) \\
 \text{Barycentric} \\
 \text{Form 1} &= w(x) \sum_{j=0}^n \frac{f_j(x) \cdot w_j}{x - x_j} \quad - \textcircled{1}
 \end{aligned}$$

Also,

$$\begin{aligned}
 1 &= \sum_{j=0}^n y_j(x) \\
 1 &= w(x) \sum_{j=0}^n \frac{w_j}{x - x_j} \quad - \textcircled{2}
 \end{aligned}$$

$$\textcircled{1} \div \textcircled{2}$$

$$\begin{aligned}
 \frac{p(x)}{1} &= \sum_{j=0}^n \frac{w_j \cdot f_j(x)}{x - x_j} = \sum_{j=0}^n \frac{w_j \cdot w(x) \cdot f_j(x)}{x - x_j} \\
 &\quad \sum_{j=0}^n \frac{w_j \cdot f_j(x)}{x - x_j} \cdot w(x)
 \end{aligned}$$

$$P(x) = \sum_{j=0}^n w_j f_j$$

$\frac{w_j f_j}{x - x_j}$   
↓  
 $\frac{w_j}{x - x_j}$

{ Barycentric Form }

### Advantages -

(3)

- ① Once I know  $w_j$ . Evaluating  $P(x)$  at  $x = x_k$  takes  $O(n)$  time.

However calculating every  $w_j$  takes  $O(n^2)$

- $O(n)$  for one  $j$ .
- $O(n^2)$  for all  $j$ .

- ② To incorporate a new data point  $(x_{n+1}, f_{n+1})$

- a) Updating an existing weight takes  $O(2n+1)$  operations.

$$w_j^{\text{old}} = \frac{1}{\prod_{\substack{k=0 \\ k \neq j}}^n (x_j - x_k)}$$

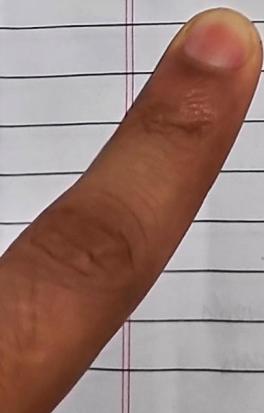
$$w_j^{\text{new}} = \frac{1}{\prod_{\substack{k=0 \\ k \neq j}}^{n+1} (x_j - x_k)} = \frac{w_j^{\text{old}}}{(x_j - x_{n+1})}$$

② Computing the new weight  $w_{n+1}$  takes  $O(2n+3)$

$$w_{n+1} = \frac{1}{\prod_{k=0}^n (x_{n+1} - x_k)}$$

③ Barycentric formulas has beautiful symmetry.

The weights  $w_j$  appear in denominator exactly as in numerator, except with factor  $f_j$ . This means any common factor in all the weights  $w_j$  may be cancelled without affecting the value of  $p(x)$ .



## LEC - 6

CLASSMATE  
Date \_\_\_\_\_  
Page \_\_\_\_\_

### NEWTON INTERPOLATING POLYNOMIAL

Example - For data set  $(x_0, f_0), (x_1, f_1), (x_2, f_2)$

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

$$p(x_0) = f_0$$

$$p(x_1) = f_1$$

$$p(x_2) = f_2$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & (x_1-x_0) & 0 \\ 1 & (x_2-x_0) & (x_2-x_0)(x_2-x_1) \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \end{bmatrix}$$

General -

$$N_0(x) = 1$$

$$N_j(x) = (x-x_0)(x-x_1)\dots(x-x_{j-1})$$

for  $j=1$  to  $N$ ,

$$N_j(x_i) = 0 \quad \text{for } i=0 : j-1 \text{ and}$$

$$j=1 : N$$

$$N_{j+1}(x) = N_j(x) \cdot (x-x_j)$$

$N_0(x), N_1(x), \dots, N_n(x) \rightarrow$  Newton's Basis

$$p(x) = a_0 N_0(x) + a_1 N_1(x) + \dots + a_n N_n(x)$$

$$P(x_j) = f_j \quad \text{for } j=0:n \quad \text{yields -}$$

$$\left[ \begin{array}{cccc|c} N_0(x_0) & 0 & 0 & \dots & 0 \\ N_0(x_1) & N_1(x_1) & 0 & \dots & 0 \\ N_0(x_2) & N_1(x_2) & N_2(x_2) & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ N_0(x_n) & N_1(x_n) & N_2(x_n) & \dots & N_n(x_n) \end{array} \right] \left[ \begin{array}{c} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{array} \right] = \left[ \begin{array}{c} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_n \end{array} \right]$$

Remark -

- Solution of lower  $\Delta$  system requires  $O(n^2)$  flops
- Can accommodate new data  $(x_{n+1}, f_{n+1})$  with addition  $O(n)$  operations
- Evaluation of  $p(x)$  at a given  $x$  require  $O(n^2)$  operations.
- Computation of  $p(x)$  requires  $O(n^2)$  operations.

Computation of  $N_j(x_j)$  may be prone OVERFLOW / UNDERFLOW.

ASK Ask remark 2

↳  $N_j, j=0:n$  only depends on less than  $j$  nodes,

↳ Also only solving so last new eqn will take  $O(n)$ .

DIVIDED DIFFERENCE -

For example - consider  $(x_0, f_0), (x_1, f_1), (x_2, f_2)$

o Newton's Polynomial is -

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

$$f[x_j] := f_j \quad \text{for} \quad j=0, 1, 2$$

$$p(x_0) = f_0, \quad p(x_1) = f_1, \quad p(x_2) = f_2$$

$$\checkmark a_0 = f[x_0] = f_0$$

$$a_1 = f[x_1]$$

$$a_2 = f[x_2]$$

$$\checkmark a_1 = \frac{f[x_1] - f[x_0]}{x_1 - x_0} = f[x_0, x_1]$$

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

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

$$(x_2 - x_1)$$

$$= \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0} = f[x_0, x_1, x_2]$$

General case -

Divided difference can be generated

Using table of divided difference.

$x$	$f$		
$(x_0)$	$f_0$		$f[x_0, x_1, x_2] = ( ) - ( )$
$x_1$	$f_1$	$f[x_0, x_1]$	$x_2 - x_0$
$(x_2)$	$f_2$	$f[x_1, x_2] \leftarrow f[x_0, x_1, x_2]$	
$x_3$	$f_3$	$f[x_2, x_3]$	$f[x_0, x_1, x_2, x_3]$
⋮			go to diagonal extreme - left of $x$

$n+1$  diagonals  $\rightarrow$  coefficients of Newton Interpolating polynomial.

Adding new data  $(x_{n+1}, f_{n+1})$ .  $\rightarrow$  Adding a new row at the bottom of the table. Additional  $O(n)$  operations.

↪ Underflow / Overflow problem is solved.

Thm -

$$\text{Let } p(x) = a_0 N_0(x) + \dots + a_n N_n(x)$$

such that  $p(x_j) = f_j$  for  $j = 0 : n$ .

then,

$$a_j = f[x_0, x_1, \dots, x_j] \quad \text{for } j = 0 : n$$

Proof - Proof using induction.

The result is true for  $n=0$ .

Assume that the result is true for Newton interpolating polynomials of degree  $\leq N-1$ .

$$\text{Let } q(x) = \sum_{j=0}^{n-1} b_j N_j(x)$$

↳ for data set  $(x_0, f_0) \dots (x_n, f_n)$

$$\text{Let } s(x) = \sum_{j=0}^{n-1} c_j N_j(x)$$

↳ for dataset  $(x_0, f_0) \dots (x_{n-1}, f_{n-1})$

By induction hypothesis  $b_{n-1} = f[x_0, x_1, \dots, x_n]$   
 $c_{n-1} = f[x_0, x_1, \dots, x_{n-1}]$

$$r(x) =$$

For the dataset  $(x_0, f_0) \dots (x_n, f_n)$

$r(x)$  be interpolating polynomial.

$$r(x) = q(x) + \frac{x-x_n}{x_n-x_0} (q(x) - s(x))$$

→  $r(x)$  satisfies for data -  $(x_1, f_1), \dots, (x_n, f_n)$

$$\text{because } q(x) - s(x) = 0$$

$$r(x) = q(x).$$

→ For  $(x_0, f_0)$

$$r(x) = q(x_0) + (q(x_0) - s(x_0))$$

$$= s(x_0) = f_0$$

→ For  $(x_n, f_n)$

$$\begin{aligned} r(x_n) &= q(x_n) + 0 \\ &= f_n \end{aligned}$$

# Hence  $r(x)$  interpolates  $(x_0, f_0), \dots, (x_n, f_n)$

Now the coefficient of  $x^n$  in  $r(x)$

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

$$a_n = f[x_0, x_1, \dots, x_n]$$

## RUNGE

## PHENOMENON

consider the runge function

$f: [-1, 1] \rightarrow \mathbb{R}$  given by

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

Then for equally spaced nodes  $(x_0, \dots, x_n)$  and values  $f_j = f(x_j)$  for  $j=0:n$ .

The interpolant  $p_n(x)$  doesn't converge to  $f(x)$ . In fact  $\max_{|x| \leq 1} |f(x) - p_n(x)| \rightarrow \infty$  as  $n \rightarrow \infty$

But for Chebyshev Nodes  $x_j = \cos\left(\frac{(2j+1)\pi}{2n+2}\right)$

for  $j = 0 : n$  are  $\max_{|x| \leq 1} |f(x) - p_n(x)| \rightarrow 0$  as  $n \rightarrow \infty$ .

The runge phenomenon is eliminated by choosing Chebyshev nodes in interpolation points in  $[-1, 1]$

## CHEBYSHEV POLYNOMIAL

Let  $\theta \in [0, \pi]$

$x \in [-1, 1]$

Define  $T_n(x) = \cos(n \cos^{-1} x)$ .

$$T_0(x) = 1$$

$$T_1(x) = x$$

$$T_2(x) = 2x^2 - 1$$

$T_n(x)$  is  $n$ -degree polynomial and is called Chebyshev Polynomial.

$T_n(x)$  satisfies -  $T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x)$

Why? If  $x = \cos \theta$

$$\cos((n+1)\theta) + \cos((n-1)\theta) = 2\cos\theta \cos n\theta$$

$$\begin{bmatrix} T_1(x) \\ T_2(x) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 2x \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$\begin{bmatrix} T_n(x) \\ T_{n+1}(x) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 2x \end{bmatrix} \begin{bmatrix} T_{n-1}(x) \\ T_n(x) \end{bmatrix}$$

$$\therefore \begin{bmatrix} T_n(x) \\ T_{n+1}(x) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 2x \end{bmatrix}^n \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$T_n(x) = \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j} x^{n-2j} (x^2 - 1)^j$$

### Chebyshev NODES -

$$T_n(x) = \cos(n \cos^{-1} x)$$

we have  $|T_n(x)| \leq 1$  for  $x$

### CHEBYSHEV

**NODES**  $\Rightarrow T_n(x_j) = 0$  for  $x_j = \cos\left(\frac{(2j+1)\pi}{2n}\right)$

$j = 0, 1, \dots, n-1$

### GAUSS

**LOBATTO**  $\Rightarrow T_n(x_j) = (-1)^j$  for  $x_j = \cos\left(\frac{j\pi}{n}\right)$

**NODES**

$j = 0, 1, \dots, n-1$

### Gauß-Lobatto Nodes

Connection of Chebyshev nodes -

We want build  $p(x)$  of degree  $\leq n$  through  $n+1$  Chebyshev nodes.

The zeros of  $T_{n+1}(x)$  will give exactly  $n+1$  points in  $[-1, 1]$  i.e. - Chebyshev nodes will be zeros of  $T_{n+1}(x)$ .

Barycentric Lagrange Interpolation with Chebyshev Nodes

$$p(x) = \frac{\sum_{j=0}^n w_j f_j / |x - x_j|}{\sum_{j=0}^n w_j / |x - x_j|}$$

where  $w_j = \frac{1}{\prod_{i \neq j} (x_j - x_i)}$  for  $j = 0 : n$

For chebyshev nodes in  $[-1, 1]$  -

$$x_j = \cos\left(\frac{j\pi}{n}\right) \quad j = 0, 1, 2, \dots, n$$

$$w_j = \begin{cases} \frac{1}{2} (-1)^j & j = 0 \text{ or } j = n \\ (-1)^j & \text{otherwise} \end{cases}$$

$$p(x) = \frac{\sum_{j=0}^n (-1)^j f_j / (x - x_j)}{\sum_{j=0}^n w_j (-1)^j / (x - x_j)}$$

$\sum'$  means that terms  $j=0$  and  $j=n$  are multiplied by  $1/2$ .

NOTE - Barycentric interpolation formula remains valid for chebyshev nodes  $[a, b]$

## Approximation

Let  $C[a, b] := \{f : [a, b] \rightarrow \mathbb{R} \mid f \text{ is continuous}\}$

For  $f \in C[a, b]$ , define

$$\|f\|_{\infty} = \max \{ |f(x)| : x \in [a, b]\}$$

$$\textcircled{1} \quad \|f\|_{\infty} = 0 \iff f = 0$$

$$\textcircled{2} \quad \|\alpha f\|_{\infty} = |\alpha| \|f\|_{\infty}$$

$$\textcircled{3} \quad \|f + g\|_{\infty} \leq \|f\|_{\infty} + \|g\|_{\infty}$$

Weierstrass approx theorem:

Let  $f \in C[a, b]$  and  $\epsilon > 0$ . Then

there is a polynomial  $p(x)$  such that

$$\|f - p\|_{\infty} \leq \epsilon.$$

(On)

$$\max \{ |f(x) - p(x)| : x \in [a, b]\} \leq \epsilon$$

Question - Does  $p_n(x)$  approx  $f(x)$  for large enough  $n$ ? In other words, does  $\|p_n - f\|_{\infty} \rightarrow 0$  as  $n \rightarrow \infty$ ?

ANS - NO, for equispaced we saw  
 $n \rightarrow \infty$  the error was  $\infty$

## Interpolation error -

Let  $f \in C[a, b]$  and  $x_0, \dots, x_n$   
be distinct nodes in  $[a, b]$ ,  $f_j = f(x_j)$   
for  $j = 0 : n$ .

Lagrange polynomial -

$$p(x) = f_0 l_0(x) + f_1 l_1(x) + \dots + f_n l_n(x)$$

Define :

$$\{ \text{LEBESGUE FUNC} \} \quad \lambda_n(x) = |l_0(x)| + \dots + |l_n(x)|$$

$$\{ \text{LEBESGUE CONST} \} \quad \Lambda_n = \|\lambda_n\|_\infty$$

Set :

$$E_n(f) = \min \{ \|f - p\|_\infty : p \in P_n \}$$

$$\text{Proof 1} - \quad |P_n(x)| \leq \Lambda_n \|f\|_\infty$$

$$\text{Proof} - \quad |P_n(x)| = \left| \sum_{j=0}^n f_j l_j(x) \right|$$

$$\leq \sum_{j=0}^n |f_j| |l_j(x)|$$

$$\leq \|f\|_\infty \sum_{j=0}^n |l_j(x)|$$

$$\leq \|f\|_\infty \Lambda_n$$

$$\|P_n(x)\| \leq \|f\|_\infty \|\lambda_n\|_\infty \\ = \|f\|_\infty \Lambda_n$$

$$|P_n(x)| \leq \|f\|_\infty \Lambda_n \quad \left. \begin{array}{l} \text{Point wise also} \\ \text{holds true} \end{array} \right\}$$

$$\text{Proof 2} - \|f - P_n\|_{\infty} \leq (1 + \lambda_n) E_n(f)$$

Proof - Let  $p^* \in P_n$  be the best uniform approximant of  $f$ .

$$\|f - p^*\|_{\infty} = E_n(f)$$

Define :  $g = f - p^*$ ,  $g(x_j) = f_j - p^*(x_j)$

The polynomial  $P_n - p^*$  is the interpolant of  $(x_j, g(x_j))$

$$P_n(x) - p^*(x) = \sum_{j=0}^n g(x_j) \cdot l_j(x)$$

$$|P_n(x) - p^*(x)| \leq \sum_{j=0}^n |g(x_j)| \cdot |l_j(x)|$$

$$\leq \|g\|_{\infty} \sum_{j=0}^n |l_j(x)| = \|g\|_{\infty} \lambda_n(x) \quad \text{--- (1)}$$

$$= \|f - p^*\|_{\infty} \lambda_n(x)$$

$$|f(x) - P_n(x)| \leq |f(x) - p^*(x)| + |p^*(x) - P_n(x)| \leq \|f - p^*\|_{\infty} + \|f - p^*\|_{\infty} \lambda_n(x)$$

$$|f(x) - P_n(x)| \leq (1 + \lambda_n) \|f - p^*\|_{\infty}$$

$$\|f - P_n\|_{\infty} \leq (1 + \lambda_n) E_n(f)$$

$$\|f - P_n\|_{\infty} \leq (1 + \lambda_n) E_n(f)$$

→ Interpolation error is at most the best possible polynomial error multiplied by  $1 + \lambda_n$ . If  $\lambda_n$  is large → worse. Small  $\lambda_n \rightarrow$  better.

→  $\lambda_n$  depends only on nodes not on  $f$ . So node choice matters.

**FACT** - For equispaced nodes, the Runge function  $f(x) = \frac{2^n}{\pi \sin(\pi x)}$  and for Chebyshev  $\lambda_n \leq \frac{2 \log(n+1)}{\pi}$

## ERROR TERM FOR SMOOTH FUNCTION

Let  $C^n[a, b]$  denote the set of  $n$  times continuously differentiable functions in  $[a, b]$ .

Theorem -

If  $f \in C^n[a, b]$  and  $p_n(x)$  be unique polynomial of degree at most  $n$  passing through  $(x_0, f(x_0)), \dots, (x_n, f(x_n))$

$$\text{error } [f(x) - p_n(x)] = \frac{f^{(n+1)}(\theta x) \cdot w(x)}{(n+1)!}$$

for some  $\theta x \in [x_{\min}, x_{\max}]$  where  $x_{\min}$  and  $x_{\max}$  are the largest and the smallest nodes in  $[x_0, x_1, \dots, x_n, x]$  and  $w(x) = \prod_{i=0}^n (x - x_i)$

Proof -

$$\text{Define- } F(t) = f(t) - p(t) - \left( f(x) - p_n(x) \right) \frac{w(t)}{w(x)}$$

$$\left. \begin{array}{l} \text{At } t = x \Rightarrow F(t) = 0 \\ \text{At } t = x \Rightarrow F(x) = 0 \end{array} \right\} \begin{array}{l} F(t) \text{ has atleast} \\ n+2 \text{ distinct roots.} \end{array}$$

Applying Rolle's Thm repeatedly  $\Rightarrow F^{(n+1)}(t)$  has atleast one zero.

$$F^{(n+1)}(t) = f^{(n+1)}(t) - 0 - \left( f(x) - p_n(x) \right) \cdot \frac{d^{n+1}(w(t))}{dt}$$

$$= f^{(n+1)}(t) - \left( f(x) - p_n(x) \right) \cdot (n+1)!$$

At  $t = \theta_x$        $\theta_x \in [x_{\min}, x_{\max}]$

$$0 = f^{(n+1)}(\theta_x) - \left( \frac{f(x) - p_n(x)}{w(x)} \right) \cdot (n+1)!$$

so,

$$f(x) - p_n(x) = \frac{f^{(n+1)}(\theta_x)}{(n+1)!} \cdot w(x)$$

$$\text{so, } |f(x) - p_n(x)| \leq \|f^{(n+1)}\|_\infty \cdot |w(x)|$$

ERROR)

Error is  $w(x)$ -dependent so we want to minimize  $|w(x)|$

Chebyshev's Thm -

Goal - Choose nodes  $x_0, \dots, x_n$  in  $[a, b]$

that minimizes  $\max_{x \in [a, b]} \prod_{j=0}^n |x - x_j|$

$\min_{x_0, \dots, x_n} \max_{x \in [a, b]} \prod_{j=0}^n |x - x_j|$

ASK - If proof is mere or not.

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

Theorem - In  $x \in [-1, 1]$

$$\min_{x_0, x_1, \dots, x_n} \max_{x \in [-1, 1]} \prod_{j=0}^n |(x-x_j)| = 2^{-n}$$

and the minimum is attained when

$$w(x) = \prod_{j=0}^n (x-x_j) = 2^{-n} T_{n+1}(x)$$

Minimum is attained when  $x_0, \dots, x_n$  are Chebyshev nodes in  $[-1, 1]$ .

Hence,

$$|f(x) - P_n(x)| \leq \|f^{(n+1)}\|_\infty$$

$$(n+1)! 2^n$$

Chebyshev Nodes in  $[a, b]$

$$x_j = \frac{a+b}{2} + \frac{b-a}{2} \cos\left(\frac{2j+1}{2n+2} \pi\right) \quad j=0 : n$$

Change of interval -

$$\phi : x \rightarrow \frac{a+b}{2} + \frac{(b-a)}{2} x \text{ carries}$$

Chebyshev nodes in  $[-1, 1]$  to Chebyshev nodes in  $[a, b]$

The function  $\phi$  does two things -

Thus for chebyshev nodes in  $[1, 1]$  we have

$$\prod_{j=0}^n |(\phi(x) - \phi(x_j))| = \left(\frac{b-a}{2}\right)^{n+1} \prod_{j=0}^n |(x-x_j)|$$

$$\leq \left(\frac{b-a}{2}\right)^{n+1} \frac{1}{2^n}$$

### CHEBYSHEV INTERPOLATION ERROR

Thm - let  $f \in C^{n+1}[a, b]$  for chebyshev nodes

$$x_j = \frac{b+a}{2} + \frac{(b-a)}{2} \cos\left(\frac{2j+1}{2n+2} \cdot \pi\right) \text{ for } j=0:n$$

$$|f(x) - p_n(x)| \leq \frac{\|f^{(n+1)}\|_\infty \cdot \|w\|_\infty}{(n+1)!}$$

$$(n+1)! = (2n+2) \cdot (2n+1) \cdots 3 \cdot 2 \cdot 1 = n! \cdot (2n+1) \cdots 3 \cdot 2 \cdot 1 = n! \cdot \left(\frac{b-a}{2}\right)^{n+1} \cdot \frac{1}{2^n}$$

$$w(x) = (b-x) + x(a-b) = x - a$$

$[a, b]$  where weighted at  $[1, 1]$  in when weighted)

## DIVIDED DIFFERENCE AT REPEATED NODES

### HERMITE INTERPOLATION

PROBLEM -

We have distinct nodes -  $x_0, x_1, \dots, x_n$

We want our polynomial  $p(x)$  to match -

$f(x_j), f'(x_j), f''(x_j), \dots$

up to  $f^{(m_j)}(x_j)$

Here  $m_j$  is how many derivative conditions at  $x_j$

This means,  $f_0 -$

Find a polynomial  $p(x)$  of least degree such that  
for  $j=0:n$ .

$$P^{(j)}(x_j) = f^{(j)}(x_j) \quad \text{for } j=0, 1, \dots, m_j-1$$

### NOW ORDINARY NEWTON POLYNOMIAL SETUP

for  $N$  distinct nodes

$(x_0, f(x_0)), (x_1, f(x_1)), (x_2, f(x_2)), \dots$

$$P(x) = f[x_0] + f[x_0, x_1] \cdot (x-x_0) + f[x_0, x_1, x_2] \cdot \frac{(x-x_0)(x-x_1)}{(x-x_2)} + \dots$$

Now what happens if nodes are very close -

$(x_0, f(x_0))$  and  $(x_0+h, f(x_0+h))$

$$\lim_{h \rightarrow 0} f[x_0, x_0+h] = f'(x_0)$$

similarly

$$f[x_0, x_0+h, x_0+h_2] = \frac{f''(x_0)}{2!}$$

$$f[x_0, \dots, \underbrace{x_0}_{\text{n times}}] = \frac{f^n(x_0)}{n!}$$

### BACK TO HERMITE

so, we can consider repeated nodes of multiplicity  $m_j$  for  $j=0:n$ , that is

$$\underbrace{x_0, \dots, x_0}_{m_0}, \underbrace{x_1, \dots, x_1}_{m_1}, \dots, \underbrace{x_n, \dots, x_n}_{m_n}$$

$$N = m_0 + m_1 + \dots + m_n$$

Then the polynomial  $p(x)$  of degree  $N-1$  can be constructed from the divided difference table for repeated nodes.

for  $j=0:n$ , set  $f_j^{(i)} = f^{(i)}(x_j)$

$$i=0, 1, \dots, m_j - 1$$

For Example - If  $m_0 = 4$  then the divided difference upper rows will look like.

$x$	$f$			
$x_0$	$f_0$			
$x_0$	$f_0$	$f'_0$		
$x_0$	$f_0$	$f'_0$	$f''_0/2!$	
$x_0$	$f_0$	$f'_0$	$f''_0/2!$	$f'''_0/3!$

$$\frac{f'''_0}{3!}$$

Example - Determine a polynomial  $p(x)$  such that

$$p(1) = 2, \quad p'(1) = 3, \quad p(2) = 6, \quad p'(2) = 7, \quad p''(2) = 8$$

$x$	$f$
1	2
1	2 3
2	6 (4) (1)
2	6 7 (5) (2)
2	6 7 8/2! (1) (-1)

$$p(x) = 2 + 3(x-1)$$

$$= 2 + 3(x-1) + 1(x-1)^2 + 2(x-1)^2(x-2) - 1(x-2)^2(x-1)^2$$

# SPLINES

CLASSMATE

Date \_\_\_\_\_  
Page \_\_\_\_\_

High order polynomial  $\rightarrow$  Exhibit more oscillations  
We would like to avoid Runge phenomenon for large dataset.

To solve this

Strategy 1 - PIECE WISE LINEAR INTERPOLATION

Strategy 2 - Piece wise polynomial interpolation

Splines Definition -

A function  $s(x)$  is a spline of degree  $K$  on  $[a, b]$

If

$$\rightarrow s \in C^{K-1}[a, b]$$

$$\rightarrow a = t_0 < t_1 < \dots < t_n = b$$

$$s(x) = \begin{cases} s_0(x) & t_0 \leq x \leq t_1 \\ s_1(x) & t_1 \leq x \leq t_2 \\ \vdots & \vdots \\ s_{n-1}(x) & t_{n-1} \leq x \leq t_n \end{cases}$$

where  $s_i(x) \in P^K$

Cubic Splines

$$S(x) = \begin{cases} S_0(x) = a_0x^3 + b_0x^2 + c_0x + d_0 & t_0 \leq x \leq t_1 \\ \vdots \\ S_{n-1}(x) = a_{n-1}x^3 + b_{n-1}x^2 + c_{n-1}x + d_{n-1} & t_{n-1} \leq x \leq t_n \end{cases}$$

which satisfies

$$S(x) = C^2[t_0, t_n] : \left. \begin{array}{l} S_i(x_i) = s_i(x_i) \\ S'_i(x_i) = s'_i(x_i) \\ S''_i(x_i) = s''_i(x_i) \end{array} \right\} i = 1, 2, \dots, n-1$$

Now,

$$\text{Given } (x_i, y_i)_{i=0}^n.$$

Task - Find  $S(x)$  such that it is a cubic spline interpolant.

 This give me  $3n-3$  equations.

In addition we have

$$S(x_i) = y_i \quad \text{for } i=0, \dots, n$$

which give  $n+1$  equations.

Total we got =  $4n-2$  equations

we want =  $4n$  degree of freedom.

Thus we have 2 degree of freedom, left.

 Number of independent params you can still freely choose after satisfying me constraints.

For 2 degree of freedom -

We can use them to define different subtypes  
of cubic splines.

$$\rightarrow s''(t_0) = s''(t_n) = 0 \quad - \text{natural cubic spline}$$

$$\rightarrow s'(t_0), s'(t_n) \text{ given} \quad - \text{clamped cubic spline.}$$

$$\rightarrow \begin{aligned} s'''(t_0) &= s'''(t_1) \\ s'''_{n-2}(t_{n-1}) &= s'''_{n-1}(t_{n-1}) \end{aligned} \quad \left. \begin{array}{l} \text{NOT a knot (MATLAB)} \\ \text{condition} \end{array} \right\}$$

## NATURAL CUBIC SPLINE

Task - Find  $s(x)$  such that it is  
a natural cubic spline.

- Let  $t_i = x_i$ ,  $i = 0, \dots, n$
- Let  $z_i = s''(x_i)$ ,  $i = 0, \dots, n$ . This means the condition that it is a natural cubic spline is simply expressed as  $z_0 = z_n = 0$
- Since  $s(x)$  is third order polynomial, we know that  $s''(x)$  is a linear spline which interpolates  $(t_i, z_i)$ .
- Hence we can first construct the linear spline interpolant  $s''(x)$  and then integrate twice to obtain  $s(x)$ .

→ The linear spline -

$$s_i''(x) = z_i \cdot \left( \frac{x - t_{i+1}}{t_i - t_{i+1}} \right) + z_{i+1} \left( \frac{x - t_i}{t_{i+1} - t_i} \right)$$

$$\rightarrow h_i = t_{i+1} - t_i \quad \text{for } i = 0, \dots, n$$

$$s''(x) = z_{i+1} \left( \frac{x - t_i}{h_i} \right) + z_i \left( \frac{t_{i+1} - x}{h_i} \right)$$

→ We can now integrate here

$$s_i(x) = \frac{z_{i+1}}{6h_i} (x - t_i)^3 + \frac{z_i}{6h_i} (t_{i+1} - x)^3 \\ + C_i(x - t_i) + D_i(t_{i+1} - x) \quad \text{--- (1)}$$

$$s_i(t_i) = s_i(t_{i+1}) = \frac{z_i}{6} h_i^2 + D_i h_i = y_i, \quad i = 0, 1, \dots, n \quad \text{--- (2)}$$

Also for continuity,

$$s_i(t_{i+1}) = y_{i+1} \Rightarrow \frac{z_{i+1}}{6} \cdot h_i^2 + C_i h_i = y_{i+1} \quad \text{--- (3)}$$

using (1) (2) (3),

$$s_i(x) = \frac{z_{i+1}}{6h_i} (x - t_i)^3 + \frac{z_i}{6h_i} (t_{i+1} - x)^3 \\ + \left( \frac{y_{i+1}}{h_i} - \frac{z_{i+1}}{6} \cdot h_i \right) (x - t_i) \\ + \left( \frac{y_i}{h_i} - \frac{z_i}{6} \cdot h_i \right) (t_{i+1} - x)$$

$$S_i'(x) = \frac{z_{i+1}}{2h_i} (x - t_i)^2 - \frac{z_i}{2h_i} (t_{i+1} - x)^2$$

$$+ \underbrace{\frac{1}{h_i} (y_{i+1} - y_i)}_{b_i} - \frac{h_i}{6} (z_{i+1} - z_i)$$

$$S_i'(t_i) = -\frac{1}{2} z_i h_i + b_i - \frac{h_i}{6} z_{i+1} + \frac{1}{6} h_i z_i$$

$$S_i'(t_{i+1}) = \frac{z_{i+1} h_i + b_i - \frac{h_i}{6} z_{i+1} + \frac{1}{6} h_i z_i}{2}$$

SKIP

$$S_{i-1}(t_i) = \frac{1}{3} z_i h_{i+1} + \frac{1}{6} h_{i+1} z_{i+1} + b_{i-1}$$

FOLLOW SLIDES 13 FOR FURTHER

FOR OF SPLINE SLIDES.

# LEC-10

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Numerical Integration - 1

some integrals are very hard to evaluate. so we approximate them.

**Problem -** Let  $f \in C[a, b]$ , compute the integral  
 $I(f) = \int_a^b f(x) \cdot dx$ .

**Strategy -** Approximate via polynomial interpolation.

→ Find an interpolating polynomial  $p_n(x)$ .

→ Evaluate  $Q_n(f) = \int_a^b p_n(x) \cdot dx$

### QUADRATURE -

Consider the nodes  $[x_0, x_1, \dots, x_n]$  in  $[a, b]$  and the Lagrange interpolating polynomial  $p_n(x) = f(x_0) l_0(x) + \dots + f(x_n) l_n(x)$

$$l_j(x) = \frac{w_j(x)}{(x-x_j)(w_j(x_j))} = \prod_{i \neq j} \left( \frac{x_j - x_i}{x_i - x_j} \right)$$

$$w(x) = (x-x_0) \dots (x-x_n)$$

$$Q_n(f) = \int_a^b p_n(x) \cdot dx = \int_a^b \sum_{j=0}^n f(x_j) l_j(x) \cdot dx$$

$$= \sum_{j=0}^n f(x_j) \int_a^b l_j(x) \cdot dx = \sum_{j=0}^n f(x_j) w_j$$

Note that,  
 $w_j$  doesn't  
depend on  $f(x)$ .

where,  $w_j = \int_a^b l_j(x) \cdot dx$  is called  $j^{\text{th}}$  quadrature weight.

$Q_n(f)$  is quadrature formula or quadrature rule.

ERROR ESTIMATE -

$$E_n(f) = I(f) - Q_n(f) . \quad \text{If } f \in C^{n+1}[a, b]$$

$$|E_n(f)| \leq \int_a^b |f(x) - P_n(x)| \cdot dx$$

$$= \frac{1}{(n+1)!} \int_a^b |f^{(n+1)}(\xi_x) w(x)| \cdot dx$$

$$\leq \frac{\|f^{(n+1)}\|_{\infty}}{(n+1)!} \int_a^b |(x-x_0)(x-x_1) \dots (x-x_n)| \cdot dx$$

For chebyshev nodes  $\rightarrow E_n(f) \leq \left( \frac{b-a}{2} \right)^{n+2} \frac{\|f^{(n+1)}\|_{\infty}}{2^{n-1} (n+1)!}$

EXACTNESS -

A quadrature rule  $Q_n(f)$  is said to have degree of freedom exactness  $m$ . If  $Q_n(f) = I(f)$  for all  $f \in P_m$

The quadrature rule

$$Q_n(f) = \int_a^b P_n(x) \cdot dx \quad \text{where}$$

degree of exactness  $\geq n$

## NEWTON COTES QUADRATURE

Quadrature rule based on an interpolating polynomial  $p_n(x)$  at  $n+1$  equally spaced nodes in  $[a, b]$  is called Newton Cotes formula of order  $n$ .

Newton Cotes not useful for large  $n$ . Because  $\bar{E}_n(f)$  is susceptible to Runge's phenomenon.

Two ways to ensure  $\bar{E}_n(f) \rightarrow 0$  as  $n \rightarrow \infty$ :

- Don't use equally spaced nodes
- Use piece wise polynomial interpolant.

### MID POINT RULE

↪ The quadrature corresponding to  $p_0(x) = f\left(\frac{a+b}{2}\right)$

↪ Newton Cotes formula of degree order  $n=0$

↪ Polynomial of degree 0

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

Now mat

$$\text{If } f(x) = x$$

$$M(x) = (b-a) \left(\frac{a+b}{2}\right) = \frac{b^2 - a^2}{2} = \int_a^b x \cdot dx$$

which shows that it is exact for polynomial of degree 1.  
Thus mid point rule is exact for polynomial of degree 1.

## ERROR IN MIDPOINT RULE (Continue)

Theorem - If  $f \in C^2[a, b]$  then there exists  $\theta \in [a, b]$  such that

$$\int_a^b f(x) dx = M(f) + \frac{f''(\theta)}{24} (b-a)^3$$

$$\text{Hence } E(f) = M(f) - M(f) = \frac{f''(\theta)}{24} (b-a)^3$$

**Proof -** Set  $w = \left( \frac{a+b}{2} \right)$ ,

Taylor's form -

$$f(x) = f(w) + \underbrace{f'(w)(x-w) + f''(\theta_x) \frac{(x-w)^2}{2}}$$

↓  
constant

$$\int_a^b f(x) dx = (b-a) \cdot f(w) + 0 + \frac{f''(\theta_x)}{2} \int_a^b (x-w)^2 dx$$

↓  
if integrate  
and then put limits

$$= M(f) + \frac{f''(\theta_x)}{24} (b-a)^3$$

Error

But still

continue

$$\frac{x^2 - w^2}{2}$$

$$\frac{(b-a)^2}{n} = \left(\frac{b+a}{2}\right) \left(\frac{B_{12} + a}{2}\right)$$

classmate

Page

## COMPOSITE MIDPOINT RULE

Previous way for one single interval.

Now,

consider  $n$  equally spaced nodes  $[x_0, \dots, x_n]$  where  
 $x_j = a + jh$ ,  $j=0:n$  and  $h = \frac{b-a}{n}$ .

Set  $w_j = \frac{x_j + x_{j+1}}{2}$  for  $j=0$  to  $n-1$ .

Then,

$$I(f) = \int_a^b f(x) \cdot dx = \sum_{j=1}^n \int_{x_{j-1}}^{x_j} f(x) \cdot dx$$

By midpoint rule,

$$\int_{x_{j-1}}^{x_j} f(x) \cdot dx \approx h \cdot f(w_{j-1}) \cdot h \cdot f(w_{j-1})$$

so  $\text{composite midpoint rule-}$

$$M_n(f) = \sum_{j=1}^n h \cdot f(w_{j-1})$$

$$= h [f(w_0) + f(w_1) + \dots + f(w_{n-1})]$$

By midpoint rule,

$$\int_{x_{j-1}}^{x_j} f(x) \cdot dx - h \cdot f(w_{j-1}) = \frac{f''(c_j)}{24} \cdot h^3$$

ERROR IN TRAPEZOIDAL & COMPOSITE MIDPOINT RULE (contd.)

For Mid point rule

$$\int_{x_j}^{x_j} f(x) dx - h f(w_{j+1}) = \frac{f''(\theta_j)}{24} \cdot h^3$$

consequently, there exists  $\theta \in [a, b]$  such that

$$E_n(f) = I(f) - M_n(f)$$

$$= \sum_{j=1}^n \frac{f''(\theta_j)}{24} \cdot h^3$$

$$= n \cdot \frac{f''(\theta)}{24} \cdot h^3$$

$$= \frac{h^2}{24} f''(\theta) \cdot (b-a)$$

Hence,

$$|E_n(f)| \leq \frac{h^2}{24} \cdot (b-a) \cdot \|f''\|_\infty$$

## TRAPEZOID RULE

The quadrature  $Q_n(f)$  is called Trapezoid rule when  $n=1$ . Then nodes are  $x_0=a$  and  $x_1=b$ , and

$$l_0(x) = \frac{x-b}{a-b} \quad \text{and} \quad l_1(x) = \frac{x-a}{b-a}$$

Using Reference Quadrature (Lagrange) -

The weights are given by -

$$w_0 = \int_a^b l_0(x) \cdot dx = \frac{b-a}{2} = \int_a^b l_1(x) \cdot dx = w_1$$

Hence,

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

Theorem - If  $f \in C^2[a, b]$  then there exists  $\theta \in [a, b]$

such that  $\int_a^b f(x) \cdot dx = T(f) - \frac{f''(\theta)}{12} \cdot (b-a)^3$

$$\text{Hence } E(f) = I(f) - T(f) = - \frac{f''(\theta)}{12} \cdot (b-a)^3$$

Proof - By MVT (mean value Theorem) of integral,

$$E(f) = \int_a^b (f(x) - p(x)) \cdot dx$$

$$= \int_a^b \frac{f''(\xi_x)}{2} (x-a)(x-b) \cdot dx$$

$$\begin{aligned}
 &= \frac{f''(\theta)}{2} \int_a^b (x-a)(x-b) \cdot dx \quad \text{for some } \theta \in [a, b] \\
 &= -\frac{f''(\theta)}{12} \cdot (b-a)^3
 \end{aligned}$$

### Composite Trapezoid Rule

Consider equally spaced nodes  $[x_0, \dots, x_n]$   
 where  $x_j = a + jh$ ,  $j = 0 : n$  and  $h = \frac{(b-a)}{n}$

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

$$= \sum_{j=1}^n \int_{x_{j-1}}^{x_j} f(x) \cdot dx$$

By Trapezoid Rule,

$$\int_{x_{j-1}}^{x_j} f(x) \cdot dx = \frac{h}{2} [f(x_{j-1}) + f(x_j)]$$

Menu

$$I(f) = \int_a^b f(x) \cdot dx = \sum_{j=1}^n \int_{x_{j-1}}^{x_j} f(x) \cdot dx$$

$$\approx \sum_{j=1}^n \frac{h}{2} [f(x_{j-1}) + f(x_j)]$$

$T_n(f)$  $I_n(f)$  -

$$T_n(f) = h \left[ \frac{f(x_0)}{2} + f(x_1) + \dots + \frac{f(x_n)}{2} \right]$$

ERROR IN COMPOSITE TRAPEZOID RULE

By trapezoidal Rule  $\rightarrow -\frac{f''(\theta_j)}{12} h^3 = E(f)$

By intermediate value theorem -

$$\sum_{j=1}^n f''(\theta_j) = n f''(\theta) \quad \text{for } \theta \in [a, b]$$

Consequently,

$$E_n(f) = I(f) - T_n(f)$$

$$= - \sum_{j=1}^n \frac{f''(\theta_j)}{12} \cdot h^3$$

$$= -n \frac{f''(\theta)}{12} \cdot h^3 = -\frac{h^2}{12} (b-a) f''(\theta)$$

$$\text{Hence } |E_n(f)| \leq \frac{h^2}{12} (b-a) \|f''\|_\infty$$

Note that error from the Trapezoid Rule  
is almost twice that of the midpoint rule.

## SIMPSON'S RULE

The Newton-Cotes quadrature  $Q_n(f)$  is called Simpson's rule when  $n=2$ . The nodes are  $x_0=a$ ,  $x_1 = \frac{a+b}{2}$  and  $x_2=b$ . Then the Simpson's rule

is given by:

$$S(f) = w_0 f(a) + w_1 f\left(\frac{a+b}{2}\right) + w_2 f(b)$$

$$\text{where } w_i = \int_a^b l_i(x) dx \quad \text{for } i=0,1,2$$

Reference -

Newton Quadrature  
rule

After many calculations,

$$S(f) = \frac{(b-a)}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right]$$

Observe Note  
mat

$$S(f) = \frac{1}{3} T(f) + \frac{2}{3} M(f)$$

↓                      ↓  
Trapezoid Rule      Midpoint Rule

Observe mat

$$S(x^j) = \int_0^1 x^j dx \quad \text{for } j = 0, 1, 2, 3.$$

Hence  $S(f)$  is exact for polynomial of degree  $\leq 3$ .

## ERROR ESTIMATION OF SIMPSON'S RULE

Theorem - suppose that  $f \in C^4[a, b]$

$$\int_a^b f(x) dx = \frac{b-a}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] + E_2(f)$$

$$\text{where } E_2(f) = -\frac{1}{90} \left( \frac{b-a}{2} \right)^5 f^{(4)}(\theta) \text{ for some } \theta \in [a, b]$$

### COMPOSITE SIMPSON'S RULE

Suppose that  $n$  is even. Set  $h = \frac{b-a}{n}$

and consider the nodes  $x_j = a + jh$  for  $j = 0 : n$ .

Then,

$$\begin{aligned} \int_a^b f(x) dx &= \int_{x_0}^{x_2} f(x) dx + \int_{x_2}^{x_4} f(x) dx + \dots + \int_{x_{n-2}}^{x_n} f(x) dx \\ &= \sum_{j=1}^{\frac{n}{2}} \int_{x_{2j-2}}^{x_{2j}} f(x) dx \stackrel{n}{\approx} \sum_{j=1}^{\frac{n}{2}} \frac{h}{3} \left[ f(x_{2j-2}) + 4f(x_{2j-1}) + f(x_{2j}) \right] \\ &\approx \frac{h}{3} \left[ f(x_0) + 2 \sum_{j=2}^{\frac{n}{2}} f(x_{2j-2}) + 4 \sum_{j=1}^{\frac{n}{2}} f(x_{2j-1}) + f(x_n) \right] \end{aligned}$$

## ERROR IN COMPOSITE SIMPSON'S RULE

$$S_n(f) = \frac{h}{3} \left[ f(x_0) + 2 \sum_{j=1}^{n/2} f(x_{2j-2}) + 4 \sum_{j=1}^{n/2} f(x_{2j-1}) + f(x_n) \right]$$

with error  $E_n^S(f) = -\frac{1}{90} \cdot h^5 \cdot \sum_{j=1}^{n/2} f^{(4)}(\theta_j)$

for some  $\theta_j \in (x_{2j-2}, x_{2j})$

$$j = 1 : (n/2)$$

By IVT there exist  $\theta \in (a, b)$  such that

$$\sum_{j=1}^{n/2} f^{(4)}(\theta_j) = \frac{n}{2} f^{(4)}(\theta)$$

$$As nh = b-a$$

$$E_n^S f(n) = -\frac{(b-a)}{180} n^4 f^{(4)}(\theta)$$