

I Example: Modulus of Elasticity of Steel

$$E = \frac{4mg\ell^3}{da^3}$$

(see Example on page 31)

II Floating Point Number System

a) Definition: $F(B, k, m, M)$

B : the base
 k : # of digits in Base expansion
 m : minimum exponent
 M : max exponent

i.e. $\pm (0d_1d_2\dots d_k)_B \times B^e$
 where $d_1 \neq 0$ & $m \leq B \leq M$

b) Example: $F(10, 2, 0, 2)$

i.e. 2 digits, base 10, $0 \leq e \leq 2$

$$= \{ .10, .11, \dots, .99 \} A \quad e=0$$

$$\{ 1.0, 1.1, \dots, 9.9 \} A \quad e=1$$

$$\{ 10, 11, \dots, 99 \} \quad -e=2$$

Types of Errors
Modeling
Discretization/Truncation
Roundoff/Data
Human Error

III. Round-off Error

a) Definition: Round off error is error induced by converting real # to float pt. equivalent

b) Example: $\sqrt{7.1} = 2.66458 \dots$

in $F(10, 2, 0, 2)$

$\sqrt{7.1} = 2.6$ ("chopping" the number = " fl_{chop} ")

$\sqrt{7.1} = 2.7$ ("rounding the number = " fl_{round} ")

c) Definition: Let p^* be any approximation of a number p

① Absolute Error: $|p^* - p|$

② Relative Error: $|p^* - p| / |p| \Rightarrow p \neq 0$

i.e. Example above using fl_{chop}

$$p^* = 2.66458 \quad p = 2.6$$

$$\Rightarrow \text{Abs Err} = |2.6 - 2.66458| = 6.458 \times 10^{-2}$$

$$\text{Rel Er} = \frac{|2.6 - 2.66458|}{|2.66458|} = .0242 \text{ or } 2.42\%$$

(2)

d) Absolute Error General Expressions

$$\textcircled{1} \frac{|fl_{\text{chop}}(y) - y|}{|y|} \leq B^{1-k} \quad (\text{proof on p 35})$$

↑
for $fl(B, k, m, M)$

$$\textcircled{2} \frac{|fl_{\text{round}}(y) - y|}{|y|} \leq \frac{1}{2} B^{1-k} \quad (\text{test for proof in HW})$$

e) Definition Machine Precision 'u'

$$u = \begin{cases} B^{1-k} & \text{chop} \\ \frac{1}{2} B^{1-k} & \text{round} \end{cases}$$

I.E single precision is $F(2, 25, -125, 128)$ Base 2

↑
'u'

(IEEE standard)

$$u = \begin{cases} 2^{1-25} = 2^{-24} & \text{chop} \\ \left(\frac{1}{2}\right) 2^{1-25} = 2^{-25} & \text{Round.} \end{cases}$$

I.E. double precision is $F(2, 53, -1021, 1024)$

∴ $u =$ _____

(f)

Definition: $B^{-(t+1)} < \left| \frac{x-y}{x} \right| < B^{-t}$

Then x & y agree to at least t
and at "most" $t+1$ "significant"
base B digits"

Example: 1.9 p 36

IV Conditioning

Solve: $x' - x = e^{-2t}$, $x(0) = 1/2$ Integrating Factor $u(t) = e^{\int -1 dt} = e^{-t}$

$$\Rightarrow e^{-t} x' - e^{-t} x = e^{-3t} \Rightarrow \int \frac{d}{dt} [e^{-t} x] = \int e^{-3t} dt$$

$$\Rightarrow -e^{-t} x = -\frac{1}{3} e^{-3t} + C$$

$$\Rightarrow x = -\frac{1}{3} x^{-2t} + C e^t \Rightarrow x(0) = -\frac{1}{3} + C = -\frac{1}{3} \Rightarrow C = 0$$

$$\therefore \boxed{x = -\frac{1}{3} x^{-2t}}$$

\Rightarrow Now perturb initial condition i.e. $x(0) = -\frac{1}{3} + \epsilon$

$$x(0) = -\frac{1}{3} + C = -\frac{1}{3} + \epsilon \Rightarrow C = \epsilon$$

$$\therefore \boxed{x = -\frac{1}{3} e^{-2t} + \epsilon e^t}$$

this tiny 'bump' causes
solution to 'blow up'

Problem is "ill-conditioned"

V Appendices

A PROOF OF General Error Expression for "Chopping"

$$\text{For } F(B, k, m, M) \quad \frac{|F_{\text{chop}}(y) - y|}{|y|} \leq B^{1-k}$$

① let $y = (d_1 d_2 \dots d_k d_{k+1} d_{k+2} \dots)_B \times B^e$

② $F_{\text{chop}}(y) = (d_1 d_2 \dots d_k)_B \times B^e$

③ $|F_{\text{chop}}(y) - y| = \underbrace{(0000 \dots 0}_{k \text{ zeros}} d_{k+1} d_{k+2} \dots)_B \times B^e$

$$= (d_{k+1} d_{k+2} \dots)_B \times B^{e-k}$$

④ $\frac{|F_{\text{chop}}(y) - y|}{|y|} = \frac{(d_{k+1} d_{k+2} \dots)_B \times B^{e-k}}{(d_1 d_2 \dots)_B \times B^e}$

$$= \frac{(d_{k+1} d_{k+2} \dots)_B \cdot B^{-k}}{(d_1 d_2 \dots)_B}$$

⑤ note: $(d_1 d_2 \dots)_B \geq (01)_B = B^{-1}$

and $(d_{k+1} d_{k+2} \dots)_B < 1_B$

$$\therefore \frac{|F_{\text{chop}}(y) - y|}{|y|} < \frac{1_B \cdot B^{-k}}{B^{-1}} = B^{1-k}$$

QED

B) Matlab eps "smallest # in Matlab"

$$\text{eps} = 2.2204... \times 10^{-16}$$

$$\log_2(\text{eps}) = -52$$

$$\therefore \text{precision} = B^{1-k} = 2^{-52} \Rightarrow k = 53$$

Recall Double Precision

$$F(2, 53, -1021, 1024)$$

↑

Matlab actually

appears to be \surd -1074, 1023