

I. Floating Pt ARITHMETIC

$$X @_{s1} y = fl(fl(x) @_{s1} fl(y))$$

① @ is a generic term for $+$, $-$, $*$, \div

② $fl(x)$ means round x to maximum allowable digits

a) Example: for $F(10, 2, -2, 2)$ Find

$$\frac{1}{7} +_{fl} \frac{1}{23} = fl(fl(\frac{1}{7}) +_{fl} fl(\frac{1}{23}))$$

$$= fl(.14 + .043)$$

$$= fl(.183) = .18$$

⇒ Note: IF I DID THIS ON CALC I'D GET .186335...

b) ARITHMETIC PROPERTIES DO NOT WORK FOR FP OPERATIONS

Example $F(10, 4, -2, 2)$

$$a) (.1329 +_{fl} 1.543) +_{fl} 23.21$$

$$= fl(1.6693) +_{fl} 23.21$$

$$= 1.669 +_{fl} 23.21 = fl(24.879) = \underline{24.89}$$

$$b) .1329 +_{fl} (1.543 +_{fl} 23.21) = \underline{24.88}$$

IF FP ARITHMETIC IS NOT DISTRIBUTIVE!!

II) Accumulation of Round-off Error

a) absolute deviation $\delta_x = fl(x) - x$

b) relative deviation $E_x = \frac{fl(x) - x}{x}$ (or $\frac{\delta_x}{x}$)

$\therefore \boxed{fl(x) = x + \delta_x}$ or $fl(x) = \boxed{x(1 + E_x)}$

c) Multiplication

$$\begin{aligned} fl(x) * fl(y) &= x(1 + E_x) y(1 + E_y) \\ &= xy(1 + E_x + E_y + E_x E_y) \end{aligned}$$

↑ this type of error will not "blow up" since E_x, E_y are small

d) division

$$\frac{fl(x)}{fl(y)} = \frac{x(1 + E_x)}{y(1 + E_y)} = \frac{x}{y} \frac{(1 + E_x)}{(1 + E_y)}$$

$$\Rightarrow \frac{1}{1 + E_y} = 1 - E_y + E_y^2 - E_y^3 \dots$$

$$\therefore \frac{fl(x)}{fl(y)} = \frac{x}{y} (1 + E_x) (1 - E_y + E_y^2 - E_y^3 \dots)$$

again these terms will not "blow up"

e) addition/subtraction

$$f(x) \pm f(y) = (x + \delta x) \pm (y + \delta y)$$

$$= (x \pm y + \delta x \pm \delta y)$$

$$= (x \pm y) \left(1 + \frac{\delta x \pm \delta y}{x \pm y} \right)$$

↑
if $x \pm y$ close to zero
this will "blow up"
("cancellation error")

Algorithms should avoid subtraction of nearly equal numbers when possible

Example: Quadratic Equation

$$ax^2 + bx + c = 0 \Rightarrow x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

① what if $b \gg 4ac \Rightarrow$ then $\sqrt{b^2 - 4ac} \approx b$

② if $b > 0$, then wish to avoid the + operation

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

CONTINUES \Rightarrow

Strategy

WHY IS THIS NOT A PROBLEM

$$\frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad \frac{(-b - \sqrt{b^2 - 4ac})}{(-b - \sqrt{b^2 - 4ac})}$$

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

$$\frac{2c}{-b - \sqrt{b^2 - 4ac}}$$

↑
We have avoided the operation

SEE Example 1.10 /
Page 45

III.) Computer Demo Example 1.12, p 48