

① WLOG. consider  $x > 0$  s.t.  $x = \left( \sum_{n=1}^{\infty} b_n 2^{-n} \right) \cdot 2^e$

$$\min \|x\| = 2^{-1} 2^e$$

consider  $b_{p+1} = 0$ :

$$x - \text{rd}(x) = \sum_{n=1}^{\infty} b_n 2^{-n} \cdot 2^e - \sum_{n=1}^p b_n 2^{-n} \cdot 2^e$$

$$= \sum_{n=p+1}^{\infty} b_n 2^{-n} \cdot 2^e$$

$$\max \|x - \text{rd}(x)\| = \sum_{n=p+2}^{\infty} 2^{-n} \cdot 2^e = 2^{-p-2} \left( \frac{1}{1-\frac{1}{2}} \right) \cdot 2^e$$

$$= 2^{-1} \cdot 2^{-p} \cdot 2^e$$

$$\therefore \frac{\max \|x - \text{rd}(x)\|}{\min \|x\|} = \frac{2^{-1} \cdot 2^{-p} \cdot 2^e}{2^{-1} \cdot 2^e} = 2^{-p}$$

consider  $b_{p+1} = 1$ :

$$x - \text{rd}(x) = \sum_{n=1}^{\infty} b_n 2^{-n} \cdot 2^e - \left( \sum_{n=1}^p b_n 2^{-n} + 2^{-p} \right) \cdot 2^e$$

$$= \sum_{n=p+1}^{\infty} b_n 2^{-n} \cdot 2^e - 2^{-p} \cdot 2^e$$

$$\max \|x - \text{rd}(x)\| = 2^{-p-1} \cdot 2^e - 2^{-p} \cdot 2^e$$

$$\therefore \frac{\max \|x - \text{rd}(x)\|}{\min \|x\|} = \frac{(2^{-p-1} - 2^{-p}) 2^e}{2^{-1} \cdot 2^e}$$

$$= |2^{-p} - 2 \cdot 2^{-p}| = |1 - 2^{-p}| = 2^{-p}$$

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

for worst case of rounding down and rounding up.



②

a) See part a) of code.

★  
NOTE! %g kills trailing zeros  
in output, but everything  
is rounded to 5 sig. digs!

b) Converges to 244.71 for  $k = 18$  (see code).

Comparison and rel. error to `math.exp()`  
in code.

c) Yes, summing right to left introduced  
less error, since we add the smallest  
values first, so that by the time  
we are adding the larger values, the sum  
of all the previous small values is of  
a similar order as the subsequent large  
values, and we don't lose information  
due to limited precision. (as much)

(Again, see code)

d)  $\left. \begin{array}{l} \text{i)} \\ \text{ii)} \\ \text{iii)} \\ \text{iv)} \end{array} \right\} \text{(see code)}$

Quickest convergence: d) iii)

Lowest Error: d) ii)

Error is worse for  $x < 0$ . compared to  $x > 0$ .

e) Perform  $e^{5.5}$ , then take the reciprocal.

(see code) You get 0 error when rounding  
to 5 sig. digs.



(3)

a)

i) if  $x$  is a machine number, then:

$$fl(x) = x$$

$$\text{and } fl(x \cdot x) = x^2 (1 + \epsilon_{x^2}) \quad \text{where } \epsilon_{x^2} \leq eps$$

$$\begin{aligned} fl(fl(x^2) \cdot x) &= x^3 (1 + \epsilon_{x^2}) (1 + \epsilon_{x^3}) \\ &= x^3 (1 + \epsilon_{x^2} + \epsilon_{x^3}) \end{aligned}$$

$$fl(fl(x^{n-1}) \cdot x) = x^n (1 + \epsilon_{x^2} + \dots + \epsilon_{x^n})$$

where  $\epsilon_{x^2}, \dots, \epsilon_{x^n}$  are all bounded by  $eps$

$$\therefore |\epsilon_{x^2} + \dots + \epsilon_{x^n}| \leq (n-1)eps$$

if  $x$  is not a machine number, then:

$$fl(x) = x(1 + \epsilon_x) \quad \text{where } |\epsilon_x| \leq eps$$

$$\begin{aligned} \text{so } fl(fl(x) \cdot fl(x)) &= fl(x^2(1 + 2\epsilon_x)) = x^2(1 + 2\epsilon_x)(1 + \epsilon_{x^2}) \\ &= x^2(1 + 2\epsilon_x + \epsilon_{x^2}) \end{aligned}$$

$$\begin{aligned} fl(fl(fl(x) \cdot fl(x)) \cdot fl(x)) &= x^3(1 + 2\epsilon_x + \epsilon_{x^2}) \cdot (1 + \epsilon_x)(1 + \epsilon_{x^3}) \\ &= x^3(1 + 3\epsilon_x + \epsilon_{x^2} + \epsilon_{x^3}) \end{aligned}$$

so for  $x^n$ :

$$= x^n (1 + n\epsilon_x + \epsilon_{x^2} + \dots + \epsilon_{x^n})$$

where  $\epsilon_x, \epsilon_{x^2}, \dots, \epsilon_{x^n}$  all bounded by  $eps$

$$\therefore |n\epsilon_x + \epsilon_{x^2} + \dots + \epsilon_{x^n}| \leq (2n-1)eps$$



ii)  $f(\ln(x)) = \ln x (1 + \epsilon_{\ln})$  where  $|\epsilon_{\ln}| \leq \epsilon_5$

assuming  $n \ln x$  calculated through multiplication:

$f(n) = n$  (assume machine number according to almighty Gabe, via Piazza)

$$\therefore \underbrace{f(f(n)f(\ln x))}_y = n \ln x (1 + \epsilon_{\ln})(1 + \epsilon_{mult})$$

$$= n \ln x (1 + \epsilon_{\ln} + \epsilon_{mult})$$

$$f(\exp(y)) = e^{n \ln x (1 + \epsilon_{\ln} + \epsilon_{mult})} (1 + \epsilon_{exp})$$

$$= e^{n \ln x} e^{n \ln x (\epsilon_{\ln} + \epsilon_{mult})} (1 + \epsilon_{exp})$$

$$= x^n (1 + n \ln x (\epsilon_{\ln} + \epsilon_{mult})) (1 + \epsilon_{exp})$$

$$= x^n (1 + \epsilon_{exp} + n \ln x (\epsilon_{\ln} + \epsilon_{mult}))$$

where  $|\epsilon_{exp}|, |\epsilon_{\ln}|, |\epsilon_{mult}|$  all  $\leq \epsilon_5$

$$\boxed{\therefore |\epsilon_{exp} + n \ln x (\epsilon_{\ln} + \epsilon_{mult})| \leq \epsilon_5 (1 + 2n \ln x)}$$

assume  $n \ln x$  calc by adding  $\ln x$ ,  $n$  times:

$$f(f(\ln x) + f(\ln x)) = (\ln x (1 + \epsilon_{\ln}) + \ln x (1 + \epsilon_{\ln})) (1 + \epsilon_1)$$

$$= 2 \ln x (1 + \epsilon_{\ln} + \epsilon_1)$$

repeating sums:

$$n \ln x (1 + \epsilon_{\ln} + \epsilon_1 + \dots + \epsilon_{n-1}) = y$$

$$f(\exp(y)) = e^{n \ln x (1 + \epsilon_{\ln} + \epsilon_1 + \dots + \epsilon_{n-1})} (1 + \epsilon_{exp})$$

by same procedure:

$$= x^n (1 + \epsilon_{exp} + n \ln x (\epsilon_{\ln} + \epsilon_1 + \dots + \epsilon_{n-1}))$$

$$\boxed{\therefore |\epsilon_{exp} + n \ln x (\epsilon_{\ln} + \epsilon_1 + \dots + \epsilon_{n-1})| \leq \epsilon_5 (1 + n^2 \ln x)}$$



Since error for repeated multiplication is bounded based on  $n$ , and assuming  $\ln x \geq 3$  calculated with multiplication, then the log-exponential method is bounded by a function of  $\ln x$ . That is repeated multiplication would be better for all  $x$  except where  $x \approx 1$  and  $|\ln x| < 1$ . Otherwise the  $\ln x$  term just increases the error bound for the log-exp method.

b)

$$\begin{aligned} \text{i) } x^{a(1+\epsilon_a)} &= x^a x^{a\epsilon_a} \\ &= x^a e^{a\epsilon_a \ln x} \\ &= x^a (1 + a\epsilon_a \ln x) \end{aligned}$$

$$\therefore f(x^a) = x^a (1 + \epsilon_p)$$

$$\boxed{\epsilon_p = a\epsilon_a \ln x}$$

← Becomes substantial when  $x \approx 0$  since  $\ln x$  blows up. Also, if  $x$  is very large.

$$\begin{aligned} \text{ii) } (x(1+\epsilon_x))^a &= x^a (1+\epsilon_x)^a \\ &= x^a (1 + a\epsilon_x) \end{aligned}$$

$$\text{So if } f(x^a) = x^a (1 + \epsilon_p)$$

$$\text{then } \boxed{\epsilon_p = a\epsilon_x}$$

Propagated error ind. of  $x$ . But, gets large for large  $a$ , which is especially important for  $x \approx 1$  where  $x^a$  doesn't overflow.



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

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

$$\epsilon_f = \left( \frac{x f'(x)}{f(x)} \right) \epsilon_x$$

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

$$\lim_{x \rightarrow 0} (\text{cond } f)(x) = \lim_{x \rightarrow 0} \frac{\ln e^{-x} - x e^{-x}}{e^{-x}} = \lim_{x \rightarrow 0} 1 - x = 1$$

and  $(\text{cond } f)'(x) =$

$$\frac{(e^{-x} - x e^{-x})(1 - e^{-x}) - (x e^{-x})(e^{-x})}{(1 - e^{-x})^2}$$

$$= \frac{e^{-x} (1 - x - e^{-x} + x e^{-x} - x e^{-x})}{(1 - e^{-x})^2}$$

$$< 0 \text{ for } x \in [0, 1]$$

$$\therefore (\text{cond } f)(x) \text{ decreasing on } [0, 1]$$

$$\text{and } \leq 1 \text{ at } x=0 \text{ so}$$

$$\boxed{(\text{cond } f)(x) \leq 1 \text{ on } [0, 1]}$$

b)  $f_A(x) \Rightarrow f(-x) = -x$  (since  $x \neq \text{machine } \#$ )

$$f(e^{-x}) = e^{-x} (1 + \epsilon_{\text{exp}})$$

$$f(1 - f(e^{-x})) = (1 - e^{-x}(1 + \epsilon_{\text{exp}}))(1 + \epsilon_{\text{sub}})$$

$$= (1 + \epsilon_{\text{sub}}) - e^{-x} (1 + \epsilon_{\text{exp}} + \epsilon_{\text{sub}})$$

$$= (1 - e^{-x}) \left[ \frac{1 + \epsilon_{\text{sub}} - e^{-x} - e^{-x}(\epsilon_{\text{exp}} + \epsilon_{\text{sub}})}{1 - e^{-x}} \right]$$

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



and from part a),  $\epsilon_f = \frac{x f'(x)}{f(x)} \epsilon_a$  where  $\epsilon_a = \frac{|x_A - x|}{|x|}$

where  $\epsilon_f = \frac{\epsilon_{\text{sub}} - e^{-x}(\epsilon_{\text{exp}} + \epsilon_{\text{sub}})}{1 - e^{-x}}$

so  $\epsilon_a$  is not bounded by  $\epsilon_{\text{ps}}$

since  $f(x_A) = f_A(x)$   
 perfect algo  $\uparrow$  actual algo  $\uparrow$  actual input

$$\frac{e^x - 2}{x}$$

$\therefore \epsilon_a = \frac{\epsilon_{\text{sub}} - e^{-x}(\epsilon_{\text{exp}} + \epsilon_{\text{sub}})}{x e^{-x}}$

and  $\epsilon_{\text{sub}}$  and  $\epsilon_{\text{exp}}$  both bounded by  $\epsilon_{\text{ps}}$   
 - RHS expression max when  $\epsilon_{\text{sub}} = -\epsilon_{\text{exp}} = \epsilon_{\text{ps}}$

$\therefore \epsilon_a = \epsilon_{\text{ps}} \left( \frac{1}{x e^{-x}} \right)$

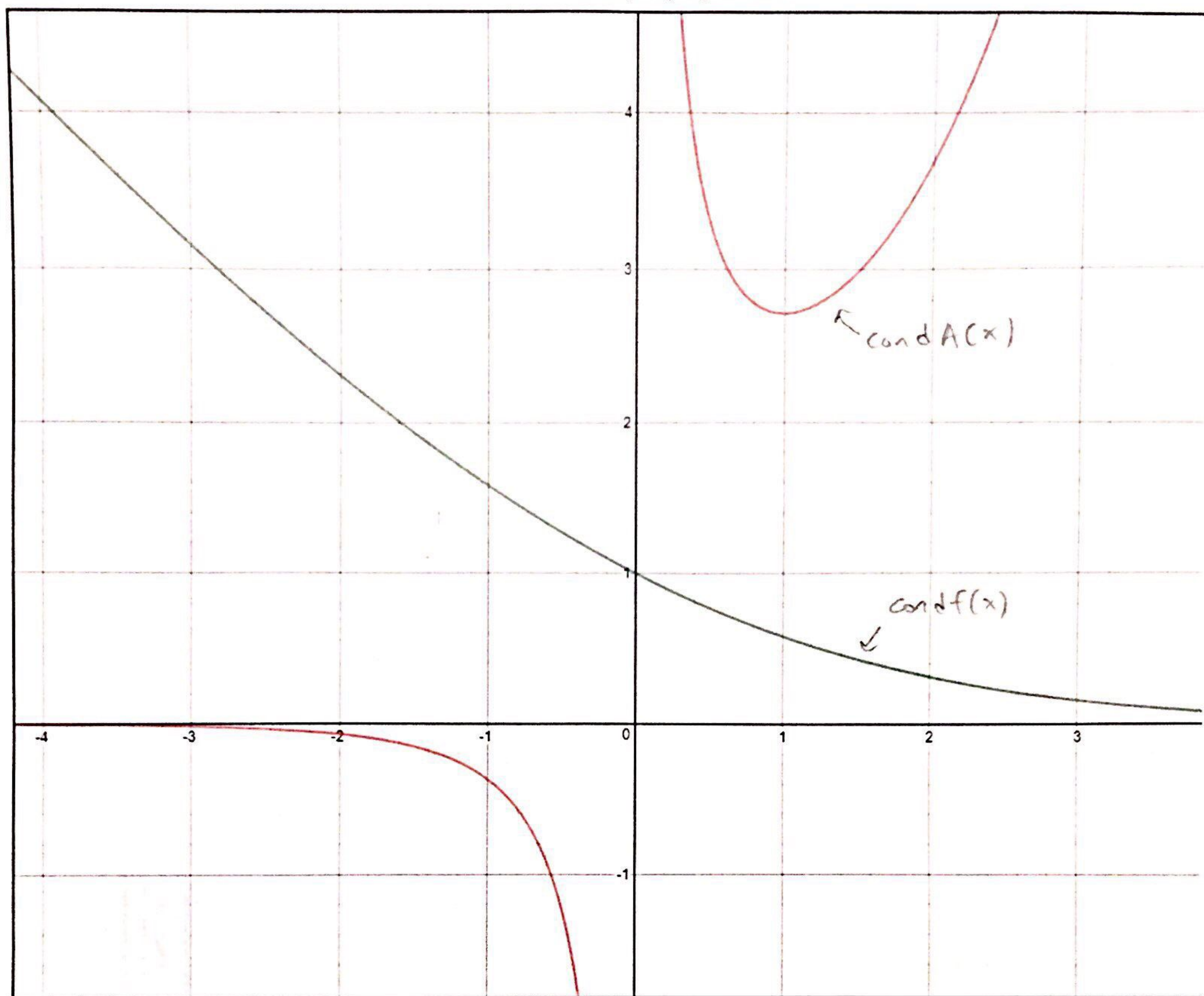
so  $\boxed{\text{cond}A(x) = \frac{e^x}{x}}$  condA

where  $\text{cond}A(1) = e \geq 1$

and  $\text{cond}A'(x) = e^x \left( \frac{x-1}{x^2} \right) < 0$  on  $[0, 1]$

$\boxed{\therefore \text{cond}A(x) \geq 1 \text{ on } [0, 1]}$





$$1 \quad \sim \quad \frac{1}{xe^{-x}} = \text{cond } A(x)$$

$$2 \quad \sim \quad \frac{(xe^{-x})}{1 - e^{-x}} = \text{cond } f(x)$$



c) See graph attached.

The cause of the poor conditioning is the subtraction in  $1 - e^{-x}$ . As  $x$  gets close to zero,  $e^{-x}$  approaches 1. As a result, two similar magnitude terms are subtracted, yielding catastrophic!!! error.

d)  $2^{-b} \leq 1 - \frac{q}{p} \leq 2^{-a}$

where  $p = 1$ ,  $q = e^{-x}$

$$\therefore 2^{-b} \leq 1 - e^{-x}$$

$$2^{-b} = 1 - e^{-x_{\min}}$$

$$e^{-x_{\min}} = 1 - 2^{-b}$$

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

For 1 bit loss!  $b = 1$

$$x_{\min} = 0.69315$$

2 bits!  $b = 2$

$$x_{\min} = 0.28768$$

3 bits!  $b = 3$

$$x_{\min} = 0.13353$$

4 bits!  $b = 4$

$$x_{\min} = 0.0645385$$



e)  $\frac{|y_A - y|}{|y|} \leq \text{cond} f(x) \left( \cancel{e^x}^6 + \text{cond} A(x) \cdot \text{eps} \right)$   
 $\uparrow$   
 assuming  $x$  is machine #

$$\therefore \frac{|y_A - y|}{|y|} \leq \left( \frac{x}{e^x - 1} \right) \left( \frac{e^x}{x} \right) \cdot \text{eps}$$

$$\frac{|y_A - y|}{|y|} \leq \frac{e^x}{e^x - 1} \text{eps}$$

plugging in  $x_{\min}$ 's from part d) and  
 using  $\text{eps} = 2^{-53}$  for double precision:

	$\frac{ y_A - y }{ y } \text{ max}$
1 bit	2 eps
2 bit	4 eps
3 bit	8 eps
4 bit	16 eps

f) Try to get rid of subtraction:

$$1 - e^{-x} \left( \frac{1 + e^{-x}}{1 + e^{-x}} \right) = \frac{1 - (e^{-x})^2}{1 + e^{-x}}$$

where  $e^{-x} = \cosh x - \sinh x$

$$= \frac{2 \sinh(x) e^{-x}}{1 + e^{-x}}$$

$$= \boxed{\frac{2 \sinh(x)}{e^x + 1}}$$

where  $\sinh(x)$  can be represented  
 by a Taylor series with all positive  
 terms (no subtraction).



or, alternatively, take Taylor series for  $e^{-x}$ :

$$1 - \left( 1 - x + \frac{x^2}{2} - \dots \right)$$

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

$$= x \left( 1 + x \left( \frac{1}{2} - (\dots) \right) \right)$$

where  $1 - e^{-x}$  can be approx. by  
N terms. Subtraction introduced in  
the above method is no longer between  
similarly sized quantities (for  $x \approx 0$ ).



⑤ (see code).

Converges to 1. This makes sense, since as  $n$  becomes larger and larger,  $\frac{1}{n}$  becomes smaller. Eventually,  $\frac{1}{n}$  becomes so small that  $1 + \frac{1}{n}$  simply evaluates to 1 within the machine precision (that is, the difference in scales of 1 and  $\frac{1}{n}$  exceeds the ability of the mantissa to capture them simultaneously). Of course, 1 to any exponent will then simply evaluate to 1, thus the loop converges to 1.

i.e. when  $n = 10^{16}$ ,  $e = 1$  b/c:

$$\text{then } e = \left(1 + 10^{-16}\right)^{10^{16}}$$

$$\text{but } 10^{-16} = 2^x$$

$$-16 \ln 10 = x \ln 2$$

$$x = \frac{-16 \ln 10}{\ln 2}$$

$$x = -53.15$$

$$\text{so } e = \left(1 + 2^{-53.15}\right)^{10^{16}}$$

$$\text{but } 1 + 2^{-53.15} \text{ is just } 1$$

with 53 bit mantissa,  $\therefore e = 1$



⑥ If we write  $x^{\frac{1}{2^i}}$  as  $e^{\frac{1}{2^i} \ln x}$  (see code)  
then take the Taylor series:

$$1 + \ln x \frac{1}{2^i} + \dots$$

we see that for  $i=52$ :

$$x^{\frac{1}{2^{52}}} = 1 + \ln x \frac{1}{2^{52}} + \dots$$

If  $x = e^n$  then

$$x^{\frac{1}{2^{52}}} = 1 + n 2^{-52} + \dots$$

but in double precision, we only have  
a 53 bit mantissa, so if  $n=1$ , then

$$= 1 + \underbrace{0.0\dots 01}_{52 \text{ dg.}} + \dots$$

$\therefore$  if  $n$  is slightly larger than 1, but less  
than 2, the information is lost when the  
first two terms are added, due to limited  
machine precision. Since the result of the  
sum is then squared 52 times, a small deviation  
from 1 is important. Yf is the same for  $e \leq x < e^2$   
since the same number is squared 52 times. We  
observe only integer powers of  $e$  as jumps.

similarly, for  $i=51$ :

$$x^{\frac{1}{2^{51}}} = 1 + \ln x 2^{-51} + \dots$$

Since  $\ln x$  is only shifted 51 places, the first binary  
decimal place for  $n = \ln x$  will not be lost due to  
precision limitations. As a result, half integer powers of  
 $e$  are also shown as jumps. Similar arguments can be  
made for other  $i$  values shown.



7

a) (see code)

b) Yes, it converges to 20 (about).

The alternative method also converges to 20.

c) The largest root becomes imaginary in each case, but the Newton method root-finder incorrectly locates a nearby real root.

It cannot find imaginary roots.

d) Same issue. Roots 16 and 17 are imaginary, but the Newton method cannot accommodate this.

$$e) i) p(\Omega_k + \delta\Omega_k) = \sum_{\substack{l=0 \\ l \neq i}}^{n-1} a_l (\Omega_k + \delta\Omega_k)^l + (a_i + \delta a_i) (\Omega_k + \delta\Omega_k)^i + (\Omega_k + \delta\Omega_k)^n$$

$$= \sum_{\substack{l=0 \\ l \neq i}}^{n-1} a_l \Omega_k^l \left(1 + \frac{\delta\Omega_k}{\Omega_k}\right)^l + (a_i + \delta a_i) \Omega_k^i \left(1 + \frac{\delta\Omega_k}{\Omega_k}\right)^i + \Omega_k^n \left(1 + \frac{\delta\Omega_k}{\Omega_k}\right)^n$$

$$= \sum_{\substack{l=0 \\ l \neq i}}^{n-1} a_l \Omega_k^l \left(1 + l \frac{\delta\Omega_k}{\Omega_k}\right) + (a_i + \delta a_i) \Omega_k^i \left(1 + i \frac{\delta\Omega_k}{\Omega_k}\right) + \Omega_k^n \left(1 + n \frac{\delta\Omega_k}{\Omega_k}\right)$$

$$= \underbrace{\sum_{\substack{l=0 \\ l \neq i}}^{n-1} a_l \Omega_k^l + a_i \Omega_k^i + \Omega_k^n}_{p(\Omega_k) = 0} + \underbrace{\sum_{\substack{l=0 \\ l \neq i}}^{n-1} a_l l \Omega_k^{l-1} \delta\Omega_k + a_i i \Omega_k^{i-1} \delta\Omega_k + n \Omega_k^{n-1} \delta\Omega_k}_{P'(\Omega_k) \delta\Omega_k}$$

$$+ \delta a_i \Omega_k^i + i \Omega_k^{i-1} \delta a_i \delta\Omega_k$$

and  $p(\Omega_k + \delta\Omega_k) = 0$  since  $\Omega_k + \delta\Omega_k$  is also a root,

$$\therefore 0 = P'(\Omega_k) \delta\Omega_k + \delta a_i \Omega_k^i \Rightarrow \frac{\delta\Omega_k}{\delta a_i} = \frac{-\Omega_k^i}{P'(\Omega_k)}$$

$$\text{and } T_{ki} = \frac{a_i \frac{\partial \Omega_k}{\partial a_i}}{\Omega_k}$$



$$\therefore \text{cond } \Omega_k(a) = \sum_{l=0}^{n-1} T_{kl}(a)$$

$$\text{cond } \Omega_k(a) = \sum_{l=0}^{n-1} \left| \frac{a_l \Omega_k^{l-1}}{P'(\Omega_k)} \right|$$

ii) (See code).

Since the condition numbers are so large for  $\Omega_k$ 's for 14, 16, 17, 20, this means that small perturbations in the coefficients  $a_l$  lead to large changes in the resulting roots.

iii) Since the problem itself is ill-conditioned, we can't use a clever algorithm to fix it. There will always be certain cases that behave poorly.



8

a)

$$g_k(y_N) = y_k$$

$$g_k(y_N(1 + \epsilon_N)) = g_k'(y_N + \Delta y_N)$$

$$= g_k(y_N) + \Delta y_N g_k'(y_N)$$

$$= y_k + y_N \epsilon_N g_k'(y_N)$$

$$= y_k \left( 1 + \frac{g_k'(y_N) y_N}{y_k} \epsilon_N \right)$$

$$= y_k (1 + \epsilon_k)$$

$$\therefore \text{cond } g_k = \frac{|\epsilon_k|}{|\epsilon_N|}$$

$$\text{cond } g_k = \left| \frac{g_k'(y_N) y_N}{y_k} \right|$$

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

$$y_{n-2} = \frac{e - \left( \frac{e - y_n}{n} \right)}{n-1}$$

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

reversed recurrence relation:

$$y_{n-1} = \left( \frac{e - y_n}{n} \right)$$

$$\therefore y_k = \frac{k!}{N!} \left( a - (b - (\dots (e - y_N))) \right)$$

so  $\left| \frac{dy_k}{dy_N} \right| = \frac{k!}{N!} \Rightarrow \left[ \text{cond } g_k = \left| \frac{k! y_N}{N! y_k} \right| \right]$



$$b) \quad \epsilon_k = \frac{g'_k y_N}{y_k} \epsilon_N$$

with  $\epsilon = \epsilon_k$  and assuming  $\epsilon_N = 1$ :

$$\epsilon_k = \frac{k!}{N!} \cdot \frac{y_N}{y_k} \epsilon_N$$

well-conditioned if:

$$\frac{k!}{N!} \cdot \frac{y_N}{y_k} \leq 1$$

and from the definition of  $y_N = \int_0^1 e^x x^N dx$

then  $y_N < y_k$  for  $N > k$

$$\therefore \left| \frac{y_N}{y_k} \right| \leq 1$$

$$\therefore \epsilon \geq \frac{k!}{N!}$$

$$\boxed{N! \geq \frac{k!}{\epsilon}}$$

c) For  $k=20$ ,  $\epsilon = 2^{-53}$

$$\frac{20!}{2^{-53}} = 2.19 \times 10^{34}$$

$$\text{so, for } \boxed{N=32} \cdot 32! = 2.63 \times 10^{35} > 2.19 \times 10^{34}$$

d) (see code and wolfram snippet)





integral from 0 to 1 of  $x^{20} \cdot e^x$



[Browse Examples](#) [Surprise Me](#)

Definite integral:

[Fewer digits](#)

[More digits](#)

☒ [Step-by-step solution](#)

$$\int_0^1 x^{20} e^x dx =$$

$$209 (4\,282\,366\,656\,425\,369 e - 11\,640\,679\,464\,960\,000) \approx 0.123803830762570$$

