

## Chapter 1 Introduction.

### 1.1 Examples? Option Pricing

Options: Simple examples are European Options.

Call, Payoff

$$p(T) = \max\{S - K, 0\}$$

American options allow early exercise. adds complexity.

when to exercise.

what value of the option today.  $V(S, 0)$

Black-Scholes Model.

1) stock price follows a geometric motion.

2) the risk-free interest rate  $r$  is constant.

3) All risk-free portfolios earn the risk-free rate. No arbitrage.

Gives a PDE (E 3.1)

$$\frac{\partial V}{\partial t} + \frac{\sigma^2 S^2}{2} \frac{\partial^2 V}{\partial S^2} + r \cdot S \cdot \frac{\partial V}{\partial S} - rV = 0$$

$V$  = value of  $S$ .

We buy options to hedge.

1) sell an option  $V$  to get cash

2) borrow  $S \cdot \frac{\partial V}{\partial S} - V$  from bank

3) buy  $\frac{\partial V}{\partial S}$  shares of stock.

Dynamic:

Hedge: buy & sell  $S$  to keep  $\frac{\partial V}{\partial S}$  amount of stock

so,  $\Delta = \frac{\partial V}{\partial S}$  is important.

For value of European options we have numerical

$$V(S, t) = S N(d_1) - K e^{-r(T-t)} N(d_2)$$

$$N(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{y^2}{2}} dy$$

$$d_{1,2} = \frac{\ln(S/K) + (r \pm \frac{1}{2} \sigma^2)(T-t)}{\sigma \sqrt{T-t}}$$

## 1.2 Numerical Issues

(1) The PDE cannot be solved by analytical means except a few cases. Approximation is needed. Sect. II b-d.

(2) Even if an analytical solution exists it might need to be approximated. eg.  $N(x)$  must be approximated. Sect II c.

(3) Real binomial method.

$$V_j^{n+1} = e^{-r\Delta t} [p V_{j+1}^n + (1-p) V_j^n]$$

$$S_j = j \cdot \Delta S$$

$$t_n = n \cdot \Delta t$$

to hedge, need  $\frac{\partial V}{\partial S}$  from discrete data II c.

(4) The volatility is an inferred quantity.

$$V(S, t) - S N(d+\delta) - K e^{-r(T-t)} N(d-\delta) = 0$$

(5) Data is given discreteness, not as functions. We approximate data by smooth functions. II b.

(6) Sometimes we must simulate process to get solutions.

eg. Monte Carlo Methods. II d.

## Character-2 : Basic Numerical Methods

### 2.1 Errors and Conditions.

#### 2.1.1 Floating point numbers

EX: 

```
float sum = 0.0f
for (int i=1 ; i<=10 ; i++)
{
    sum += 0.1f
}
cout << sum - 1.0f
```

$$\text{Result : } \sum_{i=1}^{10} 0.1 - 1 = 1.192093 \times 10^{-7}$$

Computer approximation real numbers.

A real number  $X$ .

$$X = \pm \sum_{k=1}^{\infty} \frac{d_k}{B^k} B^e$$

$B \rightarrow$  base

$d_k \rightarrow$  digits  $0 \leq d_k \leq B-1$

$$\text{ex: } 63.58 = 0.6358 \times 10^2 = 0.06358 \times 10^3$$

↑  
takes place of an exponent.

so, we required  $d_1 \neq 0$ , called normalized.

a real number is called normalized, if it is in the form:

$$\pm d_0 d_1 d_2 d_3 \dots \times 10^n$$

required  $d_0 \neq 0$ .

Converting a number to base 2 and normalizing it at the former step in storing a real number as a floating-point number in a computer.

Computers approximate  $x$  with floating-point number.

$$Fl(x) \approx x$$

$$Fl(x) = \pm s \cdot \beta^e$$

$$s \rightarrow \text{significate} \quad s = \sum_{k=1}^t \frac{\hat{d}_k}{\beta^k} \quad \text{finite sum}$$

$\beta \rightarrow \text{base}$

$$e \rightarrow \text{exponent} \quad 0 \leq \hat{d}_k \leq \beta - 1$$

normalized F.p. #,  $\hat{d}_1 \neq 0$

$$L \leq e \leq u$$

F.p. # have different properties from real #'s.

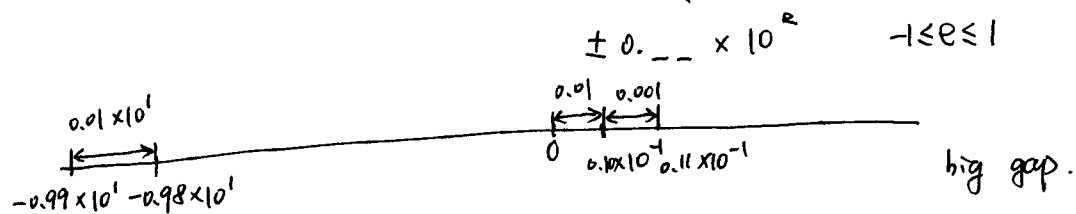
1) There is a finite # of F.p. #'s

$$Fl(x) = \pm \sum_{k=1}^t \frac{\hat{d}_k}{\beta^k} \beta^e$$

$\swarrow$                        $\downarrow$                        $\swarrow$   
 $2$                        $(\beta-1)\beta^{t-1}$                        $u+|L|+1$

$$\# = 2 \times (\beta-1) \beta^{t-1} (u+|L|+1)$$

2) They are discrete. eg:  $\beta=10, L=-1, u=1, t=2$



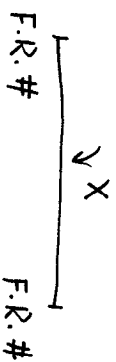
3) They are not uniformly distributed.

4) There is an upper and lower value.

eg.  $|x| > 9.9$  is as good as  $\infty$ .

5) They don't get arbitrarily close to 0.

How do we approximate  $x$  by  $fl(x)$ ?

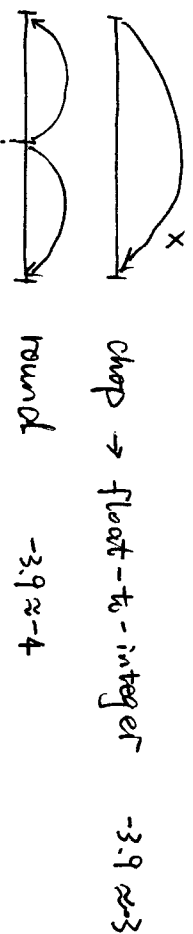


we Round

1) chop  $\pm N6$

2) round  $\pm N6$

$$fl(x) = \pm \sum_{k=1}^{\infty} \frac{d_k}{2^k} \cdot 2^e$$



Today, all machines use IEEE 754 standard.

Defines 2 precisions.

1) single precision (float in C++)

$$E_{min} = -126$$

$$E_{max} = 127$$

$$fl(x) = (-1)^s 2^{e-127} 1.M$$

$$0 < e < 255, \quad M = \sum_{k=1}^{23} \frac{d_k}{2^k}$$

$$B=2$$

$$S \rightarrow 1 \text{ bit}, \quad e \rightarrow 8 \text{ bits}$$

$$fl(x) \rightarrow 32 \text{ bits} \Rightarrow M = 23 \text{ bits}$$

2) double precision

$$E_{min} = -1022$$

$$E_{max} = 1023$$

$$fl(x) = (-1)^s 2^{e-1023} 1.M$$

$$e \rightarrow 11 \text{ bits}, \quad M = 52 \text{ bits}$$

$$\text{Total} \rightarrow 64 \text{ bits}$$

$$\text{option } 2.2204e-16 = 2^{-52}$$

maximum representable value is

$$2^{127} (2 - 2^{-23}) \approx 3.4 \times 10^{38}$$

smallest representable value is  $2^{-106} = 1.17549e^{-38}$

$$E_{s.p} = 2^{-23} = 1.19209e^{-7}$$

The largest spacing

$$2^{127} (2 - 2^{-23}) - 2^{127} (2 - 2^{-23} - 2^{-23}) = 2^{104}$$

The smallest

$$2^{-126} (1 + 2^{-23}) - 2^{-126} = 2^{-149}$$

$$\text{largest } 1.79769e308 = 2^{1023} (2 - 2^{-52})$$

$$\text{smallest } 2^{-1022} \approx 2.2251e^{-308}$$

The extra value of e flag exceptions

① NaN. (Not a number) eg.  $\frac{0}{0}$ ,  $\sqrt{-1}$ ,  
 $E=255$ ,  $M \neq 0$  (S.P.)

② INF (Infinity) eg.  $\frac{1}{0}$ ,  $E=255$ ,  $M=0$  (S.P.)

③ denormalization.

$$\cancel{Fl(x) = \pm 0.M E} \quad \cancel{Fl(x) = \pm 0.M 2^E}$$

$$Fl(x) = \pm 0.M 2^{E-27} \quad \text{S.P. } E=0$$

3.14159265357928...

$$\beta = 10 \quad 0.314159 \times 10^1 / 10$$

$$t = 6 \quad 0.314159 \times 10^0 / 10$$

$$u = 4 = -L \quad 0.314159 \times 10^{-1} / 10$$

$$0.314159 \times 10^{-2} / 10$$

$$0.314159 \times 10^{-3} / 10$$

$$0.314159 \times 10^{-4} \rightarrow \text{the largest negative denormali}$$

$$0.314159 \times 10^{-4} / 10 \quad \text{since } u = 4 = -L.$$

$$0.031416 \times 10^{-4} / 10$$

$$0.003142 \times 10^{-4}$$

### 2.1.2 Measuring Errors.

Define. Let  $p^*$  approximates  $P$ .

$$\text{Absolute Error.} \quad E_{abs} = p^* - P$$

$$\text{Relative Error.} \quad E_{rel} = E_{abs} / P \quad P \neq 0$$

$$\text{Eg: } P = 0.1, \quad p^* = 0.2$$

$$E_{abs} = p^* - P = 0.1$$

$$E_{rel} = \frac{E_{abs}}{P} = \frac{0.1}{0.1} = 1$$

$E_{rel}$  is more important, it tells us the significant of error.

Define:  $p^*$  approximates  $p$  to  $M$

significant digits (sig. figures, sig. figs)

$$\text{if } \bar{E}_{rel} \leq 0.5 \times 10^{-M} = 5 \times 10^{-(M+1)}$$

$$p = \pi, \quad p^* = 3.1415$$

$$|E_{abs}| = 9.2 \times 10^{-5}$$

$$|E_{rel}| = \frac{9.2 \times 10^{-5}}{\pi} \approx 3 \times 10^{-5} < 5 \times 10^{-5} \Rightarrow m = 4$$

$$\bar{E}_{rel} \leq 0.5 \times 10^{-M} = 5 \times 10^{-(M+1)}$$

$$Fl(x) = (1 + \varepsilon) x$$

$$\Rightarrow \left| \frac{Fl(x) - x}{x} \right| = |\varepsilon|$$

$\uparrow$   
 $E_{rel}$

what is  $\varepsilon$ ?

$$x = \sum_{k=1}^{\infty} \frac{d_k}{\beta^k} \cdot \beta^e$$

$$Fl(x) = \sum_{k=1}^t \frac{d_k}{\beta^k} \beta^e \quad \text{chopped arithmetic}$$

$$-\varepsilon = \frac{\sum_{k=1}^{\infty} \frac{d_k}{\beta^k} \beta^e - \sum_{k=1}^t \frac{d_k}{\beta^k} \cdot \beta^e}{\sum_{k=1}^{\infty} \frac{d_k}{\beta^k} \beta^e}$$

$$= \frac{\sum_{k=t+1}^{\infty} \frac{d_k}{\beta^k} \beta^e}{\sum_{k=1}^{\infty} \frac{d_k}{\beta^k} \cdot \beta^e} \leq \frac{\frac{1}{\beta^{t+1}} \sum_{k=0}^{\infty} \frac{d_{k+t+1}}{\beta^k}}{\frac{1}{\beta}}$$

$$\leq \frac{1}{\beta^{t+1}} \cdot \beta / \frac{1}{\beta} = \beta^{1-t}$$

$$|\varepsilon| \leq \beta^{1-t} = \epsilon = \text{machine epsilon (chop)}$$

$$\epsilon_{round} = \frac{1}{2} \epsilon_{chop} = \frac{1}{2} \beta^{1-t}$$

IEEE defines  $\epsilon$  as the unit round: Find  $\epsilon$ ,  $1 + \epsilon > 1$

eg. S.P.  $1.000 \dots 0$  23 bits  
 $0.000 \dots 1 = 2^{-23} = 2^{-t}$   
 $\uparrow$  23rd place  
 $1.000 \dots 1 > 1$

D.P.  $\epsilon = 2^{-52}$

$Fl(x) = (-1)^s 1.M \times 2^{e-127}$

significant figures

IEEE S.P.  $\epsilon_{s.p} = 2^{-23} = 1.192 \times 10^{-7}$   $m=6$

$\epsilon_{D.P} = 2^{-52} = 2.2 \times 10^{-16}$   $m=15$

2.1.3 Conditioning.

EX:  $p(x) = (x-1)(x-2) \dots (x-20) = x^{20} - 210x^{19} + \dots$  Roots?

Suppose we perturb 210 to  $210 + 2^{-23}$

How do the root change?

$E_{rel, \pm N} = 2^{-23}/210 = 5.7 \times 10^{-10}$   $m=8$   
 $= 0.57 \times 10^{-9}$

One root.

$x=20$  changes to  $x^* = 20.846908101$

$E_{abs, out} = 0.846908101$

$E_{rel, out} = 0.846908101/20 = 0.042345405$

$= 4.2345405 \times 10^{-2}$   $m=1$

lost 7 digits.

amplification  $\frac{E_{rel, out}}{E_{rel, in}} = \frac{0.04}{6 \times 10^{-10}} = 8 \times 10^{-7}$

Error amplified by 80 million

Example of an ill-conditioned / sensitive problem



what is an ill-conditioned problem?

when  $\frac{E_{rel, out}}{E_{rel, in}} \gg 1$

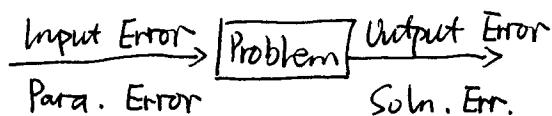
Float class  $F$ ;

FM Float  $x$ ;

FM Float  $EPS$ ;

$EPS = F \cdot \text{epsilon}(x)$

Conditioning.



Relative size of the errors gives the application

$$\left| \frac{E_{rel, solution}}{E_{rel, input}} \right|$$

if this is large, the problem is ill-conditioned.

Quantify conditions.

Suppose a solution  $X$  is a function of a parameter  $\alpha$

$$X = x(\alpha)$$

let  $\Delta\alpha = a$  in  $\alpha$

$$\Delta X = x(\alpha + \Delta\alpha) - x(\alpha)$$

Then relative change in input is  $\frac{\Delta\alpha}{\alpha}$   
output  $\frac{\Delta X}{X}$

Then 
$$\frac{\frac{\Delta X}{X}}{\frac{\Delta \alpha}{\alpha}} = \frac{\alpha}{X} \cdot \frac{\Delta X}{\Delta \alpha}$$

let  $\Delta \alpha \rightarrow 0$

$$\lim_{\Delta \alpha \rightarrow 0} \left| \frac{\frac{\Delta X}{X}}{\frac{\Delta \alpha}{\alpha}} \right| = \left| \frac{\alpha}{X} \cdot \frac{\Delta X}{\Delta \alpha} \right|$$

And define the condition number

$$K = \max \left| \frac{\alpha}{X} \cdot \frac{\Delta X}{\Delta \alpha} \right|$$

And a problem is ill-conditioned if  $K \gg 1$

what is  $\gg 1$ ?

Answer: It depends.

Note: 
$$\frac{\frac{\Delta X}{X}}{\frac{\Delta \alpha}{\alpha}} \sim K \alpha$$

In context

$$\frac{\Delta X}{X} \rightarrow \text{Erel, solution}$$

$$\frac{\Delta \alpha}{\alpha} \rightarrow \text{Erel, in}$$

$$\frac{\text{Erel, soln}}{\text{Erel, in}} \sim K \alpha$$

$$\bar{\text{Erel, soln}} \sim K \text{Erel, in}$$

If  $\text{Erel, in}$  is due to rounding

$$\bar{\text{Erel, in}} \sim \underset{\substack{\uparrow \\ \text{epsilon}}}{\epsilon}$$

$$\bar{\text{Erel, soln}} \sim K \cdot \epsilon$$

$$\text{then } \log(\bar{\text{Erel, soln}}) \sim \log(K) + \log(\epsilon)$$

Then in significant digits if

$$\log(K) + \log(\epsilon) \leq \log 0.5 - m$$

$$m \leq \log 0.5 - \log(K) - \log(\epsilon)$$

$\uparrow$                        $\underbrace{\hspace{2cm}}_{\geq 0 \text{ fixed}}$

For each of 10 in  $K$ , we lose 1 significant digit.

What is big? Depends on how many sig. digits needed.

eg.  $K=10^4$

$$\epsilon = 10^{-7} \text{ \& } \log(0.5) - \log(\epsilon) \approx 6$$

$$\text{then } m \sim 6 - 4 = 2$$

If you need  $m=4$ , then  $K=10^4$  is too big.

On the other hand. (OTOH).

if we need 15 sig. digits in input (eg IEEE D.P)

$$\text{then } m = 15 - 4 = 11 \text{ OK.}$$

EX:

$$x = e^\alpha, \quad K = \left| \frac{\alpha}{x} \cdot \frac{\Delta x}{\Delta \alpha} \right| = \left| \frac{\alpha}{e^\alpha} \cdot e^\alpha \right| = |\alpha|$$

Suppose  $\alpha = 5.5$

lose at most 1 sig. digit.

we can also do systems  $\vec{x} = \vec{x}(\vec{\alpha})$

$$\text{then } x(\vec{\alpha} + \Delta \vec{\alpha}) = x(\vec{\alpha}) + \underbrace{\left( \frac{\partial \vec{x}}{\partial \vec{\alpha}} \right)}_{\text{Jacobian}} \Delta \vec{\alpha}$$

$$\text{so } \Delta \vec{x} = J(\vec{\alpha}) \Delta \vec{\alpha}$$

$$\|\Delta \vec{x}\| = \|J \Delta \vec{\alpha}\| \leq \|J\| \cdot \|\Delta \vec{\alpha}\|$$

$$\frac{\|\Delta \vec{x}\|}{\|\vec{x}\|} \leq \frac{\|J\|}{\|\vec{x}\|} \cdot \left( \frac{\|\Delta \vec{\alpha}\|}{\|\vec{\alpha}\|} \right) \cdot \|\vec{\alpha}\|$$

relative error in parameters

$$\frac{\frac{\|\Delta \vec{x}\|}{\|\vec{x}\|}}{\frac{\|\Delta \vec{\alpha}\|}{\|\vec{\alpha}\|}} \leq \frac{\|J\| \cdot \|\vec{\alpha}\|}{\|\vec{x}\|}$$

$$\text{Refine. } K = \max \frac{\|J\| \|\vec{\alpha}\|}{\|\vec{x}\|}$$

Linear system are often ill-conditioned.

$$Ax = y$$

actually, we solve  $(A + \Delta A)x^* = y + \Delta y$

$$\text{OR } Ax^* = y + \Delta y - \Delta Ax^* \equiv y^*$$

$$\Leftrightarrow Ax^* = y^*$$

$$K = \max \frac{\|J\| \|\vec{y}^*\|}{\|x\|}$$

what is  $J$ ?

$$J = \frac{\partial \vec{x}^*}{\partial \vec{y}^*}$$

$$\vec{x}^* = A^{-1} \vec{y}^*$$

$$\frac{\partial \vec{x}^*}{\partial \vec{y}^*} = A^{-1}$$

$$\text{also, } y^* = Ax^*$$

1.20

$$Ax^* = y^*$$

$$\frac{\frac{\|\Delta x^*\|}{\|x^*\|}}{\frac{\|\Delta y^*\|}{\|y^*\|}} = \frac{\|J\| \|\vec{y}^*\|}{\|\vec{x}^*\|}$$

$$J = \frac{\partial x^*}{\partial y^*} = A^{-1}$$

$$\frac{\|A^{-1}\| \|Ax^*\|}{\|\vec{x}^*\|} \leq \frac{\|A^{-1}\| \cdot \|A\| \cdot \|\vec{x}^*\|}{\|x^*\|}$$

$$K(A) = \|A^{-1}\| \cdot \|A\|$$

EX:

$$A = \begin{vmatrix} 1 & 1.01 \\ 0.99 & 1 \end{vmatrix}$$

$$\text{Norm: } \|A\|_{\infty} = \max_i \sum_{j=1}^n \|a_{ij}\|$$

$$\text{Here } \|A\|_{\infty} = \max \{2.01, 1.99\} = 2.01$$

$$A^{-1} = \frac{1}{1 - 0.99 \times 1.01} \begin{vmatrix} 1 & -1.01 \\ -0.99 & 1 \end{vmatrix}$$

$\parallel$   
 $10^4$

$$\|A^{-1}\|_{\infty} = 10^4 \max \{2.01, 1.99\} = 2.01 \times 10^4$$

solution in single precision will have only 2 sig. figures

$$K(A) = 2.01 \times 2.01 \times 10^4 = 4 \times 10^4$$

#### 2.1.4 Stability of Algorithm

EX: Compute  $e^{-5.5}$  is this ill-conditioned?

$$K = 5.5 \quad \text{Not}$$

$$e^{-x} = 1 - x + \frac{x^2}{2} - \frac{x^3}{3!} + \dots \quad \text{Taylor Series}$$

choose  $\beta=10, t=5$

$$e^{-5.5} = 1.0000 - 5.5000 + 15.250 - 27.730 + 38.129 - \dots$$

↑ round      ↑ round

$$= 0.0026363 = 0.26363 \times 10^{-2}$$

↑

50% error, No sig. digits.

exact: 0.0040868.

solution is on size of rounding errors.

$$\text{Fixed? } \frac{1}{e^{5.5}} = \frac{1}{1.0000 - 5.5000 + \dots} = 0.0040865.$$

characteristic of algorithm.

11) subtraction of nearly equal numbers (catastrophic cancellation).

ex: exact arithmetic

$$\begin{array}{r} 151.72899 \\ - 151.71422 \\ \hline 0.01477 \end{array} \quad \text{exact}$$

approx.  $\beta=10, t=\frac{5}{5}$

$$\begin{array}{r} 151.73 \\ - 151.71 \\ \hline 0.020000 \end{array}$$

$$E_{rel} = \frac{E_{abs}}{p} = \frac{p^* - p}{p}$$

$$E_{rel} = \frac{0.02 - 0.01477}{0.01477} = 0.35 < 0.5 \times 10^0$$

↓

no sig. figures.

Rule #1 avoid catastrophic cancellation  
we often must be clever to this.

Probabilities:

(a) analytically rearrange the forward to remove subtractions.

$$\text{ex: } e^{-5.5} = \frac{1}{1 - 5.5 + \frac{5.5^2}{2} - \dots}$$

$$\begin{aligned} \text{ex: } y &= x(\sqrt{x+1} - \sqrt{x}) \\ &= x(\sqrt{x+1} - \sqrt{x}) \left( \frac{\sqrt{x+1} + \sqrt{x}}{\sqrt{x+1} + \sqrt{x}} \right) = \frac{x}{\sqrt{x+1} + \sqrt{x}} \end{aligned}$$

given different answers.

(b) make approximations, so the subtraction can be done analytically.

(2) avoid multiplying errors by large factors.

(happens in recursive algorithm)

EX:  $Ax = y$ , usually solved by Guess-Elimination

$$\begin{vmatrix} 10^{-10} & 1 \\ 1 & 0 \end{vmatrix} = \begin{vmatrix} x \\ y \end{vmatrix} = \begin{vmatrix} 1+10^{-10} \\ 1 \end{vmatrix} \quad \text{IEEE S.P.} \\ \text{exact} = \begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

approximation

$$\begin{vmatrix} 10^{-10} & 1 \\ 1 & 0 \end{vmatrix} \begin{vmatrix} x \\ y \end{vmatrix} = \begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\begin{vmatrix} 10^{-10} & 1 \\ 0 & -10^{-10} \end{vmatrix} \begin{vmatrix} x \\ y \end{vmatrix} = \begin{vmatrix} 1 \\ -10^{-10} \end{vmatrix} \Rightarrow \begin{aligned} 1-10^{10} &= -(10^{10}-1) = -10^{10} \\ y=1, x=0 &\leftarrow 10^{-10}x + 1 = 1 \end{aligned}$$

Fix: Re-order so multiplier < 1

$$\begin{vmatrix} 1 & 0 \\ 10^{-10} & 1 \end{vmatrix} \begin{vmatrix} x \\ y \end{vmatrix} = \begin{vmatrix} 1 \\ 1+10^{-10} \end{vmatrix}$$

approximation  $\begin{vmatrix} 1 & 0 \\ 10^{-10} & 1 \end{vmatrix} = \begin{vmatrix} 1 \\ 1 \end{vmatrix}$

$$\begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} \begin{vmatrix} x \\ y \end{vmatrix} = \begin{vmatrix} 1 \\ 1-10^{-10}=1 \end{vmatrix} \quad \begin{aligned} x &= 1 \\ y &= 1 \end{aligned}$$

called Guess-Elimination with partial pivoting

B) Avoid Range Errors.

EX:  $r = \sqrt{x^2 + y^2}$

$$x = 10^{20}, y = 10^{20} \Rightarrow r = \sqrt{2} \cdot 10^{20}$$

S.P. IEEE  $r = 1 \cdot N \cdot F \quad (10^{20})^2 = 10^{40} > 10^{38} \quad \text{Infinity}$

FIX: IF  $|y| > |x|$ ,  $r = |y| * \sqrt{1 + (\frac{x}{y})^2}$

else  $r = |x| * \sqrt{1 + (\frac{y}{x})^2}$

## 2.2 Solution of nonlinear algebraic equations

EX:  $V$  = value of an account

$P$  = amount deposited periodically

$i$  = interest / period

$$V = \frac{P}{i} ((1+i)^N - 1)$$

Q: what is the min. interest rate needed to have an amount  $V$  in an account after  $N$  periods if we invest  $P$ .

Find  $i$  given  $V, N, P$ .

can't do directly

we can write as a root finding problem. Find  $x \rightarrow F(x) = 0$

eg.  $V \rightarrow \frac{P}{i} ((1+i)^N - 1) = 0$

Two methods

(1) Bisection

6 totally convergent but slow.

(2) Newton's Method

Fast convergent, but not always convergent.

Bisection: Binary Search

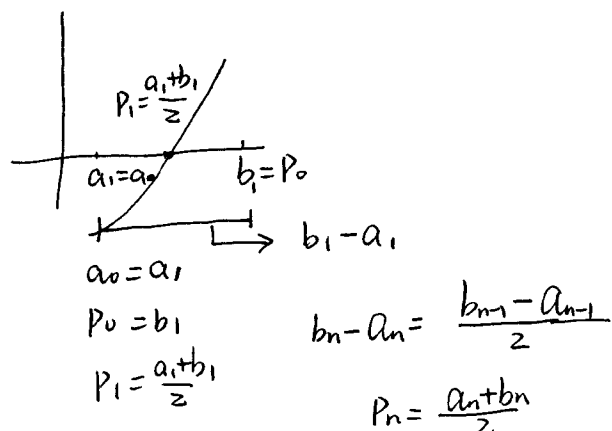
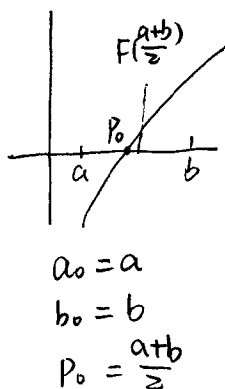
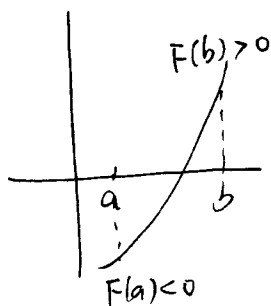
Theorem: (Intermediate value theorem)

If  $I = [a, b]$ ,  $F(x) \in C(I)$ . and  $F(a)F(b) < 0$

then  $F(x) = 0$  at some point  $P \in [a, b]$



To bisection.



$$|P_n - P| \leq \frac{b-a}{2^{n+1}} \quad n \geq 0$$

$$|P - P_n| \leq \frac{1}{2} (b_n - a_n)$$

REM: 1) clearly  $\lim_{n \rightarrow \infty} |P_n - P| = 0$

2) at each step, the error bound decreases by a factor of 2.

$$b_n - a_n = \frac{b_{n-1} - a_{n-1}}{2}$$

$$P_n = \frac{a_n + b_n}{2}$$

Proof:

$$b_n - a_n = \frac{b_{n-1} - a_{n-1}}{2} = \frac{b_{n-2} - a_{n-2}}{2^2} = \dots = \frac{b_{n-1} - a_{n-1}}{2^{n-1}} = \frac{b-a}{2^n}$$

$$|P - P_n| \leq \frac{1}{2} (b_n - a_n) = \frac{b-a}{2^{n+1}}$$

$$|P - P_n| \leq \frac{1}{2} (b_n - a_n) = \frac{1}{2} \cdot \frac{b_{n-1} - a_{n-1}}{2} = \frac{1}{4} \cdot \frac{b_{n-2} - a_{n-2}}{2} = \dots = \frac{b-a}{2^{n+1}}$$

An algorithm:

$$P_n = P_{n-1} + \frac{b_{n-1} - a_{n-1}}{2}$$

For  $i=1$  to  $N$ , Do  $P = \frac{a+b}{2}$

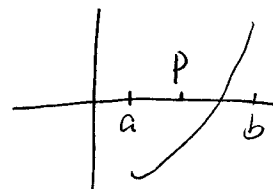
if converged exit

if  $F(a) F(p) > 0$ ,  $a = p$

else  $b = p$

endif

Next  $i$ .



$$b_n - a_n = \frac{b_{n-1} - a_{n-1}}{2}$$

$$P_n = \frac{b_n + a_n}{2}$$

$$|P - P_n| \leq \frac{b-a}{2^{n+1}}$$

Are there floating point #s? 1

How do we test convergence. 2

What is N? 3

## 1. F.P. Problems

(a) Find  $F(a)$   $F(p)$  can overflow or underflow

$$\text{Test } \text{sign}(F(a)) * \text{sign}(F(p)) > 0$$

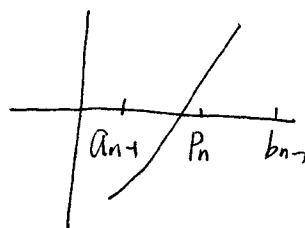
(b)  $p = \frac{a+b}{2}$

is not necessarily in  $[a, b]$

eg.  $B=10, t=3, \frac{0.981+0.982}{2} = \frac{1.964}{2} \overset{\uparrow \text{3 digits}}{=} 0.980 \neq 0.982$

normally rewrite as a correction

$$p_n = p_{n-1} + \frac{b_{n-1} - a_{n-1}}{2}$$



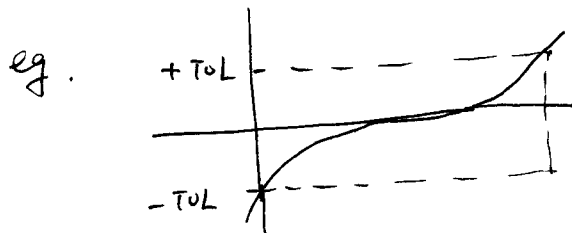
## 2. Test for convergence.

(1)  $F(p_n) \leq \text{Tol}$  Root test

(2)  $|p - p_n| \leq \text{Tol}$  absolute error test Tol  $\rightarrow$  tolerance

(3)  $\frac{|p - p_n|}{|p|} \leq \text{Tol}$  Relative error test

Root test Problems?



where is root?

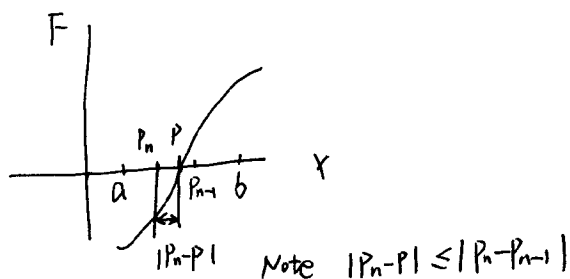
It can be anywhere

in   

where  $dF/dx = 0$

$$K = \left| \frac{\frac{dF}{dx}}{F} \right| / \left| \frac{\Delta x}{x} \right| = \left| \frac{x}{F} \cdot \frac{dF}{dx} \right| \quad F \approx 0$$

estimate error



estimate / bound on error

Test:  $|P_n - P_{n-1}| \leq \text{absolute tolerance} + \text{relative tolerance} \times |P_n|$   $|P_n - P| \leq |P_n - P_{n-1}|$   
 when  $|P_n| \approx 0$ , we have abs. err. test  $|P_n - P_{n-1}| \leq \text{absTol} + \text{relTol} \times |P_n|$   
 when  $|P_n| \gg 0$ , we have rel. err. test

what should tolerance be?

$$\epsilon \leq \text{relative tolerance} \leq 5 \times 10^{-(m+1)}$$

 $m \rightarrow$  desired # sig. figures $\epsilon \rightarrow$  machine epsilon

$$\epsilon \leq \text{relTol} \leq 5 \times 10^{-(m+1)}$$

absolute tolerance  $\approx$  smallest distance between 2 #s.absTol  $\approx$  smallest distance between 2 numbers

$$2^{-149} \approx 10^{-38}$$

3. what is  $N$ ?

$$\text{Recall } |P_n - P| \leq \frac{b-a}{2^{n+1}}$$

$$|P_n - P| \leq |P_n - P_{n-1}|$$

$$\text{so we converge if } |P_n - P_{n-1}| \leq \frac{b-a}{2^{n+1}}$$

$$\text{and again if } \frac{b-a}{2^{n+1}} \leq \text{absTol} + \text{relTol} \times |P_n|$$

$$\text{i.e. if } \frac{b-a}{\text{absTol} + \text{relTol} \times |P_n|} \leq 2^{n+1}$$

$$\text{or } n+1 \geq \log_2 \frac{b-a}{\text{absTol} + \text{relTol} \times |P_n|}$$

safe if  $N = \log_2 \frac{b-a}{\text{abstol} + \text{reltol} * |\frac{a+b}{2}|}$

$$N_{it} = \frac{\log \frac{b-a}{\text{abstol} + \text{reltol} * |\frac{a+b}{2}|}}{\log 2}$$

for  $i=1$  to  $N_{it}$

$$p = a + \frac{b-a}{2}$$

if  $\frac{b-a}{2} \leq \text{abstol} + \text{reltol} * |p|$  exit

if  $\text{sign}(F(a)) * \text{sign}(F(p)) > 0$

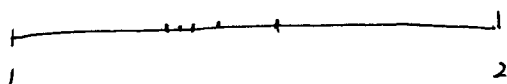
$$a = p$$

else

$$b = p$$

endif

next  $i$ .

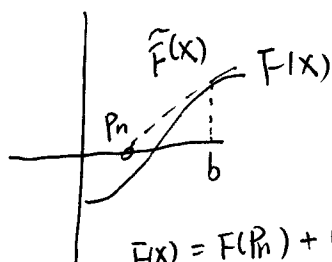


$$\frac{1}{2} \cdot \frac{1}{2} \cdots \frac{1}{2^n} = \frac{1}{2^{23}} \Rightarrow n=23$$

1.27

$$F(x) = 0$$

Newton's Method

Find roots of successive linear approximation of  $F$ 

$$F(x) = F(P_n) + (x - P_n)F'(P_n) + \frac{(x - P_n)^2}{2!} F''(\xi)$$

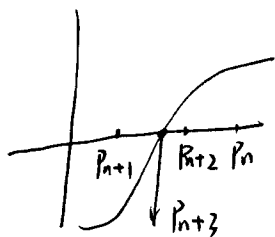
$$\tilde{F} = F(P_n) + (x - P_n)F'(P_n)$$

Find root of  $\tilde{F}$ , call it  $P_{n+1}$ 

$$F(P_n) + (P_{n+1} - P_n)F'(P_n) = 0$$

$$\text{Solve for } P_{n+1} = P_n - \frac{F(P_n)}{F'(P_n)}$$

Graphically

Given  $P = P_0 \rightarrow$  what's initial?For  $i = 1$  to  $N \rightarrow$  what's  $N$ 

$$P = P - \frac{F(P)}{F'(P)} \rightarrow \text{what if } F'(P) = 0$$

If converged exit

Next  $i$ .  $\uparrow$  How do you tell?

Theorem: let  $F \in C^2[a, b]$ , suppose  $F(p) = 0$ , for some  $p \in [a, b]$  and assume  $F'(p) \neq 0$  for  $x \in [a, b]$ , then  $\exists \delta > 0 \Rightarrow \{P_n\}_{n=0}^{\infty}$  converge to  $p$  for any  $P_0 \in [p - \delta, p + \delta]$

REM: For smooth enough  $F$ , Newton's Method converges if the initial guess is "close enough"

Proof: Let  $g(x) = x - \frac{F(x)}{F'(x)}$

Then  $P_{n+1} = g(P_n)$

Note  $P = g(P)$  ( $P$  is a fixed point of  $g$ )

$$\text{Now } g'(x) = 1 - \frac{F(x)F''}{(F')^2} - 1 = \frac{F(x)F''(x)}{(F')^2}$$

so that  $g'(P) = 0$ ,  $g(x)$  is continuous.

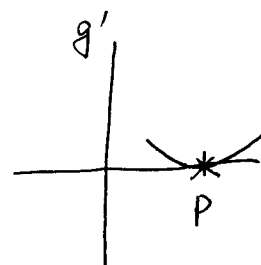
so there is an interval  $x \in [P-\delta, P+\delta]$

$$\exists g'(x) < 1$$

$\downarrow$   
so that

$$\text{i.e. } |g'| \leq k < 1, \quad \forall |x-P| \leq \delta$$

$\downarrow$   
for all



$P$  is root

in the neighborhood of  $P$   
 $g$  has to close to 0.

$$\begin{aligned} \text{Next } |g(x) - P| &= |g(x) - g(P)| \\ &= |g'(c)| |x - P| \\ &\leq k |x - P| \end{aligned}$$

choose  $x = P_n$ ,  $\Phi \quad P_{n+1} = g(P_n)$

$$\begin{aligned} \text{so } |P_{n+1} - P| &\leq k |P_n - P| \\ &\leq k^2 |P_{n-1} - P| \leq \dots \leq k^{n+1} |P_0 - P| \end{aligned}$$

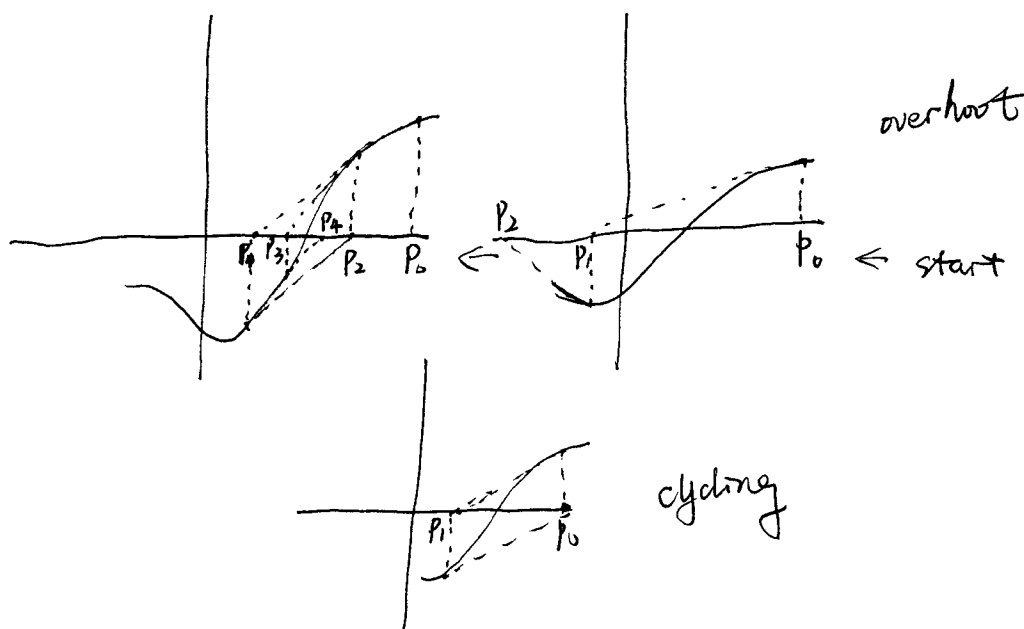
$$\begin{aligned} \exists \text{ since } k < 1, \quad \lim_{n \rightarrow \infty} |P_{n+1}| &\equiv \lim_{n \rightarrow \infty} |P_{n+1} - P| \\ &\leq \lim_{n \rightarrow \infty} k^{n+1} |P_n - P| = 0 \end{aligned}$$

so  $P_n \rightarrow P$  as  $n \rightarrow \infty$

So what is "close enough"

$$\frac{F(x) F''(x)}{(F'(x))^2} < 1$$

in practice we don't test this  
otherwise it can fail.



1.29

$$F(x) = 0$$

Given  $P = P_1$

For  $i = 1$  to  $N$

$$P = P - \frac{F(P)}{F'(P)}$$

If converge exit

Next  $i$ .

Sometimes this is modified so that  $P$  doesn't change too much in each iteration:

Newton Method

change:  $\Delta = - \frac{F(P)}{F'(P)}$  i.e.  $P = P + \Delta \Leftrightarrow P \leftarrow P + \Delta$

so required  $\frac{|\Delta|}{|p|} \leq \beta \rightarrow \text{some value}$

and update  $\beta$   
 $\Delta \leftarrow \beta |p| \text{sign}(\Delta)$

algorithm becomes

Given  $P_0$

For  $i = 1$  to  $N$

$$\Delta = F(p) / F'(p)$$

if  $|\Delta| > \beta |p|$

$$\Delta = \beta |p| \text{sign}(\Delta)$$

Endif

$$p = p + \Delta$$

if converged exit

next  $i$ .

$$\beta \sim 0.10$$

$$0.5 \times 10^{-1} = 0.05$$

Definition:

Suppose  $\lim_{n \rightarrow \infty} p_n = p$ , Let  $e_n = p_n - p$

the sequence  $\{p_n\}_{n=0}^{\infty}$  converges with order  $\alpha$  if

$$\lim_{n \rightarrow \infty} \frac{|p_{n+1} - p|}{|p_n - p|^\alpha} = \lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|^\alpha} = \lambda \rightarrow \text{constant value}$$

REM: if  $\alpha = 1 \rightarrow$  linear convergence

if  $\alpha > 1 \rightarrow$  superlinear convergence

if  $\alpha = 2 \rightarrow$  quadratic

the bound in bisection is linear  $\lambda = \frac{1}{2}$



Theorem: let  $p$  be a fixed point of  $x = g(x)$ .

suppose  $g'(p) = 0$  and  $g''$  is continuous in some open interval about  $p$ . Then

$$\exists \delta > 0 \rightarrow \text{for } p_0 \in [p - \delta, p + \delta]$$

then sequence  $\{p_n\}_{n=0}^{\infty}$  generated by  $p_{n+1} = g(p_n)$  converges quadratically.

Proof: we've already shown convergence

$$\text{write } g(x) = g(p) + g'(p)(x-p) + \frac{1}{2}g''(\xi)(x-p)^2$$

but  $g'(p) = 0$ , so

$$g(x) = \cancel{g(p)} + p + \frac{1}{2}g''(\xi)(x-p)^2$$

choose  $x = p_n$ , then  $g(p_n) = p_{n+1}$

$$p_{n+1} = p + \frac{1}{2}g''(\xi)(x-p)^2$$

$$e_{n+1} = \frac{1}{2}g''(\xi)(x-p)^2 = \frac{1}{2}g''(\xi)e_n^2$$

$$\text{so } \lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|} = \lim_{n \rightarrow \infty} \frac{1}{2}|g''(\xi)| = \frac{1}{2}|g''(p)| = \lambda$$

$\Rightarrow$  quadratic convergence.

Corollary: Newton's method converges quadratically.

$$\text{Proof: } g(x) = x - \frac{f(x)}{f'(x)}$$

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

$$= \frac{f(x)f''(x)}{(f'(x))^2} \quad \text{so } g'(p) = 0 \quad ???$$

Look at error

$$|e_{n+1}| \rightarrow \lambda |e_n|^2 \text{ as } n \rightarrow \infty$$

$$\log |e_{n+1}| \rightarrow 2 \log |e_n| + \log \lambda$$

similarity  $\log \frac{|e_{n+1}|}{|p|} \rightarrow 2 \log \frac{|e_n|}{|p|} + \log |p\lambda|$

→ exponent of the error doubles

suppose we start with  $\left|\frac{e_0}{p}\right| \approx 10^{-1}$

$$\left|\frac{e_1}{p}\right| \approx \cancel{10^{-1}} 10^{-2}$$

D.P. 4 iterations

$$\left|\frac{e_2}{p}\right| \approx \cancel{10^{-2}} 10^{-4}$$

S.P. 1 iteration

$$\left|\frac{e_3}{p}\right| \approx \cancel{10^{-4}} 10^{-8}$$

$$\left|\frac{e_4}{p}\right| \approx 10^{-16}$$

→ why stop here?  
machine  $\epsilon$

when is it converge?

$$|\Delta| = |p_{n+1} - p_n|$$

Note  $\Delta = p_{n+1} - p_n = (p_{n+1} - p) - (p_n - p) = e_{n+1} - e_n$

But  $|e_{n+1}| \rightarrow \lambda |e_n|^2 \ll \lambda |e_n|$

so  $|\Delta| \approx |e_n|$

converge test

$$|\Delta| \leq \text{abstol} + |p_{n+1}| \text{reltol}$$

$$N = 6 \sim 10$$

to get  $p_0$  use bisection

To A'-rel rel error of  $0.5 \times 10^{-1}$

2.1

## 2.3 Interpolation

Motivation: Implied/Implied volatility Black-Scholes model assumes constant volatility in asset prices & gives equation for an option price.

$$V_t + \frac{r^2 S^2}{2} \frac{\partial^2 V}{\partial S^2} + rS \frac{\partial V}{\partial S} - rV = 0$$

$r \rightarrow$  volatility

But  $\sigma$  is not measurable. What's published are  $s, t$ , and  $V$ .

At discrete times &  $s$ , so  $\sigma$  is inferred, and it's not constant.

Implied volatility problems:

Find  $\sigma$  for each  $s, t$  &  $V$  published.

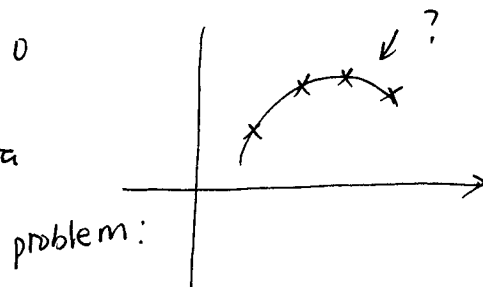
$$\underset{\substack{\downarrow \\ \text{published}}}{V(s, t)} - V_{BS}(s, t, \sigma) = 0$$

This gives  $\sigma(s_i, t_j)$ , ie. discrete values.

Then they compute

$$V_t + \frac{\sigma^2(s, t)}{2} \frac{\partial^2 V}{\partial S^2} + rS \frac{\partial V}{\partial S} - rV = 0$$

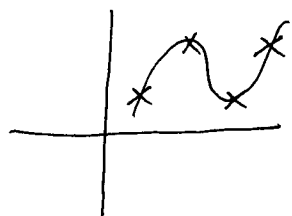
need a function form for data



How do we find a function form  
to discretize the data?

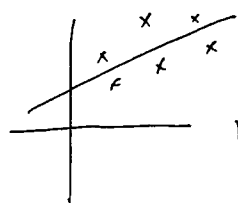
### 2.3.2 Polynomial Approximation two fundamental problems

- 1) Data is assumed to be exact. The approximation must match data.



~~Inter~~ Interpolation

- 2) Data is assumed to have errors. Try to find the best "approx." of a given form.



regression/least squares

we usually use polynomials

$$P_N(x) = \sum_{n=0}^N a_n x^n$$

Interpolation:  $Y_i = P_N(x_i) \quad x_i = 0, 1, \dots, N$

least squares: minimize  $\sum_{i=0}^N (Y_i - P_N(x_i))^2$

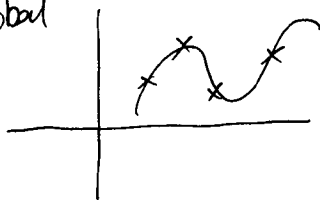
But we can also use other functional forms, eg for periodic problems

$$P_N(x) = \sum_{n=-N}^N a_n e^{inx}$$

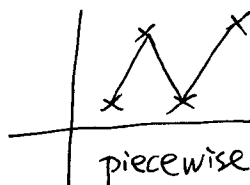
fourier series

The approx. can be global or local.

global



local



piecewise linear

note: P.W approx.'s don't have to be smoothness constraints include splines.

Smooth P.W approx's with additional

### 2.3.3 Global Polynomial Interpolation

29

Given data  $\{(x_i, y_i)\}_{i=0}^N$

The Polynomial interpolant  $P_N(x)$  satisfies

$$y_i = P_N(x_i) \quad i=0, 1, \dots, N$$

ie. with  $P_N = \sum_{n=0}^N a_n x^n$ ,  $y_i = \sum_{n=0}^N a_n (x_i^n)$   $i=0, 1, \dots, N$

↑  
min.

$N+1$  unknown

$N+1$  equations

or  $\vec{y} = M \vec{a}$

$$\vec{y} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \vec{a} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_N \end{bmatrix} \quad M_{in} = x_i^n$$

so  $\vec{a} = M^{-1} \vec{y}$

unfortunately, (for  $x_i$ 's uniformly spaced)

$$k(M) = 10^N$$

$$N \sim 3 \text{ to } 4$$

Moral: Don't do this!

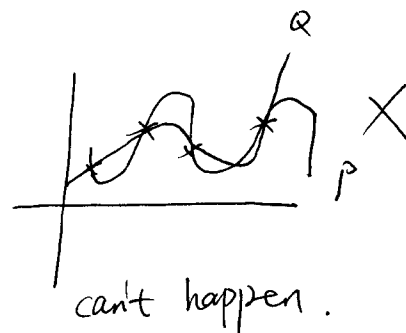
THM: Uniqueness.

Let  $P \neq Q$  be polynomials of degree  $\leq N$ .

If the nodes  $\{x_i\}_{i=0}^N$  are distinct

$$\text{then } p(x_i) = q(x_i) \quad i=0, 1, \dots, N$$

$$\text{Then } p = q, \quad \forall x$$



Proof: assume  $p \neq 0$  and define  $H(x) = p - q$ .

then  $H(x)$  is polynomial of degree  $\leq N$  and

$$H(x_i) = 0 \quad i=0, 1, \dots, N$$

$N+1$  roots, polynomial degree  $N$ .  
 the only happen if  $H(x)=0$ ,  $p=q$

We will use two forms of  $P_N$  to find polynomial interpolant

1) Lagrange

2) Newton

Lagrange: The Lagrange is

$$P_N(x) = \sum_{i=0}^N \gamma_i L_i(x)$$

↑ Lagrange interpolant polynomials

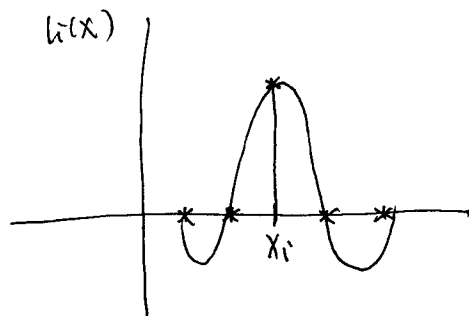
where  $L_i(x)$  satisfies

1)  $L_i(x)$  polynomials of degree  $N$ .

$$L_i(x_j) = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

$$= \delta_{ij}$$

How to construct them?



$$L_i(x) = \underbrace{(x-x_0)(x-x_1)\cdots(x-x_{i-1})(x-x_{i+1})\cdots(x-x_N)}_{N \text{ factors,}} = (x-x_0)(x-x_1)\cdots(x-x_{i-1})(x-x_{i+1})\cdots(x-x_N)$$

roots are all except  $x=x_i$ .

$$L_i(x) = \frac{(x-x_0)(x-x_1)\cdots(x-x_{i-1})(x-x_{i+1})\cdots(x-x_N)}{(x_i-x_0)(x_i-x_1)\cdots(x_i-x_{i-1})(x_i-x_{i+1})\cdots(x_i-x_N)}$$

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^N \frac{x-x_j}{x_i-x_j} \quad P_N(x) = \sum_{i=0}^N \gamma_i \left( \prod_{\substack{j=0 \\ j \neq i}}^N \frac{x-x_j}{x_i-x_j} \right)$$

EX: Find line through (1,1) (3,5)

$$P_1(x) = 11 \frac{(x-3)}{(1-3)} + 15 \frac{(x-1)}{(3-1)}$$

2.3

$$P_N = \sum_{n=0}^N a_n x^n$$

$$P_N = \sum_{j=0}^N Y_j L_j(x)$$

Indeed,  $P_N(x_i) = Y_j$ ,  $j = 0, 1, \dots, N$

$$L_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^N \frac{(x - x_i)}{(x_j - x_i)}$$

EX:

$x_j$	$Y_j$
2	6
3	11
4	18

$$P_2(x) = \frac{6(x-3)(x-4)}{(2-3)(2-4)} + \frac{11(x-2)(x-4)}{(3-2)(3-4)} + \frac{18(x-2)(x-3)}{(4-2)(4-3)}$$

Algorithm:

- Interpolate (x)
- $P = 0$
- For  $j = 0$  to  $N$ 
  - $P_i = P + Y_j * LIP(j, x)$
- next  $j$
- return  $p$ .

work:  $N \times$  to get  $p$

$N \times$  to get each  $L$

$LIP(j, x)$

$L = 1$

For  $j = 0$  to  $j-1$

$$L = (x - x_i) / (x_j - x_i) * L$$

next  $i$ .

For  $i = j+1$  to  $N$

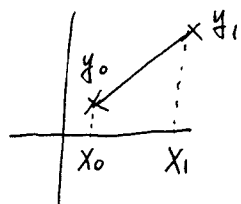
$$L = (x - x_i) / (x_j - x_i) * L$$

overall work =  $O(N^2)$

$\propto N^2$  operations.

Newton:

point slope form



$$P_1(x) = y_0 + m(x - x_0)$$

$$m = \frac{y_1 - y_0}{x_1 - x_0} \equiv y[x_0, x_1]$$

↑

called the first divided difference

$$y[x_0, x_1] = \frac{y_1 - y_0}{x_1 - x_0}$$

$$P_1(x) = y_0 + y[x_0, x_1](x - x_0)$$

identical to  $P_1(x) = y_0 \cdot \frac{x - x_1}{x_0 - x_1} + y_1 \cdot \frac{x - x_0}{x_1 - x_0}$

$$P_1(x) = y_0 + \frac{y_1 - y_0}{x_1 - x_0} (x - x_0)$$

We get higher order interpolations recursively.

DEF:  $y[x_0, x_1, \dots, x_N] \equiv \frac{y[x_1, x_2, \dots, x_N] - y[x_0, x_1, x_2, \dots, x_{N-1}]}{x_N - x_0}$

assume  $y[x_j] = y_j$

$$y[x_0, x_1, \dots, x_N] = \frac{y[x_1, x_2, \dots, x_N] - y[x_0, x_1, \dots, x_{N-1}]}{x_N - x_0}$$

EX:  $x_0, x_1, x_2,$

$$y[x_0, x_1] = \frac{y_1 - y_0}{x_1 - x_0}$$

$$y[x_1, x_2] = \frac{y_2 - y_1}{x_2 - x_1}$$

$$y[x_0, x_1, x_2] = \frac{\frac{y_2 - y_1}{x_2 - x_1} - \frac{y_1 - y_0}{x_1 - x_0}}{x_2 - x_0}$$

And we define the Newton form interpolant

$$P_N(x) = P_{N-1} + y[x_0, x_1, \dots, x_N] \prod_{i=0}^{N-1} (x - x_i)$$

$$P_N(x) = P_{N-1} + y[x_0, x_1, \dots, x_N] \prod_{i=0}^{N-1} (x - x_i)$$



Two steps:

(1) Compute divided differences

(2) Evaluate  $P_N(x)$

$$P_N(x) = P_{N-1} + y[x_0, \dots, x_N] \prod_{i=0}^{N-1} (x-x_i) \quad 33$$

EX:

j	$x_j$	$y_j$	$y[x_j, x_{j-1}]$	$y[x_0, x_1, x_2]$
0	2	6	$\frac{11-6}{3-2} = 5$	$\frac{7-5}{4-2} = 1$
1	3	11	$\frac{18-11}{4-3} = 7$	
2	4	18		

$$P_2(x) = 6 + 5(x-2) + 1 \cdot (x-2)(x-3)$$

$$P_2(2) = 6$$

$$P_2(3) = 11$$

$$P_2(4) = 6 + 10 + 2 = 18$$

$$P_N(x) = y_0 + y[x_0, x_1](x-x_0) + y[x_0, x_1, x_2](x-x_0)(x-x_1) + \dots + y[x_0, x_1, \dots, x_N] \prod_{i=0}^{N-1} (x-x_i)$$

Remarks:

(1) Order of the points is irrelevant (irrelevant).

(2) Note:

$$P_N = y_0 + y[x_0, x_1](x-x_0) + y[x_0, x_1, x_2](x-x_0)(x-x_1) + \dots + y[x_0, x_1, \dots, x_N] \prod_{i=0}^{N-1} (x-x_i)$$

Taylor Polynomial

$$f_N(x) = y(x_0) + y'(x_0)(x-x_0) + \frac{y''(x_0)}{2!}(x-x_0)^2 + \dots + \frac{y^{(N)}(x_0)}{N!}(x-x_0)^N$$

Thus

$$P_N(x) = y(x_0) + y'(x_0)(x-x_0) + \frac{y''(x_0)}{2!}(x-x_0)^2 + \dots + \frac{y^{(N)}(x_0)}{N!}(x-x_0)^N$$

(3) Recall  $y[x_0, x_1] \equiv \frac{y(x_1) - y(x_0)}{x_1 - x_0} = y'(\xi) \quad \text{MUT.}$

and so

(§)  $y[x_0, x_1, \dots, x_N] = \frac{y^{(N)}(\xi)}{N!}$

To compute,

$$P_N(x) = \sum_{i=0}^N d_i w_i(x)$$

$$w_i(x) = \prod_{j=0}^{N-1} (x - x_j)$$

$d_i \rightarrow$  divided factor

differece table  $(\vec{x}, \vec{y})$

```

For i = 0 to N
    di = yi
next i
For j = 1 to N
    For i = N to j step -1
        di = (di - di-1) / (xi - xi-j)
    next i
next j
Return d

```

Evaluate(x)

```

P = 0 d0
w = 1
For i = 1 to N
    w = w * (x - xi-1)
    P = P + di * w
next i
Return P

```

{ Newton  
Lagrange

How well does  $P_N$  approximate  $y(x)$ ?

THM: Let  $\{x_i\}_{i=0}^N$  be distinct on an interval  $[a, b]$  and let  $y \in C^{N+1}[a, b]$ ,

for each  $x \in [a, b]$ ,  $\exists \xi \in [a, b] \Rightarrow y(x) = P_N(x) + \frac{y^{(N+1)}(\xi)}{(N+1)!} \prod_{j=0}^N (x - x_j)$

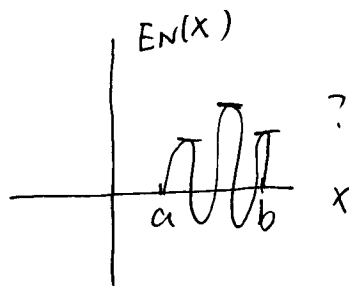
Note: The interpolation error

$$y(x) = P_N(x) + \underbrace{\frac{y^{(N+1)}(\xi)}{(N+1)!} \prod_{j=0}^N (x - x_j)}_{E_N(x)}$$

1)  $E_N(x_i) = 0$ ,  $i = 0, 1, \dots, N$

2) If  $y$  is a polynomial of degree  $\leq N$ ,  $E_N(x) = 0$ .

3)  $E_N(x)$  is a polynomial of degree  $N+1 \Rightarrow E_N$  has  $N$  extrema  
 $\Rightarrow E_N$  is very oscillatory.



Example:

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

$E_N$  is oscillatory

$$\max_x |E_N| \approx \max_x |y| \quad \max_x |E_N| \approx \max_x |y|$$

(gets worse as  $N \rightarrow \infty$ )

Note also max. error is near end points (commonly called "Runge PHENOMENON")

Consequences:

1) best to evaluate only near the center of the interval.

2) extrapolation is bad.

Fixes?

Yes, one thing we can do is to choose a better set of points.

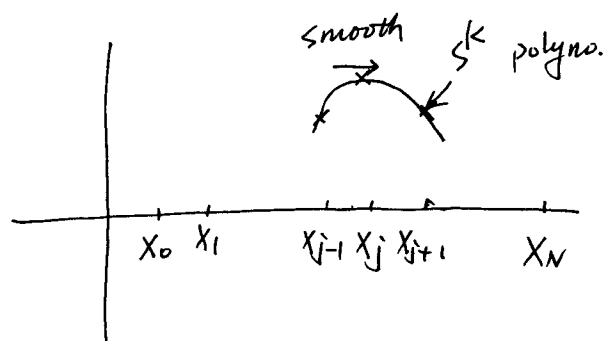
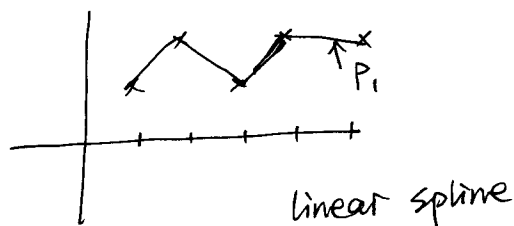
we will use piecewise lower order polynomials.

### 2.3.4. spline Interpolation

Definition: Let  $\{x_i\}_{i=0}^N$  be order and the elements distinct. The function  $s^k$  is a polynomial spline function of degree  $k$  if

- (1) on each subinterval  $[x_i, x_{i+1}]$ ,  $s^k$  is a polynomial of degree  $\leq k$
- (2)  $s$  has  $(k-1)$  continuous derivatives on the interval  $[x_0, x_N]$

EX: piecewise linear interpolation



Most common spline is  $s^3$  (cubic spline)

ie.  $s^3$  is cubic on  $[x_i, x_{i+1}]$

$s', s''$  are continuous at the  $x_i$ 's.

Derivation:

Note  $s''$  is piecewise linear

$$s''_i(x) = M_i \frac{x - x_{i+1}}{x_i - x_{i+1}} + M_{i+1} \frac{x - x_i}{x_{i+1} - x_i}$$

$\downarrow$   
 $s''$  on  $[x_i, x_{i+1}]$

let  $\Delta x_i = x_{i+1} - x_i$

$$s''_i = \frac{M_i (x_{i+1} - x)}{\Delta x_i} + M_{i+1} \cdot \frac{x - x_i}{\Delta x_i} = M_i \frac{x_{i+1} - x}{\Delta x_i} + M_{i+1} \cdot \frac{x - x_i}{\Delta x_i}$$

$$\nabla M_i = s''_i(x_i) \approx y''(x_i)$$

$$s''_i = M_i \frac{x_{i+1} - x}{\Delta x_i} + M_{i+1} \frac{x - x_i}{\Delta x_i}$$

$$M_i \approx s''_i(x_i) \approx y''(x_i)$$

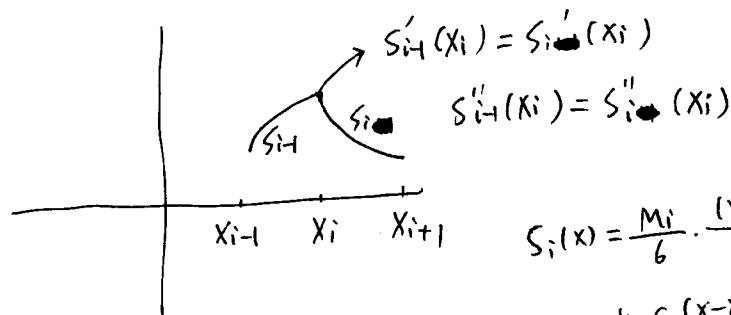
So, interpolate  $2x$ .

$$S_i(x) = \frac{M_i}{6} \cdot \frac{(x_{i+1}-x)^3}{\Delta x_i} + \frac{M_{i+1}}{6\Delta x_i} (x-x_i)^3 + C(x-x_i) + D(x_{i+1}-x)$$

$M_i, M_{i+1}, C, D$  are unknown.

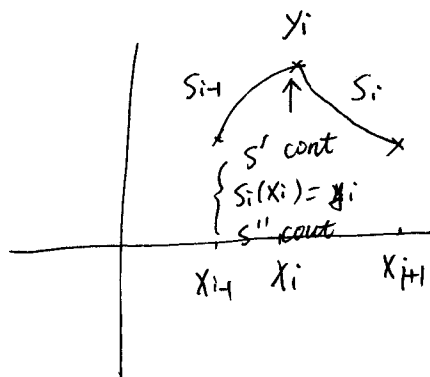
so to find  $M_i, M_{i+1}, C, D$ , from  $\boxed{S_i(x_i) = y_i = y(x_i)}$   $i=0, 1, \dots, N-1$

smoothness



$$S_i(x) = \frac{M_i}{6} \cdot \frac{(x_{i+1}-x)^3}{\Delta x_i} + \frac{M_{i+1}}{6} \cdot \frac{(x-x_i)^3}{\Delta x_i} + C(x-x_i) + D(x_{i+1}-x)$$

2.8.



$$S_i(x) = \frac{M_i}{6\Delta x_i} (x_{i+1}-x)^3 + \frac{M_{i+1}}{6\Delta x_i} (x-x_i)^3 + C_i(x-x_i) + D_i(x_{i+1}-x) \quad M_i \approx y''(x_i)$$

$$M_{i+1} \approx y''(x_{i+1})$$

$$S_i(x_{i+1}) = y_{i+1} = \frac{M_{i+1}}{6\Delta x_i} \Delta x_i^2 + C_i \Delta x_i$$

$$\Rightarrow C = \frac{1}{\Delta x_i} \left\{ y_{i+1} - \frac{M_{i+1}}{6} \Delta x_i^2 \right\}$$

$S_i(x_i) = y_i$  Gives  $D$ .

$$D = \frac{1}{\Delta x_i} \left\{ y_i - \frac{M_i}{6} \Delta x_i^2 \right\}$$

$$S_i(x_{i+1}) = y_{i+1} = \frac{M_{i+1}}{6} \Delta x_i^2 + C \Delta x_i$$

$$\Rightarrow C = \frac{1}{\Delta x_i} \left\{ y_{i+1} - \frac{M_{i+1}}{6} \Delta x_i^2 \right\}$$

$$S_i(x_i) = y_i = \frac{M_i}{6} \Delta x_i^2 + D \Delta x_i$$

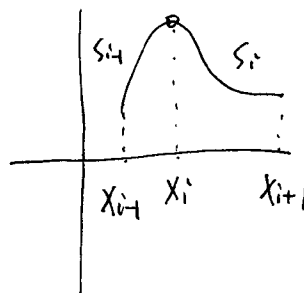
$$\Rightarrow D = \left\{ y_i - \frac{M_i}{6} \Delta x_i^2 \right\} \cdot \frac{1}{\Delta x_i}$$

$$S_i(x) = \frac{M_i}{6\Delta x_i} (x_{i+1}-x)^3 + \frac{M_{i+1}}{6\Delta x_i} (x-x_i)^3 + \frac{1}{\Delta x_i} \left( y_{i+1} - \frac{M_{i+1}}{6} \Delta x_i^2 \right) (x-x_i) + \frac{1}{\Delta x_i} \left( y_i - \frac{M_i}{6} \Delta x_i^2 \right) (x_{i+1}-x)$$

Need  $M_i$   $i=0, 1, \dots, N$

To get those use ~~condition~~ continuity of  $S'$

$$\begin{cases} S'_i(x_i) = -\frac{\Delta x_i}{3} M_i - \frac{\Delta x_i}{6} M_{i+1} + y[x_i, x_{i+1}] \\ S'_{i-1}(x_i) = \frac{\Delta x_{i-1}}{3} M_{i-1} + \frac{\Delta x_{i-1}}{6} M_i + y[x_{i-1}, x_i] \end{cases}$$



Equate

$$\begin{cases} \Delta x_{i-1} M_{i-1} + 2(\Delta x_i + \Delta x_{i-1}) M_i + \Delta x_i M_{i+1} \\ = 6 \{ y[x_i, x_{i+1}] - y[x_{i-1}, x_i] \} \quad i=1, 2, \dots, N-1 \end{cases}$$

$$\left( y[x_i, x_{i+1}] = \frac{y_{i+1} - y_i}{x_{i+1} - x_i} \right) \quad i=1, 2, \dots, N-1$$

Aside: Divide by  $\Delta x_{i-1} + \Delta x_{i+1} = x_{i+1} - x_{i-1}$

The right side is

$$6 \left( \frac{y[x_i, x_{i+1}] - y[x_{i-1}, x_i]}{x_{i+1} - x_{i-1}} \right) = 6 y[x_{i-1}, x_i, x_{i+1}] = y''(\xi) \quad \xi \in [x_{i-1}, x_i]$$

$$y[x_0, x_1, \dots, x_N] = \frac{y[x_1, x_2, \dots, x_N] - y[x_0, x_1, \dots, x_{N-1}]}{x_N - x_0}$$

We have  $N+1$  unknowns,  $N-1$  equations

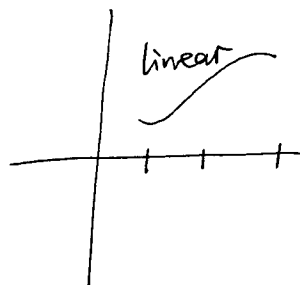
Need 2 more equations

Possibilities:

(1)  $M_0 = y''(x_0)$ ,  $M_N = y''(x_N)$  if we know  $y$ .

(2)  $M_0 = 0$ ,  $M_N = 0$

natural cubic spline.



$$M_0 = 0$$

39

$$\begin{bmatrix} 2(\Delta x_1 + \Delta x_0) & \Delta x_1 & 0 & \dots & 0 \\ \Delta x_1 & 2(\Delta x_2 + \Delta x_1) & \Delta x_2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 2(\Delta x_{N-1} + \Delta x_N) \end{bmatrix} \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_{N-1} \end{bmatrix} = \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_{N-1} \end{bmatrix}$$

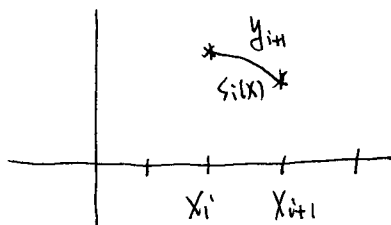
tridiagonal matrix

$$\text{let } A_i = 2(\Delta x_i + \Delta x_{i-1})$$

$$R_i = 6 \{ y[x_i, x_{i+1}] - y[x_{i+1}, x_i] \}$$

$$\begin{bmatrix} A_1 & \Delta x_1 & \dots & 0 \\ \Delta x_1 & A_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \Delta x_{N-2} & \dots & A_{N-1} \end{bmatrix} \vec{M} = \vec{R}$$

2.10



$$f_i(x) = \frac{M_i}{6\Delta x_i} (x_{i+1} - x)^3 + \frac{M_{i+1}}{6\Delta x_i} (x - x_i)^3$$

$$+ \frac{1}{\Delta x_i} \left\{ y_{i+1} - \frac{M_{i+1}}{6} \Delta x_i^2 \right\} (x - x_i)$$

$$+ \frac{1}{\Delta x_i} \left\{ y_i - \frac{M_i}{6} \Delta x_i^2 \right\} (x_{i+1} - x) \quad i = 0, 1, \dots, N-1$$

The  $M_i$ 's are the solution of  $T\vec{M} = \vec{R}$

$$\text{where } T \equiv \begin{bmatrix} A_1 & \Delta x_1 & \dots & 0 \\ \Delta x_1 & A_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \Delta x_{N-2} & \dots & A_{N-1} \end{bmatrix}$$

$$A_i = 2(\Delta x_i + \Delta x_{i-1})$$

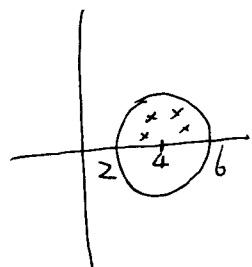
$$R_i = 6 \{ y[x_i, x_{i+1}] - y[x_{i+1}, x_i] \}$$

This is well-conditioned.

$$\text{EX: } \Delta x_i = \Delta x = \text{constant}$$

$$T = \Delta x \begin{bmatrix} 4 & 1 & \dots & 0 \\ 1 & 4 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & \dots & 4 \end{bmatrix}$$

EV's (evaluations) are ~~at~~ inside the circles with centers at  $x=4$  and radii 2 Gerschgorin



$$k(A) = \frac{\max EV}{\min EV} = \|A\|_2 \|A^{-1}\|_2$$

if  $A$  is symmetric

$$k(A) = \frac{\max EV}{\min EV} = \|A\|_2 \|A^{-1}\|_2 \quad \text{if } A \text{ symmetric}$$

$$2 \leq EV \leq 6$$

$$\Rightarrow k \leq \frac{6}{2} = 3$$

In S.P and D.P

This is well-conditioned

Solution of TRI-DIAGONAL systems (diagonal)

$$T = \begin{bmatrix} d_1 & u_1 & 0 & \dots & 0 \\ l_1 & d_2 & u_2 & & \\ \vdots & l_2 & & \ddots & \\ 0 & 0 & \dots & l_{n-1} & d_n \end{bmatrix}$$

Store 3 arrays only.

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

$$u_i \Rightarrow i = 1, 2, \dots, n-1$$

$$l_i \Rightarrow i = 1, 2, \dots, n-1$$

store 0 to  $n$ .

$$T\vec{X} = \vec{y}$$

We use a variant of guess elimination called the Thomas Algorithm.

$$\begin{bmatrix} d_1 & u_1 & & 0 \\ l_1 & d_2 & u_2 & \\ & l_2 & & \ddots \\ 0 & & \dots & l_{n-1} & d_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\begin{bmatrix} \hat{d}_1 & u_1 & & \\ & \hat{d}_2 & u_2 & \\ & & \ddots & \\ & & & \hat{d}_n \end{bmatrix} \vec{\hat{x}} = \vec{\hat{y}}$$

Forward elimination for  $i = 2$  to  $n$

$$d_i = d_i - u_{i-1} * l_{i-1} / d_{i-1} \quad *$$

$$y_i = y_i - y_{i-1} * l_{i-1} / d_{i-1} \quad *$$

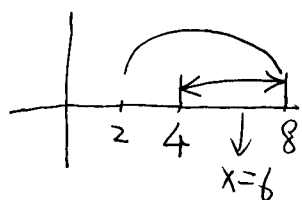
next  $i$



$$\begin{cases}
 x_n = \frac{y_n}{d_n} \\
 \text{For } i = n-1 \text{ to } 1 \text{ step } -1 & \text{backward substitution} \\
 x_i = (y_i - u_i x_{i+1}) / d_i \\
 \text{next } i
 \end{cases}$$

Note: The diagonal and right-hand side are destroyed. Save before hand if need again.

To evaluate the spline, must find the appropriate interval  $[x_i, x_{i+1}]$



$$k = \frac{N}{2}$$

$$Is\ x_k > x?$$

Use binary (Bisection)

THM: Let  $y \in C^H$  on interval  $[a, b]$ , partition  $[a, b]$  into subintervals of width  $\Delta x_i$  such that  $\Delta x \equiv \max_i \Delta x_i$  and  $\beta = \frac{\Delta x}{\Delta x_{min}}$ .

Let  $s$  be the cubic spline on  $[a, b]$ , with endpoints specified with the exactly values. Then

$$\|y^{(n)} - s^{(n)}\|_{\infty} \leq C_n \Delta x^{4-r} \|y^{(4)}\|_{\infty}$$

$$C_0 = \frac{5}{384}, \quad C_1 = \frac{1}{42}, \quad C_2 = \frac{3}{8}, \quad C_3 = \frac{1}{2}\beta + \frac{2}{\beta}$$

$$s_i(x) = \frac{M_i}{6\Delta x_i} (x_{i+1} - x)^3 + \frac{M_{i+1}}{6\Delta x_i} (x - x_i)^3$$

$$+ \frac{1}{\Delta x_i} \left( y_{i+1} - \frac{M_{i+1}}{6} \Delta x_i^2 \right) (x - x_i)$$

$$+ \frac{1}{\Delta x_i} \left( y_i - \frac{M_i}{6} \Delta x_i^2 \right) (x_{i+1} - x)$$

## Z.12 Comments on least squares approximations

A common approximation is linear least squares: for

$$P_m(x) = \sum_{k=0}^m a_k \cdot x^k$$

↓  
coefficient

$$P_m(x) = \sum_{k=0}^m a^k x^k$$

least squares: Find  $\{a_n\}_{n=0}^m \Rightarrow \sum_{i=0}^N (y_i - p(x_i))^2$  is minimum.  
 ↓  
 for a set of data  $\{(x_i, y_i)\}_{i=0}^N$

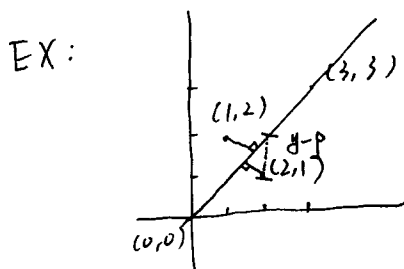
EX: linear regression:

Given  $\{(x_i, y_i)\}_{i=0}^N$  find  $a, b, \exists N$

$\sum_{i=0}^N (y_i - (a + b x_i))^2$  is minimum.

i.e.  $\nabla_{(a,b)} Q = 0$

Gives  $\begin{bmatrix} N & \sum_{i=0}^N x_i \\ \sum_{i=0}^N x_i & \sum_{i=0}^N x_i^2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^N y_i \\ \sum_{i=0}^N x_i y_i \end{bmatrix}$       $\Phi \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$



$$N=4$$

$$\sum x_i = 6$$

$$\sum x_i^2 = 14$$

$$\sum y_i = 6$$

$$\sum x_i y_i = 13$$

$$\begin{bmatrix} 4 & 6 \\ 6 & 14 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 6 \\ 13 \end{bmatrix}$$

$$\begin{bmatrix} a \\ b \end{bmatrix} = \frac{1}{4 \times 14 - 6 \times 6} \begin{bmatrix} 14 & -6 \\ -6 & 4 \end{bmatrix} \begin{bmatrix} 6 \\ 13 \end{bmatrix}$$

$$\Rightarrow P_1(x) = 0.3 + 0.8x$$

General Case:

$$p(x) = \sum_{k=0}^M a_k \phi_k(x) \quad \text{eg. } \phi_k(x) = x^k$$

least square: find  $\vec{a} \Rightarrow \sum_i (y_i - \sum a_k \phi_k(x_i))^2$  is minimum

$$Q = \sum_i (y_i - \sum a_k \phi_k(x_i))^2$$

$$\nabla_a Q = -2 \sum_i$$

$$\nabla_a Q = \frac{\partial Q}{\partial a_L} = 2 \sum_i (y_i - \sum_k a_k \phi_k(x_i)) \phi_L(x_i) = 0$$

$$- \sum_i \sum_k a_k \phi_k(x_i) \phi_L(x_i) + \sum_i y_i \phi_L(x_i) = 0$$

$$\Rightarrow \sum_{k=0}^m a_k \left( \sum_{i=0}^N \phi_k(x_i) \phi_L(x_i) \right) = \sum_{i=0}^N y_i \phi_L(x_i)$$

$$\downarrow$$

$$\Phi_{KL}$$

$$\Phi_{KL} = \sum_{i=0}^N \phi_K(x_i) \phi_L(x_i)$$

$$\text{i.e. } \Phi \vec{a} = \vec{R}$$

$$\Phi_{KL} = \sum_{i=0}^N \phi_K(x_i) \phi_L(x_i)$$

$$\vec{R}_L = \sum_i y_i \phi_L(x_i)$$

$$\Phi \vec{a} = \vec{R}$$

$$\begin{cases} \Phi_{KL} = \sum \phi_K(x_i) \phi_L(x_i) \\ \vec{R}_L = \sum y_i \phi_L(x_i) \end{cases}$$

$$\Phi_K(x) = x^K$$

For  $x \in [0, 1]$ ,  $\phi_K \in \mathcal{X}^K$ ,  $m=9$ ,  $x_i = i \cdot \Delta x$ ,  $\Delta x = \frac{1}{m}$ ,  $K(\Phi) \approx 10^{12}$

Don't use  $\mathcal{X}^K = \Phi_K$ . Instead, use orthogonal polynomials, eg. Legendre

2.15

## 2.4 Quadrature

Most integrals can only be approximated. EX: the pdf.

$$p(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-x^2} dx$$

Two basic classes of integration approximations.

1) Interpolatory quadrature.

2) Monte Carlo

state w/ #1

11) Basic idea: approximate

The integrand by a polynomial & integrate the approximation.

formal derivation:

let  $\{x_j\}_{j=0}^N$  be some points in  $[a, b]$ , then

$$F(x) = P_N(x) + \underbrace{E_N(x)}_{\text{error}} = \sum_{j=0}^N F(x_j) l_j(x) + \frac{F^{(N+1)}(\xi)}{(N+1)!} \prod_{j=0}^N (x-x_j)$$

$$I = \int_a^b F(x) dx = \underbrace{\int_a^b \sum_{j=0}^N F_j l_j(x) dx}_{\tilde{I}} + \underbrace{\int_a^b \frac{F^{(N+1)}(\xi)}{(N+1)!} \prod_{j=0}^N (x-x_j) dx}_{\text{Error}}$$

$$\tilde{I} = \sum_{j=0}^N F_j \cdot \int_a^b l_j(x) dx = \sum_{j=0}^N F_j \cdot w_j$$

$$w_j = \int_a^b l_j(x) dx = \text{Quadrature weights}$$

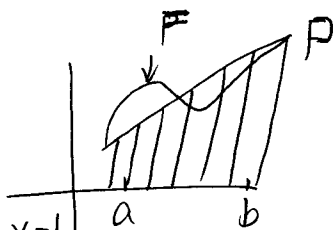
Different methods are found with different choices of  $\{x_j\}_{j=0}^N$ .  
We'll study methods where the  $x_j$ 's are uniformly spaced, called  
Newton-Cotes Methods.

They are called closed if  $a$  &  $b$  are nodes, open if  $a$  &  $b$  aren't nodes.

Examples:

11) Trapezoidal Rule

write  $P_1$  in newton form,  $x_0=a, x_1=b$



$$P_1(x) = (x-a)F[a, b] + F(a) = F(a) + (x-a) \frac{F(b)-F(a)}{b-a}$$

$$\tilde{I} = \int_a^b P_1(x) dx = \int_a^b F(a) dx + F[a, b] \int_a^b (x-a) dx$$

$$= F(a)(b-a) + F[a, b] \frac{(x-a)^2}{2} \Big|_a^b = F(a)(b-a) + F[a, b] \cdot \frac{(b-a)^2}{2}$$

let  $\Delta x = b - a$

$$\tilde{I} = F(a) \Delta x + \frac{\Delta x^2}{2} \cdot \frac{F(b) - F(a)}{b - a} = F(a) \Delta x + \frac{(F(b) - F(a)) \Delta x}{2} = \frac{\Delta x}{2} (F(b) + F(a))$$

$$\tilde{I}_{\text{trap}} = \Delta x \cdot \frac{F(b) + F(a)}{2}$$

ERROR?

$$E = \int_a^b F''(x-a)(x-b) dx = F''(\xi) \int_a^b (x-a)(x-b) dx$$

change variables:  $\begin{cases} x = a + s \cdot \Delta x & s=0 \Leftrightarrow x=a, \quad s=1 \Leftrightarrow x=b \\ x-a = s \cdot \Delta x \\ x-b = a + s \Delta x - b = -\Delta x + s \Delta x = (s-1) \Delta x \end{cases}$

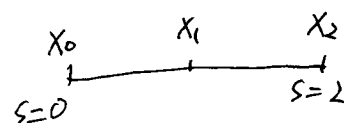
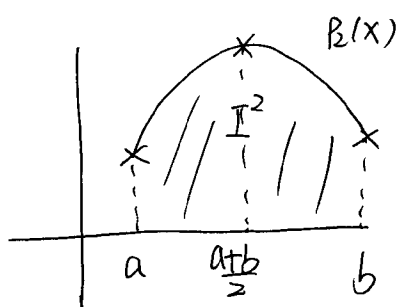
$$E = \Delta x^3 F''(\xi) \int_0^1 s(s-1) ds = \frac{\Delta x^3 F''(\xi)}{2} \left[ \frac{s^3}{3} - \frac{s^2}{2} \right]_0^1 = -\frac{\Delta x^3 F''(\xi)}{12}$$

$$E_{\text{trap}} = -\frac{\Delta x^3 F''(\xi)}{12}$$

12) Simpson's Rule

$$x - x_0 = x - a = s \cdot \Delta x$$

$$(x - x_0)(x - x_1) = s(s-1) \Delta x^2$$



$$x = a + s \cdot \Delta x$$

$$\Delta x = \frac{x - a}{2}$$

$$\tilde{I} = \int_0^2 \left( F_0 + F[x_0, x_1] \cdot s \Delta x + F[x_0, x_1, x_2] s(s-1) \Delta x^2 \right) ds \Delta x$$

$$= F_0 \Delta x \int_0^2 ds + F[x_0, x_1] \Delta x^2 \int_0^2 s ds + F[x_0, x_1, x_2] \Delta x^3 \int_0^2 s(s-1) ds$$

$$\begin{cases} \tilde{I} = 2F_0 \Delta x + 2F[x_0, x_1] \Delta x^2 + F[x_0, x_1, x_2] \Delta x^3 \left( \frac{2}{3} - \frac{2}{3} \right) \\ = \frac{2\Delta x}{6} \{ F_0 + 4F_1 + F_2 \} = \frac{b-a}{6} \{ F(a) + 4F(\frac{a+b}{2}) + F(b) \} \end{cases} \quad \text{EX:}$$

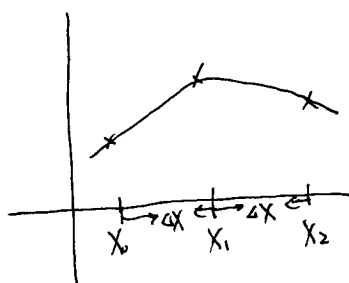
$$\tilde{I}_{\text{Simpson}} = \frac{b-a}{6} \{ F(a) + 4F(\frac{a+b}{2}) + F(b) \}$$

$$E = \frac{\Delta x^4}{3!} \int_0^2 F'' s(s-1)(s-2) ds \stackrel{?}{=} \frac{\Delta x^4}{3!} F''(\xi) \int_0^2 s(s-1)(s-2) ds \quad (r=s-1)$$

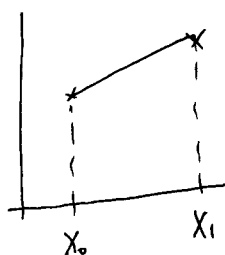
let  $r=s-1$ ,  $s=r+1$ ,  $s-2=r-1$ ,  $\Rightarrow \int_{-1}^1 (r+1)r(r-1)dr \equiv 0$

46

2.17



Simpson



TRAP

$$\int_{x_0}^{x_1} F(x) dx = \frac{\Delta x}{2} (F_0 + F_1) - \frac{\Delta x^3}{12} F''(\xi)$$

$$F_0 = F(x_0)$$

$$\int_{x_1}^{x_2} F(x) dx = \frac{\Delta x}{3} (F_0 + 4F_1 + F_2) + E$$

$$E = ?$$

Find Simpson Error.

$$\tilde{I} = \frac{\Delta x}{3} (F_0 + 4F_1 + F_2) = \frac{\Delta x}{3} \left( F_1 + (\Delta x F_1') + \frac{\Delta x^2}{2} F_1'' + \dots + 4F_1 + F_1 + \Delta x F_1' + \frac{\Delta x^2}{2} F_1'' + \dots \right)$$

Taylor Series  $\rightarrow F_0$  in terms of  $F_1$

$$\rightarrow 4F_1$$

$$\rightarrow F_2 \text{ in terms of } F_1$$

$$\tilde{I} = \frac{\Delta x}{3} (6F_1 + \Delta x^2 F_1'' + 2 \cdot \frac{\Delta x^4}{4!} F_1^{(4)}(\xi))$$

$$= 2\Delta x F_1 + \frac{\Delta x^3}{3} F_1'' + \frac{2\Delta x^5}{3 \times 4!} F_1^{(4)}(\xi)$$

$$\eta_1 - \eta_0 = \Delta x = \eta_2 - \eta_1$$

$$I = \int_{x_0}^{x_2} F(x) dx = \int_{x_0}^{x_2} \left( F_1 + (x-x_1)F_1' + \frac{(x-x_1)^2}{2} F_1'' + \dots \right) dx$$

$$= 2\Delta x F_1 + \frac{(x-x_1)^2}{2} F_1' \Big|_{x_0}^{x_2} + \frac{(x-x_1)^3}{6} F_1'' \Big|_{x_0}^{x_2} + \dots$$

$$= 2\Delta x F_1 + \frac{\Delta x^3}{3} F_1'' + \frac{\Delta x^5}{60} F_1^{(4)} + \dots$$

$$\text{so. } E = I - \tilde{I} = \Delta x^5 F_1^{(4)} \left( \frac{1}{60} - \frac{2}{3 \times 4!} \right) = -\frac{\Delta x^5}{90} F_1^{(4)}$$

$$\int_{x_0}^{x_2} F(x) dx = \frac{\Delta x}{3} (F_0 + 4F_1 + F_2) - \frac{\Delta x^5}{90} F_1^{(4)}$$

Definition: (Precision)

A quadrature has precision  $P$  if it is exact for all polynomials of

$$\text{DEG.} \leq P$$

Examples:

$$I = \int_0^1 e^x dx = e - 1 = 1.7782818 = 1.7182818$$

TRAP Rule:  $\tilde{I} = \frac{e^1 + e^0}{2} = \frac{e+1}{2} = 1.8591409$

$$|E| = \left| \frac{\Delta x}{12} e^{\xi} \right| = \frac{\Delta x}{12} e^{\xi} \quad \xi \in [0, 1]$$

$$\begin{cases} \frac{1}{12}e^0 \leq |E| \leq \frac{1}{12}e^1 \\ 0.0833 \leq |E| \leq 0.2265 \end{cases} \quad |E| = 0.140859$$

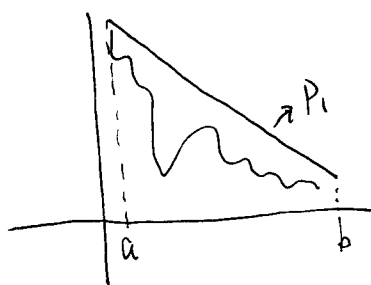
Simpson Rule:

$$\tilde{I} = \frac{1}{6}(e^0 + 4e^{0.5} + e^1) = 1.7188612$$

$$|E| = 5.8 \times 10^{-4}$$

$$\left(\frac{1}{2}\right)^5 \frac{1}{90} \leq |E| \leq \left(\frac{1}{2}\right)^5 \frac{e}{90}$$

$$3.5 \times 10^{-4} < 5.8 \times 10^{-4} < 9.4 \times 10^{-4} \quad \checkmark$$



Bad error if

(1)  $b-a$  is large

(2)  $P^{(k)}(\xi)$  is large

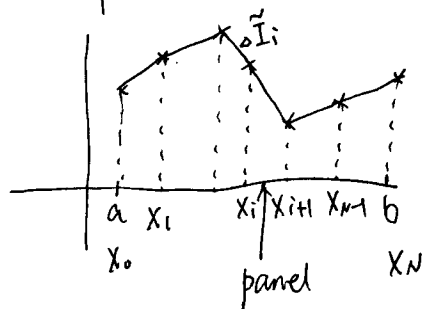
$p$  = precision

## 2.4.2 Composite Rules

If  $[a, b]$  is large, and/or  $F^{(p)}(\xi)$  is large, Newton-Cotes Rule are not accurate.

fix: use piecewise polynomial interpolations.

Composite trapezoidal.



piecewise linear

$$\tilde{I} = \sum_{i=0}^{N-1} \tilde{I}_i$$

Assume  $x_i = a + i\Delta x$   $\Delta x = \frac{b-a}{N}$

$$\tilde{I}_i = \frac{\Delta x}{2} (f(x_i) + f(x_{i+1}))$$

$$\tilde{I} = \sum_{i=0}^{N-1} \tilde{I}_i = \frac{\Delta x}{2} \sum_{i=0}^{N-1} (f(x_i) + f(x_{i+1}))$$

$$\tilde{I} = \frac{\Delta x}{2} (f_0 + 2 \sum_{i=1}^{N-1} f_i + f_N) \rightarrow \text{Composite Trap Rule}$$

$$f_i = f(x_i)$$

$$E = \sum_{i=0}^{N-1} E_i$$

$$E = -\frac{\Delta x^3}{12} \sum_{i=0}^{N-1} F''(\xi_i)$$

2.19

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

$$\tilde{I} = \sum_{i=0}^{N-1} \tilde{I}_i$$



if the spacing  $\Delta x = \text{constant}$

$$\tilde{I} = \frac{\Delta x}{2} \left\{ F_0 + 2 \sum_{j=1}^{N-1} F_j + F_N \right\}$$

trap

Algorithm:

$$S = 0$$

$$\Delta x = \frac{b-a}{N}$$

For  $j = 1$  to  $N-1$

$$x = a + j\Delta x$$

$$S = S + F(x)$$

Next  $j$

$$\tilde{I} = \frac{\Delta x}{2} * (F(a) + 2*S + F(b))$$

Error:

$$E = \sum_{i=0}^{N-1} E_i \quad E_i = -\frac{\Delta x^3}{12} F''(\xi_i) \quad \xi_i \in [x_i, x_{i+1}]$$

$$E = -\frac{\Delta x^3}{12} \sum_{i=0}^{N-1} F''(\xi_i)$$

Note: If  $F''$  is continuous,  $\min_{x \in [a,b]} F''(x) \leq F''(\xi_i) \leq \max_{x \in [a,b]} F''(x)$

Sum over  $i$ 's

$$N \min F''(x) \leq \sum_{i=0}^{N-1} F''(\xi_i) \leq N \max F''(x)$$

$$\text{or } \min F'' \leq \left( \frac{1}{N} \sum_{i=0}^{N-1} F''(\xi_i) \right) \leq \max F''$$

$\Downarrow$

hence



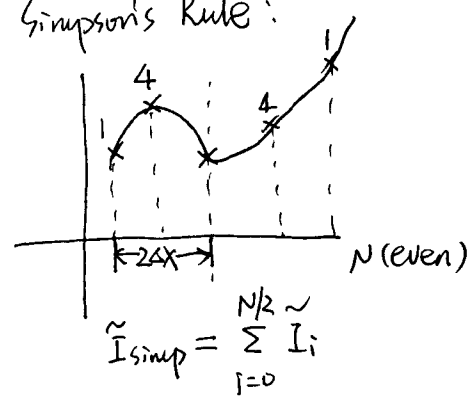
By intermediate value theorem,  $\exists u \in [a, b]$ ,  $\downarrow$  exist such that

$$\frac{1}{N} \sum_{i=0}^{N-1} F''(\xi_i) = F''(u)$$

therefore,  $E = -\frac{\Delta x^3}{12} N F''(u)$   $\Delta x \rightarrow$  spacing between points.

but  $N = \frac{b-a}{\Delta x} \Rightarrow E = -\frac{b-a}{12} \Delta x^2 F''(u)$

Simpson's Rule:



two intervals each panel

Error:

$$E_i = -\frac{\Delta x^4}{90} F^{(4)}(\xi_i) = -\frac{\Delta x^5}{90} F^{(4)}(\xi_i)$$

$$E = -\frac{b-a}{180} \Delta x^4 F^{(4)}(u)$$

$$\tilde{I}_{\text{simp}} = \frac{\Delta x}{3} \left\{ F(a) + \sum_{j=0}^{N/2-1} F_{2j} + 4 \sum_{j=1}^{N/2} F_{2j-1} + F(b) \right\}$$

Algorithm:  $se = 0$   $se \rightarrow$  sum even  
 $so = 0$   $so \rightarrow$  sum odd

$$\Delta x = \frac{b-a}{N}$$

For  $j=1$  to  $N-1$  step 2

$$X = a + j * \Delta x$$

$$so = so + F(X)$$

Next  $j$

For  $j=2$  to  $N-N$  step 2

$$X = a + j * \Delta x$$

$$se = se + F(X)$$

Next  $j$ .

$$\tilde{I} = \frac{\Delta x}{3} (F(a) + F(b) + 2 * se + 4 * so)$$

## 2.4.3 Error Estimation &amp; Adaptive Quadrature

$$I = \tilde{I} + E$$

$$E = C \Delta X^r F^{(m)}(u)$$

$\downarrow$   
 Constant

 $r=2, m=2$  Composite trap.

 $r=4, m=4$ , Composite Simpson.

For  $\Delta X$  (small),  $F^{(m)}(u) \approx \text{constant}$ .

So, as  $\Delta X \rightarrow 0$ ,  $E \rightarrow C \Delta X^r$

$\downarrow$   
 Generic

EX: Trap. Rule

$$\frac{1}{N} \sum_{i=0}^{N-1} F''(\xi_i) = F''(u) \quad \Delta X = \frac{b-a}{N} \quad \Delta X \rightarrow 0 \quad N \rightarrow \infty$$

$$\frac{1}{N} = \frac{\Delta X}{b-a}, \quad \frac{1}{N} \sum_{i=0}^{N-1} F''(\xi_i) = \frac{1}{b-a} \underbrace{\sum_{i=0}^{N-1} F''(\xi_i) \Delta X}_{\text{Riemann sum}}$$

$$\lim_{N \rightarrow \infty} \frac{1}{b-a} \sum_{i=0}^{N-1} F''(\xi_i) \Delta X = \int_a^b F''(\xi) d\xi \cdot \frac{1}{b-a}$$

$$= \frac{F'(b) - F'(a)}{b-a} = \text{constant}$$

As  $\Delta X \rightarrow 0$ ,  $F''(u) \rightarrow \text{constant}$ .

~~Error~~  $\approx$   $I \approx \tilde{I} + \underbrace{C \Delta X^r}_{E \Delta X}$

$r, I$  unknown  
one equation

$$I \approx \tilde{I}_{\Delta X} + E \Delta X$$

Redo with  $\frac{\Delta X}{2}$ .

$$I \approx \tilde{I}_{\frac{\Delta X}{2}} + C \left(\frac{\Delta X}{2}\right)^r$$

2.22.

## Error Estimation (Richardson)

51

$$I \approx \tilde{I}_{\Delta x} + \underbrace{C \Delta x^n}_E$$

 $n \rightarrow r$ 

$$I \approx \tilde{I}_{\frac{\Delta x}{2}} + C \left(\frac{\Delta x}{2}\right)^n$$

$$0 \approx (\tilde{I}_{\Delta x} - \tilde{I}_{\frac{\Delta x}{2}}) + \underbrace{C \Delta x^n}_{E_{\Delta x}} \left(1 - \frac{1}{2^n}\right)$$

$$E_{\Delta x} \doteq \frac{2^n}{2^n - 1} (\tilde{I}_{\frac{\Delta x}{2}} - \tilde{I}_{\Delta x})$$

 $f^{(m)}(u)$  doesn't change with  $\Delta x$ 

$$E_{\frac{\Delta x}{2}} \doteq \frac{1}{2^n - 1} (\tilde{I}_{\frac{\Delta x}{2}} - \tilde{I}_{\Delta x})$$

EX:  $\int_0^1 \frac{1}{1+x} dx = \ln 2$   $n=8, n=16$

Trap. Rule:  $I_8 = 0.694122$

$I_{16} = 0.693391$

$$|E_8| \doteq \left(\frac{2^2}{2^2 - 1}\right) |(0.693391 - 0.694122)| = 9.7467 \times 10^{-4}$$

Actual Error: find in the table  $9.74670 \times 10^{-4}$

$$|E_{16}| = \frac{1}{4} \times 9.7467 \times 10^{-4} = 2.4366 \times 10^{-4}$$

Also, since  $E \approx C \Delta x^r$ , we can estimate  $\Delta x$  needed to get a desired error.

$$E_n = C \Delta x_n^r$$

$$\Rightarrow \frac{E}{E_n} = \left(\frac{\Delta x}{\Delta x_n}\right)^r$$

$$E = C \cdot \Delta x^r = \text{Tol}$$

$$\left(\frac{\text{Tol}}{E_n}\right)^{\frac{1}{r}} \cdot \Delta x_n = \Delta x$$

so the new # of panels  $N = \frac{b-a}{\Delta x}$

EX: with 32 intervals simpson error was  $2.78 \times 10^{-8}$ .

Find  $N$  so trap. rule matches.

$$n=32, \quad E_{32}(\text{trap}) = \frac{1}{32} \sqrt{\frac{2.78 \times 10^{-8}}{6.10 \times 10^{-5}}}$$

$$\Delta x = 6.7 \times 10^{-4}$$

$$N = 1499$$

A simpson error control strategy.

1. compute  $\tilde{I}_{\Delta x}$  with  $\Delta x$

2. compute  $\tilde{I}_{\frac{\Delta x}{2}}$  with  $\frac{\Delta x}{2}$

3. estimate

$$E_{\frac{\Delta x}{2}} = \frac{1}{2^r - 1} (\tilde{I}_{\frac{\Delta x}{2}} - \tilde{I}_{\Delta x})$$

4. if  $|E_{\frac{\Delta x}{2}}| \leq \text{abstol} + |\tilde{I}_{\frac{\Delta x}{2}}| \cdot \text{reltol}$

return  $\tilde{I}_{\frac{\Delta x}{2}} + E_{\frac{\Delta x}{2}}$

5. repeat 1-5.

Suppose we have a function composite rule  $(F, a, b, N)$

↓

$\tilde{I}_N$

$N=2$

for  $k=1$  to  $\text{maxk}$ .

$I_c = \text{composite rule}(F, a, b, N)$

$I_f = \text{composite rule}(F, a, b, 2N)$

$$E = (I_f - I_c) / (2^r - 1)$$

if  $|E| \leq \text{abstol} + |I_f| \cdot \text{reltol}$

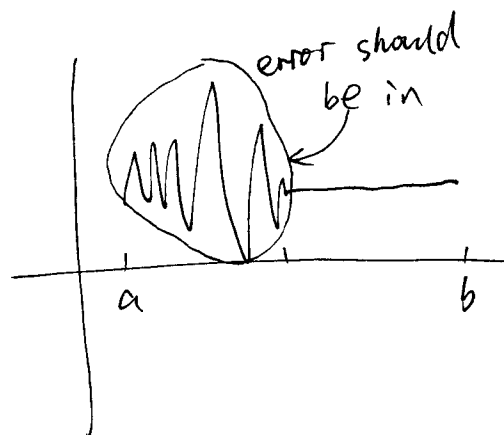
return  $I_f + E$

else

$N = 2N$

endif

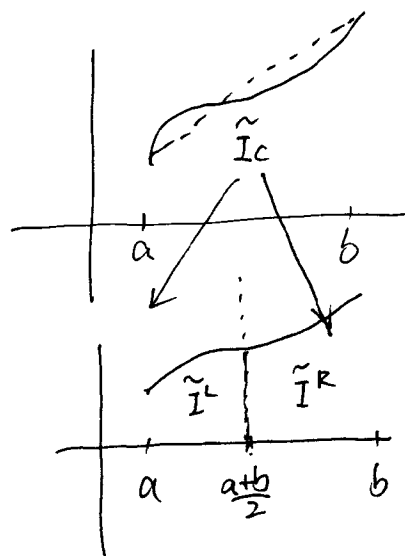
next  $k$ .



## Adaptive Quadrature

$$E_c \equiv I - \tilde{I}_c = \frac{2^r}{2^r - 1} (\tilde{I}_f - \tilde{I}_c)$$

$$\text{if } |E_c| \leq \text{tol?}$$



$$\tilde{I}_f = \tilde{I}^L + \tilde{I}^R$$

(2.24)

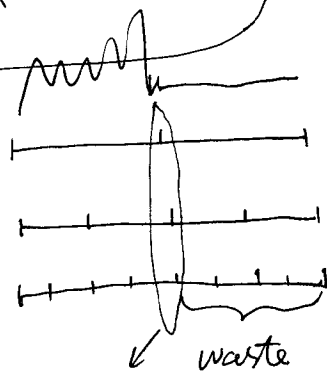
$$\max_{x \in [a, b]} |y^{(r)}(x) - s^{(r)}(x)| \leq C (\Delta x^{4-r}) \|y^{(4)}\|_\infty$$

error  $\leq C \Delta x^{4-r}$

$\Delta x^3$

$E_s \sim \Delta x^4$

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



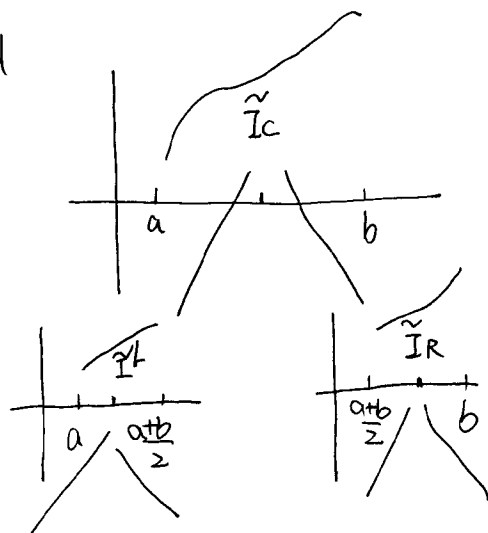
$$n=2, \tilde{I}_2$$

$$n=4, \tilde{I}_4$$

$$E_4 \doteq \frac{1}{2^n - 1} (\tilde{I}_4 - \tilde{I}_2) \quad n=8, \tilde{I}_8 \nparallel \text{repeat}$$

recalculate same point over and over.

Instead



$$E_f \approx \frac{1}{2^r - 1} (\tilde{I}^L + \tilde{I}^R - \tilde{I}_c)$$

repeat

(top)

at each level, we test

$$\frac{1}{2^r} (\tilde{I}^L + \tilde{I}^R - \tilde{I}_c) \leq \text{Tol}.$$

if yes, return  $\tilde{I}^L + \tilde{I}^R + E_f$

if no, we split each half.  $\nparallel$  redo.

on each half we allow  $\frac{\text{tol}}{2}$  error.

so total is approximately

$$\frac{\text{tol}}{2} + \frac{\text{tol}}{2} = \text{tol}.$$

Assume we have a quadrature function  $Q(F, a, b)$

54

Adaptive- $Q(F, a, b, abstol)$

$$I_c = Q(F, a, b)$$

$$I_f = Q(F, a, \frac{a+b}{2}) + Q(F, \frac{a+b}{2}, b)$$

$$E = (I_f - I_c) / (2^r - 1)$$

if  $|E| > abstol$

$$I_L = \text{adaptive-}Q(F, a, \frac{a+b}{2}, abstol/2)$$

$$I_R = \text{adaptive-}Q(F, \frac{a+b}{2}, b, abstol/2)$$

return  $I_L + I_R$

else

return  $I_f + E$

endif

#### 2.4.4 Infinite intervals

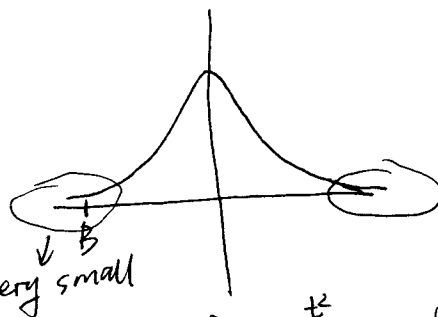
$$E(X): p(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt$$

suppose  $x > 0$ .

$$p(x) \geq \frac{1}{2}$$

$$\text{Need } \int_{-\infty}^B e^{-\frac{x^2}{2}} dx = \epsilon$$

$$\text{find } B, \text{ so } \int_{-\infty}^B e^{-\frac{x^2}{2}} dx \leq \frac{\epsilon}{2} \xrightarrow{\text{machine epsilon}}$$

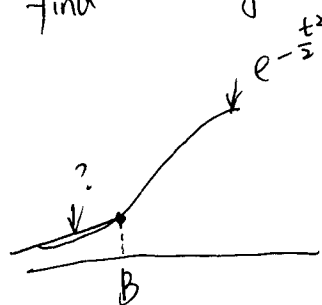


$$\text{find } B, \text{ so } \int_{-\infty}^B e^{-\frac{t^2}{2}} dt \ll \int_B^x e^{-\frac{t^2}{2}} dt$$

↓

can be ignored.  $B?$

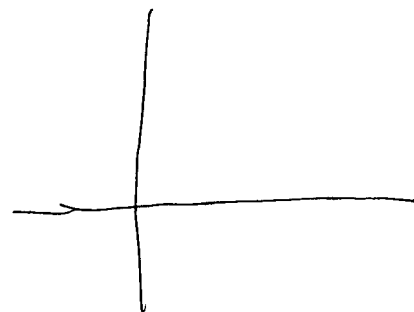
$$\text{find an integral } \Rightarrow \int_{-\infty}^B e^{-\frac{t^2}{2}} dt \leq \int_{-\infty}^B (?) dt < \frac{\epsilon}{2}$$



$$e^{-\frac{B^2}{2}} e^{-\frac{x-B}{2}}$$

$$e^{-\frac{B^2}{2}} \int_{-\infty}^B e^{-\frac{(x-B)^2}{2}} dx = \frac{\epsilon}{2}$$

$$\Rightarrow 2e^{-B^2/2} = \frac{\epsilon}{2} \Rightarrow e^{-\frac{B^2}{2}} = \frac{\epsilon}{4}$$



$$B = -\sqrt{-2 \log \frac{\epsilon}{4}}$$

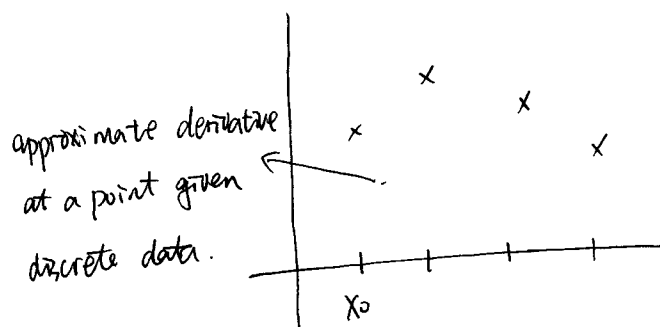
$B \approx -6$  for single precision.

But:  $\int_{-\infty}^0 e^{-\frac{t^2}{2}} dt = \frac{1}{\sqrt{2}}$

$$\frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sqrt{2}} + \int_0^x e^{-\frac{t^2}{2}} dt \right) = p(x)$$

2.26

## 2.6 Differentiation

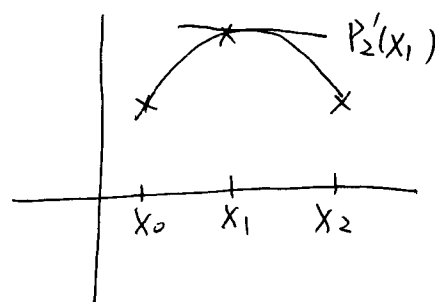


Idea. Approximate the function by an interpolant  $M_0 \approx y''(x_0)$  and differentiate the interpolant.

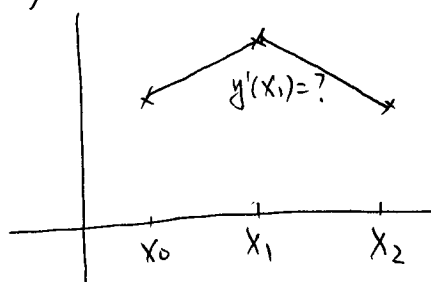
Given a set of point  $\{x_i\}_{i=0}^N$  approx  $y \approx p_N(x)$

$$p_N(x_i) = y_i, \quad y \approx p_N(x)$$

$$\left. \frac{d^k y}{dx^k} \right|_{x_j} \approx \left. \frac{d^k p_N(x)}{dx^k} \right|_{x_j}$$



EX: Three points



$$P_1(x) = y_0 + y[x_0, x_1](x - x_0)$$

$$P'_1(x) = y[x_0, x_1] = \frac{y_1 - y_0}{x_1 - x_0}$$

$$P_1(x) = y_1 + y[x_1, x_2](x - x_1)$$

$$P'_1(x) = y[x_1, x_2] = \frac{y_2 - y_1}{x_2 - x_1}$$

can  $\delta y_j^+ = \frac{y_{j+1} - y_j}{\Delta x_j}$  first order forward derivative approximation

$\delta y_j^- = \frac{y_j - y_{j-1}}{\Delta x_{j-1}}$  backward difference.

Quadratic

$$P_2(x) = y_0 + (x - x_0) F[x_0, x_1] + (x - x_0)(x - x_1) F[x_0, x_1, x_2]$$

$$P'_2(x) = F[x_0, x_1] + [(x - x_1) + (x - x_0)] F[x_0, x_1, x_2]$$

$$P'_2(x_1) = F[x_0, x_1] + \Delta x_0 F[x_0, x_1, x_2]$$

Assume: The  $x_j$ 's are uniformly spaced.

$$\Delta x_j = \Delta x = \text{constant}$$

$$\begin{aligned} P'_2(x_1) &= \frac{y_1 - y_0}{\Delta x} + \Delta x \left[ \frac{\frac{y_2 - y_1}{\Delta x} - \frac{y_1 - y_0}{\Delta x}}{2\Delta x} \right] \\ &= \frac{y_1 - y_0}{\Delta x} + \frac{\Delta x}{2} \cdot \frac{y_2 - 2y_1 + y_0}{\Delta x^2} \\ &= \frac{y_1 - y_0}{\Delta x} + \frac{y_2 - 2y_1 + y_0}{2\Delta x} \\ &= \frac{2y_1 - 2y_0 + y_2 - 2y_1 + y_0}{2\Delta x} = \frac{y_2 - y_0}{2\Delta x} \end{aligned}$$



Define:

$$\delta^0 y_j = \frac{y_{j+1} - y_j}{2\Delta x}$$

centered  
operator

and note  $\delta^0 y_j = (\delta y_j^+ + \delta y_j^-) / 2$

Error?

$$y(x) = P_N(x) + \frac{y^{(N+1)}(\xi)}{(N+1)!} \prod_{i=0}^N (x - x_i)$$

$$y(x) = P_N(x) + \underbrace{\frac{y^{(N+1)}(\xi)}{(N+1)!} \prod_{i=0}^N (x - x_i)}_{E_N(x)}$$

$$\frac{d^k y}{dx^k} = \frac{d^k P_N(x)}{dx^k} + \underbrace{\frac{d^k E_N(x)}{dx^k}}_{E_N^{(k)}(x)}$$

$$E_N^{(k)} = \frac{y^{(N+1)}(\xi)}{(N+1)!} \frac{d^k}{dx^k} \prod_{i=0}^N (x - x_i)$$

for uniform spaced points,  $x_i$  ordered

Let  $x = x_0 + s * \Delta x$

note:  $\frac{d^k}{dx^k} = \frac{d^k}{ds^k} \cdot \frac{1}{\Delta x^k}$

$x_i = x_0 + i * \Delta x$

$x - x_i = (s - i) \Delta x$

$$E_N^{(k)} = \frac{y^{(N+1)}(\xi)}{(N+1)!} \cdot \frac{1}{\Delta x^k} \left( \frac{d^k}{ds^k} \prod_{i=0}^N (s - i) \right) \Delta x^{N+1}$$

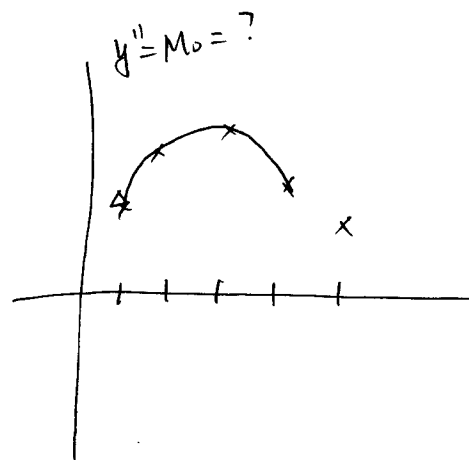
$$= \left[ \frac{y^{(N+1)}(\xi)}{(N+1)!} \frac{d^k}{ds^k} \prod_{i=0}^N (s - i) \right] \Delta x^{N-k+1}$$

n derivative  $\Rightarrow N - k + 1 = n$ .

$N - k + 1 = 2$

$k = 2$

$N = 3$



Definition:  $E_N^{(k)}$  is of order  $r$  if  $E_N^{(k)} = O(\Delta x^r)$

58

In general, we use one of two ways to find error

(1) A-priori

$$\text{use } E_N^{(k)} = \left[ \frac{y^{(n+1)}(\xi)}{(n+1)!} \frac{dk}{s^k} \prod_{i=0}^n (s-i) \right] \Delta x^{n-k+1}$$

(2) A-posteriori

(3) Use Taylor's theorem.

$$\text{EX. } \delta_{y_j}^0 = \frac{y_{j+1} - y_{j-1}}{2\Delta x} \approx y'(x_j)$$

$$y_{j+1} = y_j + \Delta x y'_j + \frac{\Delta x^2}{2} y''_j + \dots$$

$$y_{j-1} = y_j - \Delta x y'_j + \frac{\Delta x^2}{2} y''_j + \dots$$

$$y_{j+1} - y_{j-1} = \left( 2\Delta x y_j^{(1)} + \frac{2\Delta x^3}{3!} y_j^{(3)} + \dots \right) = 2\Delta x y'_j + 2\Delta x^3 \frac{y_j^{(3)}}{3!} + \dots$$

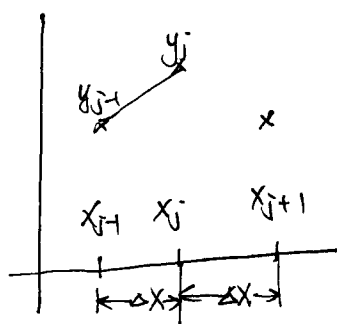
$$\delta_{y_j}^0 = y'_j + \frac{\Delta x^2}{2 \cdot 6} y'''(\xi)$$

so  $\delta^0$  is order 2.

3.1

58

## Finite differentiates



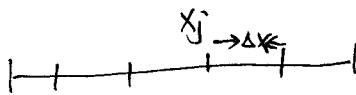
$$\delta^- y_j \equiv \frac{y_j - y_{j-1}}{\Delta x}$$

$$\text{and } y'_j \equiv y'(x_j) = \delta^- y_j + O(\Delta x)$$

$$\delta^+ y_j \equiv \frac{y_{j+1} - y_j}{\Delta x}$$

$\delta^+, \delta^-$  are operators

on a grid



Define the shift operator  $S y_j = y_{j+1}$

$$\text{so } S^2 y_j = S(S y_j) = y_{j+2}$$

$$S^{-1} y_j = y_{j-1}$$

$$\text{then } S S^{-1} y_j = y_j$$

$$\text{define } I y_j = y_j$$

$$\Rightarrow S S^{-1} = I$$

$$\delta^- y_j = \frac{y_j - y_{j-1}}{\Delta x} = \frac{I y_j - S^{-1} y_j}{\Delta x} = \frac{I - S^{-1}}{\Delta x} y_j$$

$$\delta^- = \frac{I - S^{-1}}{\Delta x}$$

$$\delta^+ = \frac{S - I}{\Delta x}$$

Note:

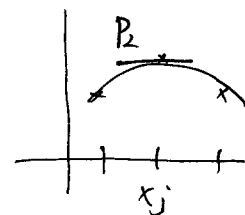
$$\frac{\delta^+ + \delta^-}{2} = \left( \frac{S - I}{\Delta x} + \frac{I - S^{-1}}{\Delta x} \right) \frac{1}{2} = \frac{1}{2} \cdot \frac{S - S^{-1}}{\Delta x} \equiv \delta^*$$

Also,  $\frac{\delta^0 y_j - \delta^0 y_{j+1}}{\Delta x} = \delta^- \delta^0 y_j$

$$\frac{d^2}{dx^2} \approx \delta^- \delta^0, \text{ or } (\delta^0)^2, \text{ or } (\delta^+ \delta^-)$$

$$\delta^+ \delta^- = \frac{(S-I) \times (I-S^-)}{\Delta x \Delta x} = \frac{S-2I+S^{-1}}{\Delta x^2} = \frac{y_{j+1} - 2y_j + y_{j-1}}{\Delta x^2}$$

$$y_{j+1} = y_j + \Delta x y_j' + \frac{\Delta x^2}{2} y_j'' + \dots$$



call  $D = \frac{d}{dx}$

$$y_{j+1} = I y_j + \Delta x D y_j + \frac{\Delta x^2}{2} D^2 y_j + \dots$$

$$S y_j = \underbrace{(I + \Delta x D + \frac{1}{2} \Delta x^2 D^2 + \dots)}_{e^{\Delta x D} = e^{\Delta x \frac{d}{dx}}} y_j$$

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

So  $S = e^{\Delta x D}$

$$\delta^+ = D + E_{\text{error}}$$

$$\frac{\delta - I}{\Delta x} = \frac{e^{\Delta x D} - I}{\Delta x} = \frac{(I + \Delta x D + \frac{\Delta x^2}{2} D^2 + \dots) - I}{\Delta x}$$

$$\delta^+ = D + \left( \frac{\Delta x}{2} D^2 + \dots \right) = D + O(\Delta x)$$

1<sup>st</sup> derivative

$$y_j' = \delta^+ y_j + O(\Delta x)$$

$$y_j' = \delta^- y_j + O(\Delta x)$$

$$y_j = \delta^0 y_j + O(\Delta x^2)$$

$\delta^0 y_j \rightarrow$  centered

$$y_j' = \frac{-y_{j+2} - 4y_{j+1} - 3y_j}{2\Delta x} + O(\Delta x^2)$$

↑

$$X \quad y_j' = -y_{j+2} + 8y_{j+1} + 8y_j \quad \text{one-sided}$$

2<sup>nd</sup> derivative

$$y_j'' = \delta^+ \delta^- y_j = \frac{y_{j+1} - 2y_j + y_{j-1}}{\Delta x^2}$$

$$y_j'' = \frac{-y_{j+2} + 16y_{j+1} - 30y_j + 16y_{j-1} - y_{j-2}}{12\Delta x^2} + O(\Delta x^4)$$

61

3.3.

## 2.7 Integration of odes

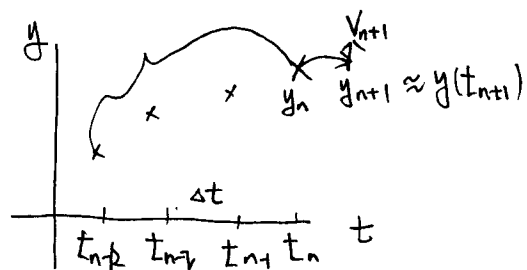
Initial Value Problem IVP

$$\begin{cases} \vec{y}' = \frac{d\vec{y}}{dt} = \vec{F}(t, \vec{y}) & t \geq 0 \\ \vec{y}(0) = \vec{y}_0 \end{cases}$$

EX: Continuous bond price  
 $r(t)$  the interest rate  
 $R(t)$  coupon payment.

$$\frac{dv}{dt} = r(t)V - R(t) \equiv F(t, V)$$

We will look at a class of step-by-step methods  $t_n = n \cdot \Delta t$



if  $k > 0$ , the method is called multistep.

if  $k = 0$ , the method is called single step.

Example: Euler's method [Forward Euler]

$$\Delta y = \frac{dy}{dt} \cdot \Delta t = F \Delta t$$

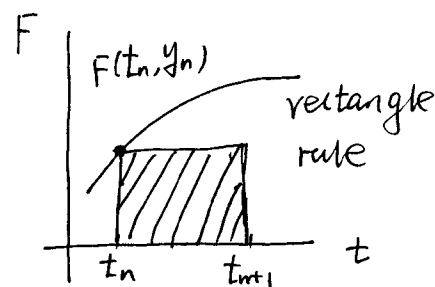
Derivative 1:  $y' = F(t, y)$

$$\int_{t_n}^{t_{n+1}} y' dt = \int_{t_n}^{t_{n+1}} F(t, y) dt$$

$$y_{n+1} - y_n = \int_{t_n}^{t_{n+1}} F(t, y) dt$$

$$y_n \equiv y(t_n)$$

$$y_{n+1} = y_n + \underbrace{\int_{t_n}^{t_{n+1}} F(t, y) dt}_{\text{approx. by quadrature.}}$$



$$F(t, y(t)) = \underbrace{F(t_n, y_n)}_{P_0} + \frac{F'(\xi)}{1!} (t - t_n)$$

$$F(t, y(t)) = \underbrace{F(t_n, y_n)}_{P_0} + \frac{F'(\xi)}{1!} (t - t_n)$$

$F'(\xi) = y''(\xi)$

$$y_{n+1} = y_n + \int_{t_n}^{t_{n+1}} F(t_n, y_n) dt + y''(\xi) \int_{t_n}^{t_{n+1}} (t - t_n) dt$$

$$y_{n+1} = y_n + F(t_n, y_n) \Delta t + y''(\xi) \cdot \frac{\Delta t^2}{2}$$

Drop  $\Delta t^2$  term and let  $y_{n+1}$  satisfy

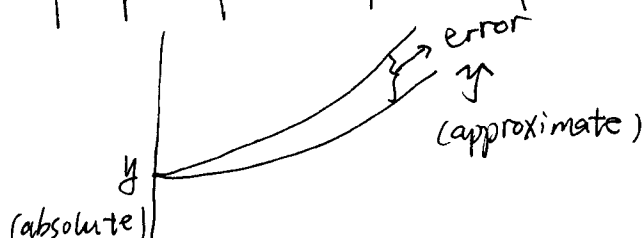
$$y_{n+1} = y_n + \Delta t F(t_n, y_n)$$

Example:  $\begin{cases} y' = y \\ y(0) = 1 \end{cases} \quad y = e^t$

Euler:  $y_{n+1} = y_n + \Delta t y_n = (1 + \Delta t) y_n$

choose  $\Delta t = 0.1$ ,  $y_{n+1} = 1.1 y_n$

(approx. of $y_n$ )			
n	$t_n$	$y_n$	$y_n$
0	0	1	1
1	0.1	1.1	1.105
2	0.2	1.21	1.221
3	0.3	1.331	1.350
4	0.4	1.464	1.492



Derivation 2.

$$y_{n+1} = y_n + y' \Delta t + \frac{y''(\xi)}{2} \cdot \frac{\Delta t^2}{2}$$

$$y' = F(t_n, y_n)$$

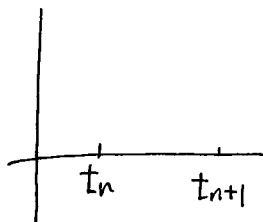
$$y_{n+1} = y_n + F(t_n, y_n) \Delta t + \frac{y''(\xi)}{2} \Delta t^2$$

$$\tilde{y}_{n+1} = y_n + F(t_n, y_n) \Delta t \quad FE.$$

Derivation 3.

$$y' = F(t, y)$$

$$S_t^T y_n = F(t_n, y_n) + \text{Error}$$



$$\frac{y_{n+1} - y_n}{\Delta t} = F(t_n, y_n) + E$$

$$y_{n+1} = y_n + \Delta t F(t_n, y_n) + (\Delta t E) \rightarrow \text{flow away}$$

$$\tilde{y}_{n+1} = y_n + \Delta t F(t_n, y_n)$$

we can find the error over one step a-posteriori

$$y_{n+1} = y_n + (\Delta t \times F(t_n, y_n))$$

$$\tilde{y}_{n+1} = y_n + \Delta t F(t_n, y_n)$$

substitute exact solution.

$$y_{n+1} = y_n + \Delta t F(t_n, y_n) + \tau$$

$$y_{n+1} = y_n + y' \Delta t + \frac{y''(\xi)}{2} \cdot \frac{\Delta t^2}{2}$$

$$= y_n + \Delta t F(t_n, y_n) + \tau$$

$$\Rightarrow \tau = \frac{y''(\xi) \cdot \Delta t^2}{2}$$

↓

truncation error

Definition: The local truncation error is the error created over one time step.

Definition: A step by step method is of order  $r$  if  $\tau = O(\Delta t^{r+1})$

Forward Euler is order 1,  $\tau \approx O(\Delta t^2)$

Definition: A method is consistent if it is at least order 1.

Remark: Euler's Method

$$y_{n+1} = y_n + \Delta t F(t_n, y_n) + y''(\xi) \cdot \frac{\Delta t^2}{2}$$

$$\frac{y_{n+1} - y_n}{\Delta t} = F(t_n, y_n) + \frac{1}{\Delta t} \tau \quad \tau = y''(\xi) \cdot \frac{\Delta t^2}{2}$$

and limit,  $\lim_{\Delta t \rightarrow 0} y' = F + 0$

Want to know.

$$e_n = y_n - \hat{y}_n$$

THM: let  $\Delta t = \frac{T}{N}$ ,  $t_n = n \Delta t$

$$\text{let } y_{n+1} = y_n + \Delta t F(t_n, y_n)$$

$$\text{suppose } \left| \frac{\partial F}{\partial y} \right| \leq L, \left| \frac{\partial F}{\partial t} \right| \leq R, |F| \leq Z.$$

then  $y_N \rightarrow y(T)$  as  $N \rightarrow \infty$  ( $\Delta t \rightarrow 0$ )



3.5.

$$y' = F(t, y)$$

Forward E.  $y_{n+1} = y_n + \Delta t F(t_n, y_n)$

Proof of THM:

$$y_{n+1} = y_n + \Delta t F(t_n, y_n) + \tau_n$$

$$- y_{n+1} = y_n + \Delta t F(t_n, y_n)$$

---

 ~~$e_{n+1} = y_{n+1}$~~   $e_{n+1} = y_{n+1} - y_{n+1}$

$$= e_n + \Delta t [F(t_n, y_n) - F(t_n, y_n)] + \tau_n$$

$$|e_{n+1}| \leq |e_n| + \Delta t |F(t_n, y_n) - F(t_n, y_n)| + |\tau_n|$$

By MVT.

$$|F(t_n, y_n) - F(t_n, y_n)| \leq \left( \frac{\partial F}{\partial y} \right) |y_n - y_n|$$

$$\leq L$$

 $\tau \sim$  truncation error

$$|e_{n+1}| \leq |e_n| + \Delta t L |e_n| + |\tau_n|$$

$$\underbrace{\hspace{1cm}}_{(1+\Delta t L) |e_n|}$$

$$\tau = \max_n |\tau_n|$$

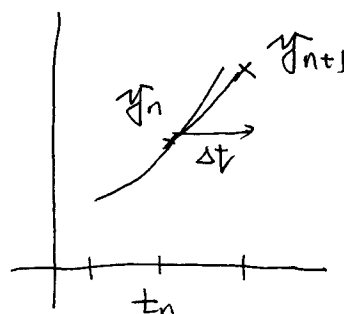
$$|e_{n+1}| \leq (1+\Delta t L) |e_n| + \tau$$

$$\leq (1+\Delta t L) \{ (1+\Delta t L) |e_{n-1}| + \tau \} + \tau$$

$$\leq (1+\Delta t L)^2 |e_{n-1}| + (1+\Delta t L) \tau + \tau$$

$$\dots \leq \underbrace{(1+\Delta t L)^{n+1}}_0 |e_0| + \tau \sum_{j=0}^{n+1} (1+\Delta t L)^j$$

$$|e_n| \leq \tau \sum_{j=0}^{n+1} (1+\Delta t L)^j$$



Recall  $\sum_{j=0}^{N-1} \epsilon^j = \frac{\epsilon^N - 1}{\epsilon - 1}$

so  $|e_N| \leq \tau \frac{(1+\Delta t L)^N - 1}{\Delta t L - 1}$

~~$1 + \Delta t L \leq e^{\Delta t L}$~~

$|e_N| \leq \frac{\tau}{\Delta t} (e^{L T} - 1) / L \Rightarrow |e_N| \leq \frac{\tau}{\Delta t} \frac{e^{L T} - 1}{L}$

But  $\tau = \max_{\xi} \frac{1}{2} |y''(\xi)| \Delta t^2$

$|e_N| \leq \frac{1}{2} \max |y''(\xi)| \Delta t \cdot \frac{e^{L T} - 1}{L}$

But  $y'' = F' = \frac{dF}{dt} = \frac{\partial F}{\partial t} + \frac{\partial F}{\partial y} \cdot \frac{dy}{dt} \quad (F = \frac{dy}{dt})$   
 $= \frac{\partial F}{\partial t} + \frac{\partial F}{\partial y} \cdot F$

so,  $y'' = \frac{\partial F}{\partial t} + \frac{\partial F}{\partial y} F \leq K + LZ$  just a number

so,  $|e_N| \leq \Delta t \left( \frac{1}{2} (K + LZ) \frac{e^{L T} - 1}{L} \right)$

$\lim_{\substack{N \rightarrow \infty \\ \Delta t \rightarrow \frac{T}{N}}} |e_N| = 0$

Also,  $|e| = O(\Delta t)$   
 $(\tau = O(\Delta t^2))$

RMK: consistency is a necessary condition for convergence.

Except: we assume that

$(1 + \Delta t L)^N \cdot |e_0| \equiv 0$

so, really,  $|e_N| \leq e^{L T} |e_0| + O(\Delta t)$

if  $L T$  is large, this doesn't tell us much.  
 need condition on growth of errors.

Test problem:

$$y' = \lambda y \quad \lambda \in \mathbb{C}$$

$$(\vec{y}' = \vec{F}(t, \vec{y}))$$

↓ linearize

$$\vec{y}' = A(t) * \vec{y}$$

↓ matrix

$$\text{freeze } t \quad \vec{y}' = A \vec{y}$$

↓ diagonalize

$$A = S \Lambda S^{-1}$$

$$\vec{y}' = S \Lambda S^{-1} \vec{y}$$

$$S^{-1} \vec{y}' = \Lambda S^{-1} \vec{y}$$

define  $\vec{w} = S^{-1} \vec{y} \Rightarrow \vec{w}' = \Lambda \vec{w}$

$$w_i' = \lambda w_i$$

Now let  $u$  = solution with no error initial

$v$  = solution with initial error

$$u' = \lambda u$$

$$v' = \lambda v$$

$$\hat{y}' = (u-v)' = \lambda \hat{y}$$

↑  
 $u-v$  error

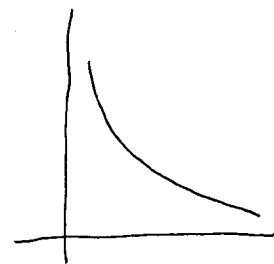
$$\hat{y}(0) = e_0$$

and scale  $y = \frac{\hat{y}}{e_0}$

$$\begin{cases} y' = \lambda y \\ y(0) = 1 \end{cases} \rightarrow \text{the problem we have to test}$$

Apply Euler's method

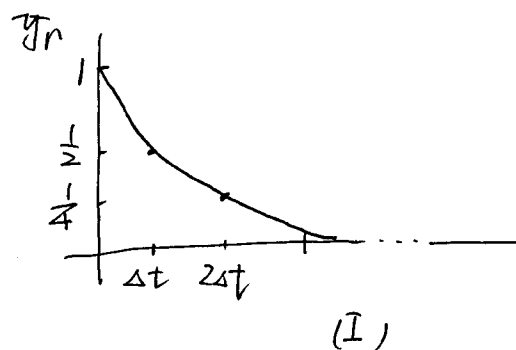
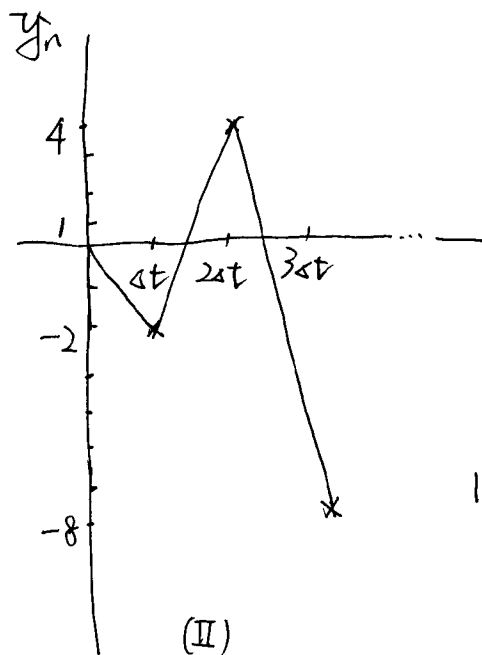
$$\begin{aligned} y_{n+1} &= y_n + \Delta t (\lambda y_n) \\ &= (1 + \lambda \Delta t) y_n \end{aligned}$$



Assume  $\lambda$  is real and negative ( $\lambda < 0$ )

CASE I: choose  $\lambda \Delta t$ , such that  $1 + \lambda \Delta t = \frac{1}{2}$

CASE II: choose  $\Delta t$ ,  $\Rightarrow 1 + \lambda \Delta t = -2$



$n \rightarrow \infty$   
 $|y_n| \rightarrow \infty$  blows up.

To Not blow up,

$$|1 + \lambda \Delta t| \leq 1$$

Definition: A step-by-step method is absolutely stable if

$$|y_{n+1}| \leq |y_n|$$

Forward Euler is absolutely stable if  $|1 + \lambda \Delta t| \leq 1$