

Numerical Analysis

Spring 2025

Homework #2 Math and Floating-point Number Representation

Instructor: *Kwangki Kim*

Name : 장준영

Student ID# : 12191529

Problem 1.1 (FORMAL TAYLOR SERIES FOR f ABOUT c) What is the Taylor series of the function

$$f(x) = 3x^5 - 2x^4 + 15x^3 + 13x^2 - 12x - 5$$

at the point $c = 2$?

$$\text{sol.) } f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(c)}{n!} (x-c)^n$$

$$\textcircled{1} f(2) = 3 \cdot 32 - 2 \cdot 16 + 15 \cdot 8 + 13 \cdot 4 - 12 \cdot 2 - 5 = 207$$

$$\textcircled{2} f'(x) = 15x^4 - 8x^3 + 45x^2 + 26x - 12$$

$$\rightarrow f'(2) = 15 \cdot 16 - 8 \cdot 8 + 45 \cdot 4 + 26 \cdot 2 - 12 = 396$$

$$\textcircled{3} f''(x) = 60x^3 - 24x^2 + 90x + 26$$

$$\rightarrow f''(2) = 60 \cdot 8 - 24 \cdot 4 + 90 \cdot 2 + 26 = 590$$

$$\textcircled{4} f^{(3)}(x) = 180x^2 - 48x + 90$$

$$\rightarrow f^{(3)}(2) = 180 \cdot 4 - 48 \cdot 2 + 90 = 714$$

$$\textcircled{5} f^{(4)}(x) = 360x - 48$$

$$\rightarrow f^{(4)}(2) = 360 \cdot 2 - 48 = 672$$

$$\textcircled{6} f^{(5)}(x) = 360$$

$$\rightarrow f^{(5)}(2) = 360$$

$$\textcircled{7} f^{(6)}(x) = 0$$

⋮

$$\Rightarrow f_T(x) = f(2) + f'(2)(x-2) + \frac{f''(2)}{2!}(x-2)^2 + \frac{f^{(3)}(2)}{3!}(x-2)^3 + \frac{f^{(4)}(2)}{4!}(x-2)^4 + \frac{f^{(5)}(2)}{5!}(x-2)^5$$

$$= 207 + 396(x-2) + \frac{590}{2!}(x-2)^2 + \frac{714}{3!}(x-2)^3 + \frac{672}{4!}(x-2)^4 + \frac{360}{5!}(x-2)^5$$

Problem 1.2 (TAYLOR'S THEOREM FOR $f(x)$)

- (a) Derive the Taylor series for e^x at $c = 0$, and prove that it converges to e^x by using Taylor's Theorem.
- (b) Derive the formal Taylor series for $f(x) = \ln(1+x)$ at $c = 0$, and determine the range of positive x for which the series represents the function.

(a) $f(x) = e^x$ 74 349,

$$f'(x) = e^x ; f''(x) = e^x \dots$$

$$f_T(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} x^n = \sum_{n=0}^{\infty} \frac{1}{n!} x^n$$

$$(\text{Taylor}) R_{n+1}(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} x^{n+1} \quad (\because \xi \in (0, x))$$

$$\text{74741} \quad f^{(n+1)}(\xi) = e^{\xi} \leq e^{|x|}$$

$$\text{74741} \quad |R_{n+1}(x)| \leq \frac{e^{|x|} x^{n+1}}{(n+1)!}$$

$$\therefore \lim_{n \rightarrow \infty} \frac{e^{|x|} x^{n+1}}{(n+1)!} = 0 \quad \text{74741} \quad e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

(b) $f(x) = \ln(1+x)$

$$f'(x) = \frac{1}{1+x} \rightarrow f'(0) = 1$$

$$f''(x) = -\frac{1}{(1+x)^2} \rightarrow f''(0) = -1$$

$$f^{(3)}(x) = \frac{2}{(1+x)^3} \rightarrow f^{(3)}(0) = 2$$

\vdots

$$f^{(n)}(x) = \frac{(-1)^{n-1}}{(1+x)^n} (n-1)! \rightarrow f^{(n)}(0) = (-1)^{n-1} (n-1)!$$

$$\therefore f_T(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} x^n = \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} x^n$$

$$\therefore (\text{74741}) : 0 < x \leq 1$$

Problem 1.3 (A technique of calculating the number of terms to be used in a series by just making the $(n+1)$ st term less than some tolerance) If the logarithmic series

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \cdots = \sum_{k=1}^{\infty} (-1)^{k-1} \frac{x^k}{k} \quad (-1 < x \leq 1)$$

is to be used for computing $\ln 2$ with an error of less than $\frac{1}{2} \times 10^{-6}$, how many terms will be required?

$$\text{Sol)} \ln(1+x) = \sum_{k=1}^n (-1)^{k-1} \frac{x^k}{k} + (-1)^n \frac{x^{n+1}}{n+1}$$

$$\therefore R_{n+1}(x) = \left| (-1)^n \frac{x^{n+1}}{n+1} \right| = \left| \frac{x^{n+1}}{n+1} \right| \quad (-1 < x \leq 1)$$

$$\leq \frac{1}{n+1} < \frac{1}{2} \times 10^{-6} \quad ; \quad n > 1,999,999$$

답: 2,000,000 개 이상

Problem 1.4 (*Floating-Point Representation and Errors*) If x, y , and z are machine numbers in a 32-bit word-length computer, what upper bound can be given for the relative roundoff error in computing $z(x + y)$?

sol) 32 bit = 1 (sign) + 8 (exponent) + 23 (mantissa)

1개의 hidden bit 가 존재한다고 가정하자.

총 유효숫자의 개수는 24개 이므로, 1개 연산에서의 오차 $\epsilon \leq 2^{-24}$

연산은 총 2번 이루어진다:

$$\textcircled{1} x+y \rightarrow |\epsilon_1| \leq 2^{-24}$$

$$\textcircled{2} z(x+y) \rightarrow |\epsilon_2| \leq 2^{-24}$$

따라서 전체 상대오차는 미약 $|\epsilon| \leq 2 \cdot 2^{-24} = 2^{-23}$ 의 upper bound를

갖는다고 볼 수 있다.

Problem 1.5 (*Floating-Point Representation and Errors*) Critique the following attempt to estimate the relative roundoff error in computing the sum of two real numbers, x and y . In a 32-bit word-length computer, the calculation yields

$$\begin{aligned} z &= \text{fl}[\text{fl}(x) + \text{fl}(y)] \\ &= [x(1 + \delta) + y(1 + \delta)](1 + \delta) \\ &= (x + y)(1 + \delta)^2 \\ &\approx (x + y)(1 + 2\delta) \end{aligned}$$

Therefore, the relative error is bounded as follows:

$$\left| \frac{(x + y) - z}{(x + y)} \right| = \left| \frac{2\delta(x + y)}{(x + y)} \right| = |2\delta| \leq 2^{-23}$$

Why is this calculation not correct?

sol) 위 계산이 정확하지 않은 이유는 다음과 같다.

■ 모든 Roundoff Error를 하나의 δ 로 가정함.

■ 그러나, 현실에서는 연산의 각 단계마다 다른 오차가 발생한다.

$$\bullet \text{fl}(x) = x(1 + \delta_1)$$

$$\bullet \text{fl}(y) = y(1 + \delta_2)$$

$$\bullet \text{fl}(\text{sum}) = \{x(1 + \delta_1) + y(1 + \delta_2)\}(1 + \delta_3) = z$$

$$\text{즉, } z = \{x(1 + \delta_1) + y(1 + \delta_2)\}(1 + \delta_3) = \{(x + y) + (\delta_1 x + \delta_2 y)\}(1 + \delta_3)$$

따라서 상대오차는 :

$$\left| \frac{(x + y) - z}{(x + y)} \right| = \left| \frac{(x + y)\delta_3 + \delta_1 x(1 + \delta_3) + \delta_2 y(1 + \delta_3)}{(x + y)} \right| = \left| 1 + \frac{\delta_1 x + \delta_2 y}{x + y}(1 + \delta_3) \right|$$

$$\blacksquare |\delta_i| < \epsilon : |\delta_1 x + \delta_2 y| \leq \epsilon(|x| + |y|)$$

$$\therefore \left| 1 + \frac{\delta_1 x + \delta_2 y}{x + y}(1 + \delta_3) \right| \leq \epsilon + \epsilon(1 + \epsilon) \cdot \frac{|x| + |y|}{|x + y|}$$

■ 일반적인 경우 : $x \approx y$

$$\therefore \epsilon + \epsilon(1 + \epsilon) \cdot \frac{2x}{2x} = \epsilon + \epsilon + \epsilon^2 \approx 2\epsilon$$

Problem 1.7 (Exactly how many significant binary digits are lost in the subtraction $x - y$ when x is close to y ?) In the subtraction $37.593621 - 37.584216$, how many bits of significance will be lost?

(Hint) Let x and y be normalized floating-point machine numbers, where $x > y > 0$. If $2^{-p} \leq 1 - (y/x) \leq 2^{-q}$ for some positive integers p and q , then at most p and at least q significant binary bits are lost in the subtraction $x - y$. **Proof:** We prove the second part of the theorem and leave the first as an exercise. To this end, let $x = r \times 2^n$ and $y = s \times 2^m$, where $\frac{1}{2} \leq r, s < 1$. (This is the normalized binary floating-point form.) Since $y < x$, the computer may have to shift y before carrying out the subtraction. In any case, y must first be expressed with the same exponent as x . Hence, $y = (s2^{m-n}) \times 2^n$ and

$$x - y = (r - s2^{m-n}) \times 2^n$$

The mantissa of this number satisfies the equations and inequality

$$r - s2^{m-n} = r \left(1 - \frac{s2^m}{r2^n} \right) = r \left(1 - \frac{y}{x} \right) < 2^{-q}$$

Hence, to normalize the representation of $x - y$, a shift of at least q bits to the left is necessary. Then at least q (~~spurious~~^{가짜}) zeros are supplied on the right-hand end of the mantissa. This means that at least q bits of precision have been lost. ■

sol) $x = 37.593621$, $y = 37.584216$ 이라 하자.

$$1 - \frac{y}{x} = 1 - \frac{37.584216}{37.593621} = 0.0002501754$$

$$\therefore 2^{-p} \leq 0.0002501754 \leq 2^{-q}$$

$$-\log_2(0.0002501754) \approx 11.96 \quad \text{이므로} \quad p=12, q=11$$

답: 최소 11 비트, 최대 12비트 손실