

I. Convert the decimal integer 477 to a normalized FPN with $\beta = 2$.

$$477 = 2^8 + 2^7 + 2^5 + 2^4 + 2^3 + 2^2 + 2 + 1 = (110111111)_2 = (1.10111111)_2 \times 2^8.$$

II. Convert the decimal fraction $3/5$ to a normalized FPN with $\beta = 2$.

$$\begin{aligned}\frac{3}{5} \times 2 &= \frac{6}{5} = 1 + \frac{1}{5} \\ \frac{1}{5} \times 2 &= \frac{2}{5} \\ \frac{2}{5} \times 2 &= \frac{4}{5} \\ \frac{4}{5} \times 2 &= \frac{8}{5} = 1 + \frac{3}{5}\end{aligned}$$

So we have $\frac{3}{5} \times 2^4 = (1001)_2 + \frac{3}{5}$, thus $\frac{3}{5} = (0.10011001\dots)_2 = (1.00110011\dots) \times 2^{-1}$.

III. Prove $x_R - x = \beta(x - x_L)$.

Since $x = 1 \times \beta^e$, we have

$$\begin{aligned}x_R &= \left(1 + \frac{1}{\beta^{p-1}}\right) \times \beta^e \\ x_L &= \left(\beta - \frac{1}{\beta^{p-1}}\right) \times \beta^{e-1} = [(\beta - 1) + \frac{\beta - 1}{\beta} + \dots + \frac{\beta - 1}{\beta^{p-1}}] \times \beta^{e-1}.\end{aligned}$$

So we have the equation

$$(x_R - x) = \frac{1}{\beta^{p-1}} \times \beta^e = \frac{1}{\beta^{p-1}} \times \beta^{e-1} \times \beta = \beta(x - x_L).$$

IV. Find out the two normalized FPNs adjacent to $x = 3/5$ under the IEEE 754 single-precision protocol. What is $fl(x)$ and the relative roundoff error?

By problem II, we have $\frac{3}{5} = (1.00110011\dots)_2 \times 2^{-1}$, under the IEEE 754 single-precision protocol, $p = 24$, so we have

$$\begin{aligned}x_L &= 2^{-1} \times (1.00110011001100110011001)_2 \\ x_R &= 2^{-1} \times (1.00110011001100110011010)_2.\end{aligned}$$

Then calculate $3/5 - x_L$ and $x_R - 3/5$, we have

$$\begin{aligned}\frac{3}{5} - x_L &= 2^{25} \times (1.00110011\dots)_2 \\ x_R - \frac{3}{5} &> x_R - 2^{-1} \times (1.0011001100110011001100111)_2 = 2^{-25} \times (10.1)_2 > \frac{3}{5} - x_L.\end{aligned}$$

So $fl(\frac{3}{5}) = x_L = 2^{-1} \times (1.00110011001100110011001)_2$, round error $e = 2^{-25} \times (1.00110011\dots)_2 = 2^{-25} \times 3/5$.

V. What would the unit roundoff be if the IEEE 754 single-precision protocol dropped excess bits?

$$\epsilon_u = \max_{x \in \mathcal{R}(\mathcal{F})} \frac{|fl(x) - x|}{\beta^{e_x}} = \frac{(1.11\dots)_2 \times 2^{-24} \times \beta^{e_x}}{\beta^{e_x}} = 2^{-23}.$$

VI. How many bits of precision are lost in the subtraction $1 - \cos x$ when $x = \frac{1}{4}$?

Because $2^{-6} \leq 1 - \cos(\frac{1}{4}) \leq 2^{-5}$, subtraction $1 - \cos(\frac{1}{4})$ will lost at least 5, at most 6 bits of precision.

VII. Suggest at least two ways to compute $1 - \cos x$ to avoid catastrophic cancellation caused by subtraction.

To avoid catastrophic cancellation, we can use multiplication to calculate $1 - \cos(x)$ that is

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

Since multiplication is accurate, we can avoid catastrophic cancellation.

The second method is to avoid the result of addition close to 0. By using Taylor series we can avoid catastrophic cancellation, that is

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