

~~1. Condition~~

4. Condition

$$(ComA)f(x) = \left| \frac{x f'(x)}{x} \right|$$

a)  $f(x) = 1 - e^{-x}$   $f'(x) = e^{-x}$

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

Now  $e^x - 1 = \sum_{n=1}^{\infty} \frac{x^n}{n!}$ , so  $e^x - 1 > x$  on  $[0, 1]$

$\lim_{x \rightarrow 0} (Comf)(x) < 1$  on  $[0, 1]$

Wie  $(Comf)(x) = \left( \frac{e^x - 1}{x} \right)^{-1} = \left( \sum_{n=1}^{\infty} \frac{x^{n-1}}{n!} \right)^{-1} = \frac{1}{1 + \frac{x + x^2}{2} + \dots} < 1$  on  $[0, 1]$

b)  $f_L(e^{-x}) = e^{-x} (1 + \epsilon)$

$$\begin{aligned} f_L(1 - f_L(e^{-x})) &= f_L(1 - e^{-x} (1 + \epsilon)) \\ &= f_L((1 - e^{-x}) (1 - \frac{e^{-x} \epsilon}{1 - e^{-x}})) \\ &= (1 - e^{-x}) (1 - \frac{\epsilon}{e^x - 1}) (1 + \epsilon) \\ &= (1 - e^{-x}) (1 - \frac{\epsilon}{e^x - 1}) \end{aligned}$$

$$f_A(x) = f_L(1 - f_L(e^{-x})) = (1 - e^{-x}) (1 - \frac{e^{-x} \epsilon}{e^x - 1})$$

$$f(x_A) = 1 - e^{-x_A} = (1 - e^{-x}) (1 - \frac{e^{-x} \epsilon}{e^x - 1})$$

$1 - e^{-x_A} = 1 - e^{-x} - \epsilon \Rightarrow e^{-x_A} = e^{-x} + \epsilon = e^{-x} (1 + e^x \epsilon)$   
 $\Rightarrow x_A = -\ln(1 + e^x \epsilon) \approx x - e^x \epsilon$  (since  $\ln(1 + \epsilon) \approx \epsilon$ )

$x_A = x - e^x \epsilon$

$|x_A - x| = e^x \epsilon$

$\lim_{\epsilon \rightarrow 0} |x_A - x| \leq e^x \epsilon$

$(ComA)(x) = \frac{1}{e^x} \frac{|x_A - x|}{x}$

$(ComA)(x) = \frac{e^x}{x}$

$\lim_{x \rightarrow 0} (ComA)(x) = \frac{e^x}{x} = \sum_{n=0}^{\infty} \frac{x^{n-1}}{n!}$   
 $\lim_{x \rightarrow 0} (ComA)(x) > 1$  on  $[0, 1]$  since the first term is  $\frac{1}{x}$ , which is always  $> 1$  on  $[0, 1]$ , and all others are positive

→ from b)  $f_A(x) = (1 - e^{-x}) \left(1 - \frac{\epsilon}{1 - e^{-x}}\right)$

so lose bits:  $\frac{\epsilon}{1 - e^{-x}} \rightarrow$  <sup>at most</sup>  $\frac{\epsilon}{1 - e^{-x}} \leq 2^b$  if  
 since the 1st  $\epsilon$  multiplies  $\epsilon$  and the carrying error, first bit is at  $2\epsilon$

so if we are willing to lose at most  $b$  bits, the minimum  $X$  satisfies

$$\frac{1}{1 - e^{-x}} = 2^b$$

$$1 - 2^{-b} = e^{-x}$$

$$X_{\min} = -\ln(1 - 2^{-b})$$

(f) → avoid subtraction & evaluation

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

this way around is ok, but should not pose particular problems for small  $x$

so for 2 bits:  $X_{\min}^{(b=2)} = -\ln(1 - \frac{1}{2}) = \ln(2)$

so for 1 bit of precision left, the minimum  $X = \ln(2)$

2 bits:  $X_{\min}^{(b=2)} = \ln(1 - \frac{1}{4}) = \ln(\frac{4}{3})$

3 bits:  $X_{\min}^{(b=3)} = \ln(1 - \frac{1}{8}) = \ln(\frac{8}{7})$

4 bits:  $X_{\min}^{(b=4)} = \ln(1 - \frac{1}{16}) = \ln(\frac{16}{15})$

e) relative error  $\tilde{\epsilon} \in \frac{|f_A(x) - f(x)|}{f(x)} \leq \frac{\frac{\epsilon}{1 - e^{-x}}}{\frac{e^x - 1}{e^x}} \leq \frac{\epsilon}{1 - \exp(-2^{-b})} \leq \frac{\epsilon}{2^{-b}}$

so the relative error  $\tilde{\epsilon}$  is bounded by

$$\tilde{\epsilon} \leq 2^b \epsilon$$

as required

$X = \ln(2), b=1: \tilde{\epsilon} \leq 2\epsilon$

$X = \ln(\frac{4}{3}), b=2: \tilde{\epsilon} \leq 4\epsilon$

$X = \ln(\frac{8}{7}), b=3: \tilde{\epsilon} \leq 8\epsilon$

$X = \ln(\frac{16}{15}), b=4: \tilde{\epsilon} \leq 16\epsilon$