

## FRIDAY

- **QUIZ TODAY** - quiz sheets at back of room.
- **MATLAB WARM-UPS** - go through worksheets carefully **due MON.**

## LAST DAY

- running Matlab's GE  $\rightarrow$  "backslash"  $\vec{x} = [A] \backslash \vec{b}$
- the GOOS, the BOS & the UGLY.  $\rightarrow$  1<sup>st</sup> computing report.

- **finite-precision** arithmetic leads to errors

• we want to find  $\vec{x}$  so that  $[A]\vec{x} = \vec{b}$

• so we use GE:  $\vec{x}_{GE} = [A] \vec{b}$

• if we then check  $\vec{x}_{GE}$ :  $[A]\vec{x} - \vec{b} = \text{error} = 0$

• expect these errors to as matrix size  $\sqrt{\text{arithmetic}}$

• demo #1  $\rightarrow$  for typical "GE" matrices **ERROR GROWTH** seems

$$\text{obs: } \text{rms}([A]\vec{x} - \vec{b}) = O(10^{-10} - 10^{-5}), \quad N = (10)$$

• demo #2  $\rightarrow$  there are small matrices ( $N < 10$ )

for which the unusually

MATLAB FLOATING-POINT (from )

$$\text{sign} \rightarrow \pm (d_1 d_2 d_3 \dots d_N) \times 10^P \leftarrow \text{exponent}$$

mantissa

• largest magnitude  $\pm \text{realmax} \sim 1.79 \times 10^{+308}$   
violating this is

• smallest magnitude  $\pm \text{realmin} \sim 2.22 \times 10^{-308}$   
violating this is

• estimate of #  $10^{72} - 10^{87}$

UNAVOIDABLE ERRORS IN COMPUTER ARITHMETIC

• addition of floating point

$$a + b = (+D_a \times 10^{P_a}) + (+D_b \times 10^{P_b})$$

• need to add of (say  $P_a > P_b$ )

• example, say  $P_b = P_a - 5$

$$\begin{array}{r}
 + \quad \underbrace{a.aaa \quad aaaa \quad aaaa \quad aaaa}_{5 \text{ for } 10^{P_a - 5}} \quad \left| \quad \underbrace{6666 \quad 6}_{\text{last digits used}} \times 10 \\
 \quad \times 10
 \end{array}$$

- can simulate smaller  $\sqrt{\text{an}}$  calculator

- add  $(\frac{1}{3} * 10^5) + (\frac{1}{7} * 10^3)$

```
>> (1/3)*1e5
```

```
ans =
```

```
3.333333333333333e+04
```

keeping only 5 digits (IEEE std uses *rounding\**)

$$\begin{array}{r}
 3.3333 \times 10^4 \\
 + \quad \quad \quad 3 \times 10^3 \\
 \hline
 3.3476 \times 10^4
 \end{array}$$

if carry,

(\* rounding occurs in binary.)

- multiplication has less of a last digit problem

$$\begin{array}{r}
 a.aaa \dots aaaa \times 10^{P_a} \\
 \times \quad b.bbb \dots bbbb \times 10^{P_b} \\
 \hline
 c.ccc \dots ccc \times 10^{P_a+P_b}
 \end{array}$$

if carry  $\rightarrow$  then

possible  
or

(similar for division)

- week 03 quiz - artificial truncation error

PRACTICE + BRING CALCULATOR

- subtraction is interesting for nearly equal numbers

$$\begin{array}{r}
 1.11a \text{ } aaaa \dots aaaa \times 10^p \\
 - 1.11b \text{ } bbbb \dots bbbb \times 10^p \\
 \hline
 0.00c \text{ } cccc \dots cccc \times 10^p \\
 c.ccc \text{ } cccc \dots c \times 10^p
 \end{array}$$

↑  
"padded" zeros

answer contains only sig digits ( )

- in sequential calculations  $(a-b) * d \rightarrow$   
only sig digits (even though  
- digits)

- in sequential calculations, need to round-off at step

- NOTE: can " " lost digits by addition to magnitude #

- example.

$$\begin{array}{r}
 \begin{array}{cccc}
 6 & 6666 & 6.666 & 6666 & 666 \\
 & & a.aaa & aaaa & aaaa
 \end{array} \\
 \hline
 6 & 666c & cccc & cccc & ccc
 \end{array}$$

erroneous  
xxxxx  
\* 10<sup>p</sup>  
\* 10<sup>p</sup>

irrelevant.

## ABSOLUTE vs RELATIVE ERRORS (p17)

- for  $x^*$  an approximated value for  $x$  (exact)

$$\text{absolute error} = | - |$$

$$\text{relative error} = \frac{| - |}{| |}$$

- what we really mean by finite-precision truncation error is really error of  $10^{-16}$

```
>> factorial(19)
```

```
ans =
```

```
1.216451004088320e+17
```

- in last day's GE residual experiment, it might have been better to use error relative to  $\text{rms}(\vec{r})$ .

That is,  $\text{relative residual error} = \frac{\text{rms}([A]\vec{x} - \vec{b})}{\text{rms}(\vec{r})}$

(but the experiment was designed to have  $\text{rms}(\vec{r}) \approx 1$ )

SIMPLE



## STRATEGIES for REDUCING FINITE-PRECISION ERROR

- recall the calculus limit  $\lim_{\theta \rightarrow 0} \frac{1 - \cos \theta}{\theta} = -$

```
>> calcLim = @(th) (1-cos(th))./(th.^2)
```

```
calcLim =
```

```
function_handle with value:
```

```
@(th)(1-cos(th))./(th.^2)
```

1) there is an analytical 'fix' for this -

$$\frac{1 - \cos \theta}{\theta^2} = \frac{1}{2} \left( \frac{\theta/2}{\theta/2} \right)^2$$

2) reducing # of operations?

• how many operations to evaluate polynomial

$$p(x) = \sum_{k=0}^N c_k \cdot x^k$$

$\underbrace{\hspace{10em}}_{\text{additions}} \quad \underbrace{\hspace{10em}}_{\text{multiplications per term}}$

$$( + + + \dots + ) = + ( )$$

• NESTED EVALUATION (p24)

$$p(x) = ( \dots ( (c_N \cdot x + c_{N-1}) + c_{N-2} ) + \dots + c_0 )$$

additions + multiplications

3) respecting magnitudes.

$$S = (1 \dots) e^{17} + 1$$

- SORTES SUMMATION ( , but )
- add , magnitude quantities
- Matlab's sort command + abs.



muraki

these notes are for the use of SFU students in MACM 316 (spring 2019) & SFU copyright applies