

- 1) Representing reals in a floating point number system
- 2) rounding error / relative error
- 3) stability of formulae

①  $R_b(t, s)$

$\uparrow$  base       $\uparrow$  size of mantissa       $\uparrow$  size of exponent

ex: Represent  $(72.6)_{10}$  in the FP # system

$b=3$        $t=8$        $e = [-2, +4]$

$\uparrow$  range of exponent

$(72)_{10} \rightarrow (2200)_3$

numerator	denominator	quotient	remainder
72	3	24	0
24	3	8	0
8	3	2	2
2	3	0	2

Stop when quotient is 0

$$(.6)_{10} \rightarrow (. \overline{1210})_3 \quad \leftarrow \text{repeats infinitely}$$

	multiplier	base	product	integral	fraction
Repeats!	.6	3	1.8	1	.8
	.8	3	2.4	2	.4
	.4	3	1.2	1	.2
	.2	3	0.6	0	.6
	.6	3	1.8	1	.8

Stop when fraction is 0 or cycle restarts

$$\therefore (72.6)_{10} = (2200.\overline{1210})_3$$

$$\Rightarrow (\underbrace{.2200\overline{1210}}_{\text{mantissa}})_3 \times \underbrace{3^4}_{\text{base}} \quad \begin{array}{l} \uparrow \\ 4 \rightarrow (11)_3 \\ \uparrow \\ \text{size of exponent is 2} \\ \uparrow \\ \text{exponent} \end{array}$$

2) Absolute error (AE) =  $x - f_1(x)$

Relative error (RE) =  $\frac{x - f_1(x)}{x}$

$$x = (0.2200\overline{1210})_3 \times 3^4$$

$$f_1(x) = (0.22001210)_3 \times 3^4$$

$$AE = (0.00000000\overline{1210})_3 \times 3^4$$

3)  $1 - \cos x$

when  $x \rightarrow 0$ ,  $\cos x \rightarrow 1$  the formula suffers from

catastrophic / subtractive cancellation

$$\text{Alternative: } 1 - \cos x \times \frac{1 + \cos x}{1 + \cos x}$$

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

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

$$[\sin^2 x + \cos^2 x = 1]$$

This formula is fine when  $x \rightarrow 0$ , but not when  $x \rightarrow \pi$  (since  $\cos x \rightarrow -1$ )

$1 - \cos x$  is fine as  $x \rightarrow \pi$

Solution: Depending on the value of  $x$ , choose a formula that doesn't suffer from catastrophic / subtractive cancellation

$$\left\{ \begin{array}{l} 1 - \cos x \quad x \rightarrow \pi \\ \frac{\sin^2 x}{1 + \cos x} \quad x \rightarrow 0 \end{array} \right.$$