

Sam Moore

(1)

$$x = \pm \left(\sum_{\ell=1}^{\infty} b_{-\ell} 2^{-\ell} \right) \cdot 2^e, \quad b_i \in (0, 1)$$

If $b_{p+1} = 0$, the rounding is akin to truncation:

$$x - \text{trunc}(x) = \pm \left(\sum_{\ell=p+1}^{\infty} b_{-\ell} 2^{-\ell} \right) \cdot 2^e$$

$$\Rightarrow \max |x - \text{trunc}(x)| = \left| \sum_{\ell=p+1}^{\infty} 2^{-\ell} \right| \cdot 2^e \quad [b_{p+1} = 0, \text{ not } 1]$$

$$= \left(2^{-(p+2)} + 2^{-(p+3)} + \dots \right) \cdot 2^e$$

$$\Rightarrow \max \left| \frac{x - \text{trunc}(x)}{x} \right| = \frac{2^{-p-1} \cdot 2^e}{2^{-1} \cdot 2^e} = 2^{-p}$$

$$\Rightarrow \left| \frac{x - \text{trunc}(x)}{x} \right| \leq 2^{-p}$$

If $b_{p+1} = 1$, add 1 to b_p , and then truncate.

$$X = \pm \left(\sum_{\ell=1}^{\infty} b_{-\ell} 2^{-\ell} \right) \cdot 2^e$$

$$\Rightarrow X^* = \pm \left(b_{-1} 2^{-1} + b_{-2} 2^{-2} + \dots + (b_p + 1) 2^{-p} + b_{p+1} 2^{-(p+1)} + \dots \right) \cdot 2^e$$

$$\Rightarrow \max |X^* - \text{trunc}(X^*)| = \left| \sum_{\ell=p+1}^{\infty} 2^{-\ell} \right| \cdot 2^e = (2^{-p-1} + 2^{-p-2} + \dots)$$

$$\Rightarrow \max \left| \frac{X^* - \text{trunc}(X^*)}{X^*} \right| = \frac{(2^{-p} \cdot 2^{-p-1}) 2^e}{2^{-1} \cdot 2^e} = 2^{-p}$$

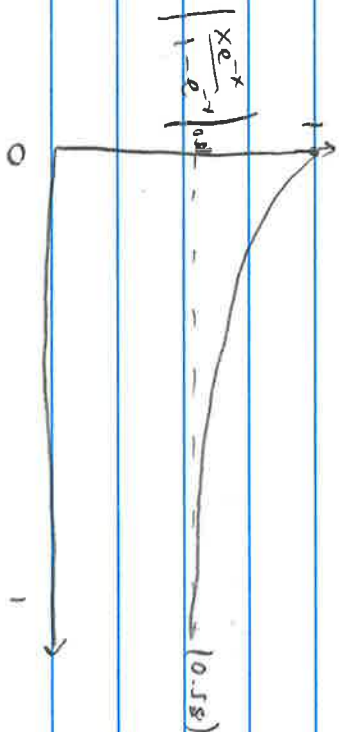
Thus $\left| \frac{x - \text{rd}(x)}{x} \right| \leq 2^{-p}$ by combining both cases

4 a) $\text{Cond } f(x) = \left\| \frac{x f'(x)}{f(x)} \right\|$

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

$\Rightarrow (\text{Cond } f)(x) = \left\| \frac{x e^{-x}}{1 - e^{-x}} \right\|$

Noting that $\lim_{x \rightarrow 0} \left\| \frac{x e^{-x}}{1 - e^{-x}} \right\| = 1$, then it is clear that on the interval $[0, 1]$, $\left\| \frac{x e^{-x}}{1 - e^{-x}} \right\| < 1$



\Rightarrow the problem is well-conditioned.

This method
is roughly correct,
but wasn't working
for me

b) First, note that negation yields no error. Then, applying Step 2, $f_A(x) = 1.0 - e^{-x} (1 + \epsilon_{\text{exp}})$, where we note the error yielded by the exp function, and so that as $1.0 \in \mathbb{R}(\mu_9)$, there is no round-off error introduced here).

From lecture: $|\epsilon_{x+y}| \leq \left| \frac{x}{x+y} \right| |\epsilon_x| + \left| \frac{y}{x+y} \right| |\epsilon_y|$
in this case $\Rightarrow |\epsilon_{x+y}| \leq \left| \frac{e^{-x}}{1 - e^{-x}} \right| \epsilon_{\text{exp}}$

Then: $\left(1.0 + \frac{e^{-x}}{1 - e^{-x}} \epsilon_{\text{exp}} \right) (1 + \epsilon_{\text{rd}}) \approx 1 - e^{-x_A}$

$\approx \left(1 + \frac{e^{-x}}{1 - e^{-x}} \epsilon_{\text{exp}} + \epsilon_{\text{rd}} + \frac{e^{-x}}{1 - e^{-x}} \epsilon_{\text{exp}} \epsilon_{\text{rd}} \right) = 1 - e^{-x_A}$
 $\approx \left(1 + \frac{\epsilon_{\text{exp}}}{e^{x-1}} (1 + \epsilon_{\text{rd}}) + \epsilon_{\text{rd}} \right) \approx 1 - e^{-x_A}$

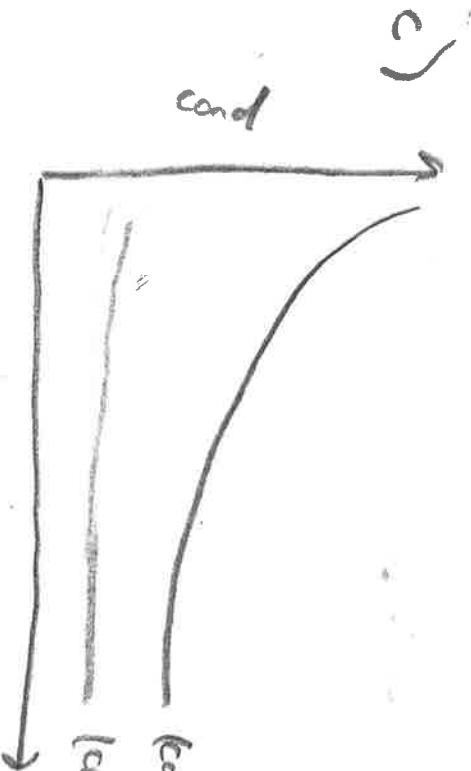
$\Rightarrow x_A = \ln(1 + \delta) \approx 4\delta + O(\delta^2)$ customarily δ_{small}

$$4b) f_A(x) = [1 - e^{-x}(1+\varepsilon)](1+\varepsilon) \\ \approx (1 - e^{-x})\left(1 - \frac{e^{-x}}{1-e^{-x}}\varepsilon\right) \approx (1 - e^{-x})\left(1 + \frac{\varepsilon}{1-e^{-x}}\right)$$

$$\Rightarrow |f(x_A) - f(x)| = f(x) \frac{\varepsilon}{1-e^{-x}} \approx |f'(x)| |x - x_A|$$

$$\Rightarrow \left| \frac{x - x_A}{x} \right| \approx \left| \frac{f(x)}{x f'(x)} \frac{\varepsilon}{1-e^{-x}} \right| \quad \leftarrow \text{or something like this}$$

$$\Rightarrow (C_{\text{cond}} A)(x) \approx \frac{1}{\varepsilon} \left| \frac{x - x_A}{x} \right| \approx \frac{e^x}{x}$$



(Using Debnas for
Convenience)

Note that when x is small, $1 - e^{-x}$ gives a large error
[as $x \rightarrow 0, e^{-x} \rightarrow 1$, and the subtraction of two numbers
that are very close gives a big error].

$$d) 1 - e^{-x} < 2^{-6} \Rightarrow x = \ln(1 - 2^{-6})$$

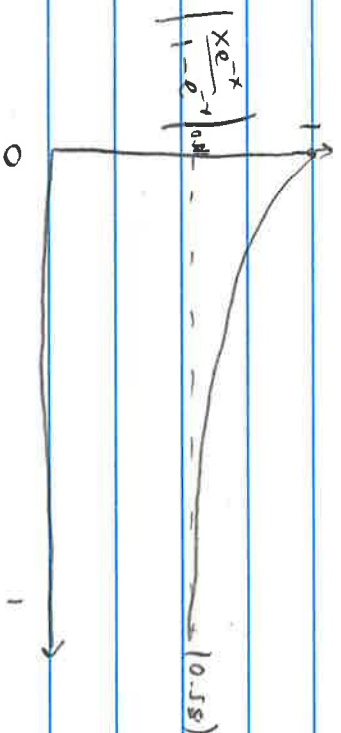
$$\text{rel-err} = \frac{\varepsilon}{1 - e^{-x}} = 2^6 \varepsilon$$

$n=1$	$\rightarrow x = 0.693$	rel-err = 2^{-51}
$n=2$	$\rightarrow x = 0.288$	rel-err = 2^{-49}
$n=3$	$\rightarrow x = 0.134$	rel-err = 2^{-48}
$n=4$	$\rightarrow x = 0.065$	rel-err = 2^{-45}

4 a) $\text{cond } f(x) = \left\| \frac{x f'(x)}{f(x)} \right\|$
 $f(x) = 1 - e^{-x}$; $f'(x) = e^{-x}$

$$\Rightarrow (\text{cond } f)(x) = \left\| \frac{x e^{-x}}{1 - e^{-x}} \right\|$$

Noting that $\lim_{x \rightarrow 0} \left\| \frac{x e^{-x}}{1 - e^{-x}} \right\| = 1$, then it is clear that on the interval $[0, 1]$, $\left\| \frac{x e^{-x}}{1 - e^{-x}} \right\| < 1$



\Rightarrow the problem is well-conditioned.

This method
is roughly correct,
but wasn't working
for me

b) First, note that negation yields no error. Then, applying Step 2, $f_A(x) = 1.0 - e^{-x} (1 + \epsilon_{\text{exp}})$, where we note the error yielded by the exp function, and so that w.l.o.g. $x \in \mathbb{R}^+$, there is no round-off error introduced here.

From technique: $|\epsilon_{x+y}| \leq \left| \frac{x}{x+y} \right| |\epsilon_x| + \left| \frac{y}{x+y} \right| |\epsilon_y|$
 in this case $\Rightarrow |\epsilon_{x+y}| \leq \left| \frac{e^{-x}}{1 - e^{-x}} \right| \epsilon_{\text{exp}}$

Then: $\left(1.0 + \frac{e^{-x}}{1 - e^{-x}} \epsilon_{\text{exp}} \right) (1 + \epsilon_{\text{rd}}) = 1 - e^{-x_A}$

$$= \left(1 + \frac{e^{-x}}{1 - e^{-x}} \epsilon_{\text{exp}} + \epsilon_{\text{rd}} + \frac{e^{-x}}{1 - e^{-x}} \epsilon_{\text{exp}} \epsilon_{\text{rd}} \right) = 1 - e^{-x_A}$$

$$= \left(1 + \frac{\epsilon_{\text{exp}}}{e^x - 1} (1 + \epsilon_{\text{rd}}) + \epsilon_{\text{rd}} \right) = 1 - e^{-x_A}$$

$$\Rightarrow x_A = \ln(1 + \delta) \approx \delta + 0(\delta^2) \quad \text{assuming } \delta \text{ small}$$

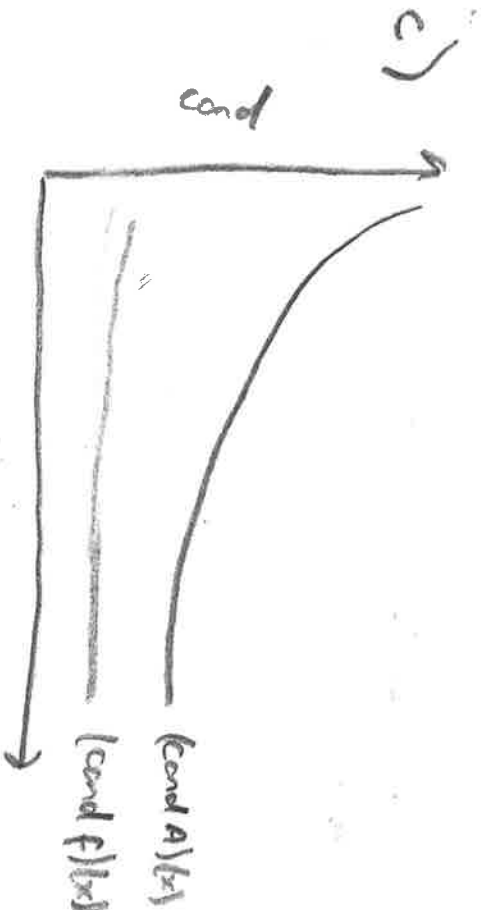
$$4b) f_A(x) = [1 - e^{-x}(1+\varepsilon)]/(1+\varepsilon) \\ \approx (1 - e^{-x})(1 - \frac{e^{-x}}{1-e^{-x}} \varepsilon) \approx (1 - e^{-x})(1 + \frac{2}{1-e^{-x}}x)$$

$$\Rightarrow |f(x_A) - f(x)| = f(x) \frac{\varepsilon}{1-e^{-x}} \approx |f'(x)| |x - x_A|$$

$$\Rightarrow \left| \frac{x - x_A}{x} \right| \approx \left| \frac{f(x)}{x f'(x)} \frac{\varepsilon}{1-e^{-x}} \right|$$

or something like
close to this

$$\Rightarrow |(Cond A)(x)| \approx \frac{1}{\varepsilon} \left| \frac{x - x_A}{x} \right| \approx \frac{e^x}{x}$$



(Using DeMorgan's for
Convenience)

Note that when x is small, $1 - e^{-x}$ gives a large error
[as $x \rightarrow 0, e^{-x} \rightarrow 1$, and the subtraction of two numbers
that are very close gives a big error].

$$d) 1 - e^{-x} < 2^{-b} \Rightarrow x = \ln(1 - 2^{-b})$$

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

$n=1$	$\rightarrow x=0.693$	$rel-err = 2^{-1}$
$n=2$	$\rightarrow x=0.288$	$rel-err = 2^{-2}$
$n=3$	$\rightarrow x=0.134$	$rel-err = 2^{-3}$
$n=4$	$\rightarrow x=0.065$	$rel-err = 2^{-4}$

$$b) x^a = e^{a \ln x}$$

$$i) e^{a \ln x} = e^{a(1+\varepsilon) \ln x} = e^{a \ln x + a \varepsilon \ln x}$$

$$\approx x^a (1 + a \varepsilon \ln x)$$

$$\Rightarrow \xi = (a \ln x) \varepsilon_a$$

\rightarrow if $a \ln x \gg 1$,
 ξ is significant.

$$ii) \ln x = \ln(x(1+\varepsilon)) = \ln x + \varepsilon_x$$

$$\Rightarrow e^{a \ln x} = \exp[a \ln x + a \varepsilon_x]$$

$$\Rightarrow \xi = a \varepsilon_x \quad \rightarrow \quad x^a (1 + a \varepsilon_x)$$

\rightarrow if $a \gg 1$,
 ξ is significant.

$$7e) i) p(x) = \sum_{k=0}^n a_k x^k = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1} + x^n$$

Let $p(x)$ be perturbed st. $p_{\text{per}}(x) = a_0 + \dots + (a_i + \delta a_i) x^i + \dots + a_{n-1} x^{n-1} + x^n$

$$\Rightarrow p_{\text{per}}(x + \delta x) = a_0 + \dots + (a_i + \delta a_i) (x + \delta x)^i + \dots + (a_{n-1}) (x + \delta x)^{n-1} + (x + \delta x)^n$$

$$\Rightarrow p_{\text{per}}(x + \delta x) = p(x + \delta x) + \delta a_i (x + \delta x)^i$$

$$\Rightarrow p_{\text{per}}(x - \delta x) - p(x + \delta x) = \delta a_i (x + \delta x)^i$$

$$\Rightarrow \left| \frac{p(x + \delta x) - p(x)}{\delta x} \right| = \left| \frac{a_i (x + \delta x)^i}{\delta x} \right|$$

$$\delta x \rightarrow 0 \quad \Rightarrow |p'(x)| = \left| \frac{a_i x^i}{\delta x} \right|$$

$$\text{From } K \left| \frac{\delta a_i}{a_i} \right| = \left| \frac{\delta x}{x} \right| \quad \Rightarrow K_i = \left| \frac{a_i x^{i-1}}{p'(x)} \right|$$

$$\text{Now } |Cond R_n(a_i)| = \left| \frac{a_i R_n^{i-1}}{p'(R_n)} \right|$$

$$\Rightarrow |Cond R_n(a_i)| = \sum_{i=0}^{n-1} \left| \frac{a_i R_n^{i-1}}{p'(R_n)} \right|$$

ii) See code.

iii) I doubt that it could, as there condition number are very large, so the problem isn't well-conditioned

$$8) a) y_{n+1} = e - (n+1)y_n$$

$$\Rightarrow \boxed{y_n = \frac{e - y_{n+1}}{n+1}}$$

Then for the mapping $g_k(y_n) = y_k$:

$$g_n(y_n) = y_n$$

$$g_{n-1}(y_n) = \frac{e - y_n}{n}$$

$$g_{n-2}(y_n) = \frac{e(n-1) + y_n}{n(n-1)}$$

One could formally prove by induction that

$$g_{n-a}(y_n) = \frac{e(n-1)^{a-1} + y_n}{(n-a-1)!}$$

[though the pattern is rather obvious]

$$\text{Furthermore, } g_k'(y_n) = \frac{(-1)^{n-k} k!}{n!}$$

$$\Rightarrow (c \text{ and } g_k)(y_n) = \left| \frac{y_n k!}{y_k n!} \right|$$

$$\text{As } k < n, \quad \max \left| \frac{y_n}{y_k} \right| = 1$$

$$\Rightarrow \boxed{(c \text{ and } g_k)(y_n) = \frac{k!}{n!}}$$

b) ~~Let ε_n~~

$$b) \varepsilon_N = 1, \quad \left| \frac{\varepsilon_N}{\varepsilon_N} \right| = \left| \frac{y_N g_N}{y_N} \right| \Rightarrow \varepsilon_N = \frac{y_N g_N}{y_N}$$

This is well-conditioned if $\varepsilon_N (= \varepsilon) \leq 1$.

$$\Rightarrow \frac{y_N g_N}{y_N} = 1 \quad \text{for minimum}$$

$$\Rightarrow \frac{K^1 y_N}{N^1 y_N} = 1 \quad \text{The nobig Mat } \frac{y_N}{y_N} \leq 1, \text{ for min. } \frac{y_N}{y_N} = 1$$

$$\Rightarrow \varepsilon \geq \frac{K^1}{N^1} \Rightarrow \left[N^1 \geq \frac{K^1}{\varepsilon} \right]$$

$$c) \varepsilon = \varepsilon_{ps}, \quad k = 20$$

$$N^1 = \frac{K^1}{\varepsilon_{ps}} = \frac{20^1}{2} \approx 2.19 \times 10^{34}$$

$$\text{By inspection: } 31^1 = 8.26 \times 10^{33} \\ 32^1 = 2.67 \times 10^{35}$$

$$\Rightarrow \left[N^1 \approx 32 \right]$$

$$d) x = 20, N = 32$$

$$\text{yields } y_{20} = 0.1238038 \quad [\text{see code}]$$

$$\text{Using Mathematica, } y_{20} = \int_0^1 e^{x^{20}} dx \approx 0.12380 \\ (\text{not accurate enough } B. \text{ to see relative error})$$

Using quad function, see code

$$\text{rel. error} = \underline{4.48 \times 10^{-16}}$$

Good