

as we know, representing a number in computer with floating point representation will cause some round off errors.

For example in previously discussed 6 bit floating number we don't have number 11.75. instead we have 8 10 ^{11.75} 12 14

So if we want to store 11.75 in computer with in a 6 bit floating point number representation, we will have some round off errors (as previously discussed, the relative error will be always less than ϵ_m (relative machine precision). ϵ_m is equal to 2^{-t} in which t is number of bits used to represent mantissa (23 for 32 bit, and 2 for my 6 bit floating point number)

x $\xrightarrow[\text{will turn it to}]{\text{Saving in Computer}}$ $x(1+\varepsilon)$; $|\varepsilon| \leq \varepsilon_m$

Note that the relative error is less than ε_m

$$\left| \frac{x - x(1+\varepsilon)}{x} \right| = \varepsilon$$

★ So whenever we want to save a value in computer (be it an input value or the result of a calculation), we will have the roundoff error

Error of addition/subtraction

$$x_1 \longrightarrow \text{rd}(x_1) = x_1(1+\varepsilon_1) = \tilde{x}_1$$

$$x_2 \longrightarrow \text{rd}(x_2) = x_2(1+\varepsilon_2) = \tilde{x}_2$$

$$\tilde{x}_1 + \tilde{x}_2 \longrightarrow \text{rd}(\tilde{x}_1 + \tilde{x}_2) = (\tilde{x}_1 + \tilde{x}_2)(1+\alpha)$$

$$(x_1(1+\varepsilon_1) + x_2(1+\varepsilon_2))(1+\alpha)$$

$$= (x_1 + x_2 + x_1 \varepsilon_1 + x_2 \varepsilon_2) (1 + \alpha)$$

$$= x_1 + x_2 + x_1 \varepsilon_1 + x_2 \varepsilon_2 + \alpha(x_1 + x_2)$$

$$\text{relative Error} = \frac{rd(\tilde{x}_1 + \tilde{x}_2) - (x_1 + x_2)}{(x_1 + x_2)}$$

$$= \frac{x_1 \varepsilon_1 + x_2 \varepsilon_2 + \alpha(x_1 + x_2)}{(x_1 + x_2)}$$

$$\text{Relative } E = \frac{x_1}{x_1 + x_2} \varepsilon_1 + \frac{x_2}{x_1 + x_2} \varepsilon_2 + \alpha$$

if $x_1 \approx -x_2$ then we will have large relative errors

Error of multiplication/division

$$x_1 \rightarrow \tilde{x}_1 = rd(x_1) = x_1(1 + \varepsilon_1)$$

$$x_2 \rightarrow \tilde{x}_2 = rd(x_2) = x_2(1 + \varepsilon_2)$$

$$x_1 x_2 \rightarrow rd(\tilde{x}_1 \tilde{x}_2) = x_1 x_2 (1 + \varepsilon_1)(1 + \varepsilon_2)(1 + \alpha)$$

$$(x_1 + x_1 \varepsilon_1)(x_2 + x_2 \varepsilon_2)(1 + \alpha)$$

$$= x_1 x_2 + x_1 x_2 \varepsilon_1 + x_1 x_2 \varepsilon_2 + \alpha x_1 x_2$$

$$= x_1 x_2 (1 + \varepsilon_1 + \varepsilon_2 + \alpha)$$

$$\left\{ \begin{array}{l} \text{relative} \\ \text{error} \end{array} = \frac{rd(x_1 x_2) - x_1 x_2}{x_1 x_2} = 1 + \varepsilon_1 + \varepsilon_2 + \alpha \right.$$

For division

$$x_1 \leadsto \tilde{x}_1 = rd(x_1) = x_1 (1 + \varepsilon_1)$$

$$x_2 \leadsto \tilde{x}_2 = rd(x_2) = x_2 (1 + \varepsilon_2)$$

$$\frac{x_1}{x_2} \rightarrow rd\left(\frac{\tilde{x}_1}{\tilde{x}_2}\right) = \left(\frac{x_1 (1 + \varepsilon_1)}{x_2 (1 + \varepsilon_2)}\right)(1 + \alpha)$$

$$\begin{array}{l} \text{relative} \\ \text{error} \end{array} = \frac{\frac{x_1}{x_2} - rd\left(\frac{\tilde{x}_1}{\tilde{x}_2}\right)}{\frac{x_1}{x_2}} = \frac{\frac{x_1}{x_2} - \frac{1 + \varepsilon_1}{1 + \varepsilon_2} (1 + \alpha) \left(\frac{x_1}{x_2}\right)}{\frac{x_1}{x_2}}$$

$$= 1 - \frac{1 + \varepsilon_1}{1 + \varepsilon_2} (1 + \alpha) = 1 - \frac{1 + \varepsilon_1}{1 + \varepsilon_2} - \alpha \frac{1 + \varepsilon_1}{1 + \varepsilon_2}$$

