

$$1 \quad 477 = 111011101_2 \\ = (1.11011101)_2 \times 2^8$$

$$\begin{array}{r} 1 \overline{) 0.6} \\ 1 \overline{) 1.2} \\ 0 \overline{) 0.4} \end{array}$$

$$\begin{array}{r} 1 \overline{) 477} \\ 0 \overline{) 238} \end{array}$$

$$2 \quad \frac{2}{5} = \frac{4}{10} = 0.4 \\ 0.4 = (0.100110011001 \dots)_2 \\ = (1.00110011 \dots)_2 \times 2^{-1}$$

$$\begin{array}{r} 0 \overline{) 0.8} \\ 1 \overline{) 1.6} \\ 1 \overline{) 1.2} \end{array}$$

$$\begin{array}{r} 1 \overline{) 59} \\ 1 \overline{) 29} \end{array}$$

$$\begin{array}{r} 0 \overline{) 0.4} \\ 0 \overline{) 0.6} \\ 1 \overline{) 1.6} \end{array}$$

$$\begin{array}{r} 0 \overline{) 14} \\ 1 \overline{) 7} \end{array}$$

$$\begin{array}{r} 0 \overline{) 0.6} \\ 1 \overline{) 1.6} \end{array}$$

$$\begin{array}{r} 1 \overline{) 3} \\ 1 \overline{) 1} \end{array}$$

3 $x = \beta^e, e \in \mathbb{Z}, L < e < U, x_R - x = \beta(x - x_L)$
 Let $x = (1.000 \dots 00)_\beta \times \beta^e$
 $x_L = [(\beta - 1) (\beta - 1) (\beta - 1) \dots (\beta - 1)]_\beta \times \beta^{e-1}$
 $x_R = (1.000 \dots 01)_\beta \times \beta^e$
 $x_R - x = (0.000 \dots 01)_\beta \times \beta^e$
 $= \beta^{e-p}$
 $x - x_L = (0.000 \dots 01)_\beta \times \beta^{e-1}$
 $= \beta^{e-p}$
 $\therefore x_R - x = \beta(x - x_L)$

4 $x = (1.00110011 \dots)_2 \times 2^{-1}$
 $x_L = (1.0011001100110011001100)_2 \times 2^{-1}$
 $x_R = (1.0011001100110011001101)_2 \times 2^{-1}$
 $x - x_L = (1.1001100 \dots)_2 \times 2^{-23} = \frac{8}{5} \times 2^{-23}$
 $x_R - x_L = 2^{-22}$
 $x_R - x = (x_R - x_L) - (x - x_L)$
 $= 2^{-22} - \frac{8}{5} \times 2^{-23}$
 $= \frac{2}{5} \times 2^{-23}$

$\therefore x - x_L > x_R - x, f(x) = x_R$

Relative roundoff error is $\frac{|x_R - x|}{|x|} = \frac{2}{5} \times 2^{-23}$

5. from $\epsilon_M = \beta^{1-p}$ and IEEE 754 single-precision protocol,
 $\beta = 2, p = 24, \epsilon_M = \beta^{-23}$
 $\therefore \epsilon_u = \epsilon_M = 2^{-23}$

6. from Theorem 4.49, $1 > \cos x, x = \frac{1}{4}$

$$\beta^{-6} \leq 1 - \frac{\cos x}{1} \leq \beta^{-5}$$

$$2^{-6} \leq 0.0310875... \leq 2^{-5}$$

\therefore When $\beta = 2$, $1 - \cos x$ lost at most 6 and at least 5 significant digits.

7. Two ways compute $1 - \cos x$ to avoid catastrophic cancellation

(1) Taylor's expansion:

$$1 - \cos x = 1 - \left(1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots\right)$$

$$= \frac{x^2}{2!} - \frac{x^4}{4!} + \frac{x^6}{6!} - \dots$$

(2) Trigonometric identity

$$1 - \cos x = 2 \sin^2 \frac{x}{2}$$

8 (1) $(x-1)^\alpha$

$$\alpha \neq 0, C_f(x) = \left| \frac{x f'(x)}{f(x)} \right|$$

$$= \left| \frac{x \alpha (x-1)^{\alpha-1}}{(x-1)^\alpha} \right|$$

$$= \left| \frac{\alpha x}{x-1} \right|$$

\therefore When $x \rightarrow 1, C_f(x) \rightarrow +\infty$

$\alpha = 0, C_f(x) = 0$, never be large

(2) $\ln x$

$$C_f(x) = \left| \frac{1}{\ln x} \right|$$

$\therefore x \rightarrow 0^+, C_f(x) \rightarrow +\infty$

(3) e^x

$$C_f(x) = \left| \frac{x e^x}{e^x} \right| = |x|$$

$\therefore x \rightarrow +\infty, C_f(x) \rightarrow +\infty$

(4) $0 < \cos x$

$$L_f(x) = \left| -\frac{x}{\sqrt{1-x^2} \arccos x} \right|$$

$\therefore x = \pm 1, L_f(x) \rightarrow +\infty$

94) $f(x) = 1 - e^{-x}, x \in [0, 1]$

$$L_f(x) = \left| \frac{x e^{-x}}{1 - e^{-x}} \right|$$

$$= \left| \frac{x}{e^x - 1} \right|$$

$\lim_{x \rightarrow 0} L_f(x) = 1, L_f(x) \text{ decrease in } x \in [0, 1]$

$\therefore L_f(x) \leq 1$

● (2) $|f(x_A) - f(x)| = |f'(s)| |x - x_A| < \epsilon_n, s \text{ between } x_A \text{ and } x$

$$|x - x_A| < \frac{\epsilon_n}{|f'(s)|} = \frac{\epsilon_n}{e^{-s}} \leq e \epsilon_n$$

$$\text{Cond}_A(x) = \frac{1}{\epsilon_n} \min_{\{x_A\}} \frac{|x - x_A|}{x}$$

$$< \frac{e}{x}$$

(3)

10) $q(x) = \sum_{i=0}^n a_i x^i, a_n = 1, a_0 \neq 0, a_i \in \mathbb{R}$

$$q(r) = \sum_{i=0}^n a_i r^i$$

$$r^j + \sum_{i=0}^{j-1} i a_i r^i + \frac{\partial r}{\partial a_j} = 0, j = 1, 2, \dots, n-1$$

$$\nabla r = - \frac{(1, r, \dots, r^{n-1})}{q'(r)}$$

$$\text{Cond}_r(a) = \frac{\|a\|_1 \|\nabla r\|_1}{|r|}$$

$$= \frac{\sum_{i=0}^{n-1} |a_i| \sum_{i=0}^{n-1} |r^i|}{|r|^n}$$

In Wilkinson example, $q(x) = \prod_{i=1}^n (x - i)$

$$\sum_{i=0}^{n-1} |a_i r^i| \geq - \sum_{i=0}^{n-1} a_i r^i = r^n = n^n$$

$$\sum_{i=0}^{n-1} |a_i r^i| = |r| |q'(r)|$$

$$= n |q'(n)|$$

$$< n^2 n!$$

11 $\beta=2, p=2, L=-1, U=1$

$$a = (1.0)_2 \times 2^0, b = (1.1)_2 \times 2^0$$

$$\frac{a}{b} = (0.101010\dots)_2$$

$$f\left(\frac{a}{b}\right) = \sqrt{2}[(0.101)_2]$$

$$= (0.10)_2 \times 2^0$$

$$= 0.5$$

$$\text{Relative error} = \left| \frac{0.5}{\frac{1}{\sqrt{2}}} - 1 \right| = 0.25$$

$$\text{Machine precision: } \epsilon_u = 2^{-2} = 0.25$$

12 $\epsilon_M = 2^{-23}$

$$2^7 \epsilon_M = 2^{-16} \approx 1.5259 \times 10^{-5} > 2 \times 10^{-6}$$

\therefore We can't compute the root with absolute accuracy $< 10^{-6}$

13 calculate $s(x) = ax^3 + bx^2 + cx + d$ on $[x_i, x_{i+1}]$ by $s(x_i), s(x_{i+1}), s'(x_i), s'(x_{i+1})$.

$$\text{Coefficient matrix: } \begin{pmatrix} x_i^3 & x_i^2 & x_i & 1 \\ x_{i+1}^3 & x_{i+1}^2 & x_{i+1} & 1 \\ 3x_i^2 & 2x_i & 1 & 0 \\ 3x_{i+1}^2 & 2x_{i+1} & 1 & 0 \end{pmatrix}$$

It has large condition number when x_i and x_{i+1} are close enough, so the result will be very inaccurate.