

Round-off errors

Round-off errors occur when you do not have enough digits to store the answer.

ex) 8-bit integer : $\frac{7}{2} = 3.5$ \leftarrow not enough digits to store 3.5

In double precision the precision is 2^{-52}

Roughly 2.22×10^{-16} in base-10

Called machine precision & given by `eps` in Matlab.

ex) Addition problem w/ a machine precision of $10^{-3} = 0.001$ (you can store 3 sig-figs)

Exact: $0.99 + 0.0044 + 0.0042 = 0.9986$

3-digit w/ rounding

$$(0.99 + 0.0044) + 0.0042 = 0.994 + 0.0042 = 0.998$$

$$0.99 + (0.0044 + 0.0042) = 0.99 + 0.0086 = 0.999$$

Frequent in subtraction and division.

Re-arranging the calculation can reduce these errors.

e.g) Roots of $x^2 - bx + 1 = 0$, b is large

$$r = \sqrt{b^2 - 4} \quad x_1 = \frac{b+r}{2} \quad x_2 = \frac{b-r}{2}$$

If b is large then $b+r$ are close to each other

If $b = 110$ $\Rightarrow r = (110^2 - 4)^{1/2} = 109.9818 \dots$

True answer is $109.99 \dots$ and $0.00909166 \dots$

Now consider a machine with $\epsilon = 10^{-3}$ (3 sig figs)
using chopping

$$\Rightarrow b = 110 \quad r = 109$$

$$x_1 = \frac{b+r}{2} = \frac{(110 + 109)}{2} = \frac{219}{2} \stackrel{\text{actually } 109.5}{=} 109 \leftarrow 3 \text{ digits correct}$$

$$x_2 = \frac{b-r}{2} = \frac{(110 - 109)}{2} = \frac{1}{2} = 0.5 \leftarrow 0 \text{ digits correct}$$

$$\text{but: } x_2 = \frac{b-r}{2} \frac{b+r}{b+r} = \frac{b^2 - r^2}{2(b+r)} = \frac{4}{2(b+r)} = \frac{2}{b+r}$$

$$\text{Then } x_2 = \frac{2}{b+r} = \frac{2}{110+109} = \frac{2}{219} = 0.00913 \leftarrow \text{one digit \& much closer}$$

Why we care: When solving $Ax = \underline{b}$
we will decompose A ,

To do this you do repeated operation
on \underline{A} , These use $+$, $-$, $*$, $/$

Certain algorithms are less susceptible to
round off errors (more stable) &
this give better answers,