

Charles Maher

i.)

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

Case 1:

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

$$x - \text{round}(x) = (0.0b_{p+2} b_{p+3} \dots) \times 2^e 2^{-p}$$

$$|x - \text{round}(x)| < (.0\overline{0})_2 \times 2^{e-p} = 2^{e-p-1}$$

Note smallest possible mantissa $\rightarrow (0.b_1 b_2 \dots)_2 \geq (0.b_1 00 \dots)_2 = (0.1)_2$

Case 2:

$$b_{p+1} = 1 \quad \text{round}(x) = (0.b_1 \dots b_p)_2 2^e + 2^{-p} \times 2^e$$
$$x - \text{round}(x) = (0.1b_{p+2} \dots)_2 2^{e-p} - 2^{e-p} = \frac{1}{2} + \frac{1}{2}(0.b_{p+2} \dots)_2 2^{e-p} 2^{e-p}$$
$$= \left(\frac{1}{2} + \frac{1}{2}(0.b_{p+2} \dots) - 1\right) 2^{e-p}$$

$$\text{and } 0 \leq (0.b_{p+2} \dots)_2 \leq 1$$

$$|x - \text{round}(x)| = \frac{1}{2}(1 - (0.b_{p+2} \dots)_2) 2^{e-p} < \frac{1}{2} 2^{e-p} \cdot 2^{e-p-1}$$

Using the above mantissa lower bound:

$$\frac{|x - \text{round}(x)|}{|x|} < \frac{2^{e-p-1}}{2^{e-1}} = 2^{-p} \quad \square.$$

Charles Maher

2.)

(code used to generate numbers attached separately)

a.)

(first 31 terms in attached notebook)

→ sum of these 31 terms is 244.71

b.)

(partial sums in notebook)

converges at $k=17$

$$\text{magnitude of error} \rightarrow \frac{|e^{5.5} - 244.71|}{e^{5.5}} = 7.3838 \times 10^{-5}$$

c.)

(partial sums in notebook)

also converges to 244.71, takes all 31 terms

relative error is the same because the converged value is the same

d.) (partial sums in notebook)

i.) converges to 3.8363×10^{-3} at $k=26$

$$\frac{|e^{-5.5} - 3.8363 \times 10^{-3}|}{e^{-5.5}} = 0.06128$$

ii.) converges to 4×10^{-3} at $k=31$

$$\frac{|e^{-5.5} - 4 \times 10^{-3}|}{e^{-5.5}} = 0.02123 \quad \leftarrow \text{least error}$$

iii.) converges to 0 at $k=18$, relative error is 1 \leftarrow fastest converging

iv.) converges to 1×10^{-2} at $k=31$

$$\frac{|e^{-5.5} - 1 \times 10^{-2}|}{e^{-5.5}} = 1.4469$$

- all are less accurate/slower converging than the $e^{5.5}$ case

e.) To compute e^{-SS} more accurately:

$$\frac{1}{e^{SS}} \rightarrow \frac{1}{244.71} = 0.0040865$$

$$\frac{|e^{-SS} - 0.0040865|}{e^{-SS}} = 6.6419 \times 10^{-5}$$

↑ orders of magnitude
less than error of
other methods explored

Charles Maher
3.)

a) let $x, n \in \mathbb{R}(p, q)$

i.) $x \rightarrow f(x \cdot x) = x^2(1+\varepsilon) \rightarrow f(x^3(1+\varepsilon)) = x^3(1+\varepsilon)^2 \sim x^3(1+2\varepsilon)$
 $\dots \rightarrow x^n(1+(n-1)\varepsilon)$

$$|\varepsilon_{xn}| \leq (n-1) \text{eps}$$

ii.) $f(\ln(x)) = \ln(x)(1+\varepsilon)$

$$f(n \ln(x)(1+\varepsilon)) = n \ln(x)(1+\varepsilon)(1+\varepsilon) \sim n \ln(x)(1+2\varepsilon)$$

$$\begin{aligned} f(e^{n \ln(x)(1+2\varepsilon)}) &= e^{n \ln(x)} e^{2n \ln(x)\varepsilon}(1+\varepsilon) \\ &\sim e^{n \ln(x)} (1 + 2n \ln(x)\varepsilon)(1+\varepsilon) \\ &= e^{n \ln(x)} (1 + (2n \ln(x) + 1)\varepsilon) + O(\varepsilon^2) \end{aligned}$$

$$|\varepsilon_{e^{n \ln(x)}}| \leq (2n \ln(x) + 1) \text{eps}$$

\rightarrow if $(n-1) \leq 2n \ln(x) + 1$

then repeated multiplication
is less error prone than
the exp-log method

D.)

i.) $x^{(a+\varepsilon_a)} \rightarrow x^a + \varepsilon_a x^a \log(x)$
 $= x^a (1 + \varepsilon_a \log(x))$

ii.) $(x + \varepsilon x)^a = x^a + a \varepsilon x x^{a-1} = x^a (1 + \frac{a \varepsilon x}{x})$

Scenario 1 is not likely to ever have
substantially high error, but Scenario 2
can high error if $x \rightarrow 0$

4.)

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

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

↳ problem condition

$$\frac{d}{dx} \left(\frac{x}{e^{x-1}} \right) = -\frac{e^x(x-1)+1}{(e^{x-1})^2} \rightarrow \text{derivative is negative semi-definite}$$

$$\rightarrow \max \left| \frac{x}{e^{x-1}} \right| \text{ on } [0,1] \text{ is at } 0$$

$$\lim_{x \rightarrow 0} \left| \frac{x}{e^{x-1}} \right| = \lim_{x \rightarrow 0} \left| \frac{1}{e^x} \right| = 1$$

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

$$b.) -x \in \mathbb{R}(p,q) \Rightarrow f_1(e^{-x}) = e^{-x}(1 + \varepsilon_{\text{exp}}) \rightarrow f_1(1 - e^{-x}(1 + \varepsilon_{\text{exp}}))$$

$$\hookrightarrow f_1(1 - e^{-x}(1 + \frac{e^{-x}}{1-e^{-x}} \varepsilon_{\text{exp}})) = 1 - e^{-x}(1 + \frac{e^{-x}}{1-e^{-x}} \varepsilon_{\text{exp}})(1 + \varepsilon_{\text{rd}})$$

$$\sim 1 - e^{-x}(1 + e^{-x} \frac{\varepsilon_{\text{exp}} + \varepsilon_{\text{rd}}}{1-e^{-x}})$$

$$= 1 - e^{-x} + (1 - e^{-x}) \varepsilon_{\text{rd}} + e^{-x} \varepsilon_{\text{exp}} = f_n(x)$$

$$f(x_n) = 1 - e^{-x_n}$$

→ setting equal and solving for x_n :

$$(1 - e^{-x})(1 + \varepsilon_{\text{rd}}) + e^{-x} \varepsilon_{\text{exp}} = 1 - e^{-x_n}$$

$$e^{-x} + \varepsilon_{\text{rd}}(1 - e^{-x}) + e^{-x} \varepsilon_{\text{exp}} = -e^{-x_n}$$

$$- \ln(e^{-x}(1 - \varepsilon_{\text{rd}} + \varepsilon_{\text{exp}} + \varepsilon_{\text{rd}} e^{-x})) = x_n$$

$$|x_n - x| = \ln(1 + (e^{-x} - 1) \varepsilon_{\text{rd}} + \varepsilon_{\text{exp}})$$

$$\leq |(1 + e^{-x}) \varepsilon_{\text{rd}} + \varepsilon_{\text{exp}}| \leq |2 + e^{-x}| \text{ cps}$$

$$\text{cond}(A) = \frac{1}{\text{cps}} \frac{|2 + e^{-x}|}{|x|} \text{ cps} = \frac{2 + e^{-x}}{|x|}$$

always
↳ true on $[0,1]$

$$\frac{2 + e^{-x}}{|x|} > 1 \rightarrow 2 + e^{-x} > x \rightarrow e^{-x} > x - 2 \Rightarrow \text{cond}(A) > 1 \text{ on } [0,1]$$

c.) (plot in separate pdf)

Poor conditioning result of subtraction

d.)

I) 1 bit $\rightarrow x = e^{-x}; y = 1$

$$2^{-1} < 1 - \frac{x}{4} \rightarrow \frac{1}{2} < 1 - e^{-x} \leftrightarrow e^{-x} < \frac{1}{2}$$

$$\begin{aligned} x &> -\ln(\frac{1}{2}) \\ x &> 0.6931 \end{aligned}$$

(at least, 0.6931)

II) 2 bits \rightarrow

$$\frac{1}{4} < 1 - e^{-x} \rightarrow e^{-x} < \frac{3}{4}$$

$$x > -\ln(\frac{3}{4})$$

$$x > 0.2877$$

(at least, 0.2877)

III) 3 bits $\rightarrow x > -\ln(\frac{7}{8})$

$$x > 0.1335$$

(at least, 0.1335)

IV) 4 bits $\rightarrow x > -\ln(\frac{15}{16})$

$$x > 0.0645$$

(at least, 0.0645)

e.) (rel. error in output) $\leq (\text{cond } f)(x)(\epsilon_{\text{PS}} + (\text{cond } A)(x)\epsilon_{\text{PS}})$

I) $\epsilon_y \leq \left| \frac{x}{e^x - 1} \right| (\epsilon_{\text{PS}} + \frac{2+e^x}{x} \epsilon_{\text{PS}})$
 $\leq .6931 \epsilon_{\text{PS}} + 4 \epsilon_{\text{PS}} = 4.6931 \epsilon_{\text{PS}}$

$\epsilon_y \leq 4.6931 \epsilon_{\text{PS}}$

II) $\epsilon_y \leq 10.863 \epsilon_{\text{PS}}$

III) $\epsilon_y \leq 22.9347 \epsilon_{\text{PS}}$

IV) $\epsilon_y \leq 46.9681 \epsilon_{\text{PS}}$

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

no subtraction \rightarrow likely better
conditioned
but more "expensive"

S.)

at 10^{16} we observe that the limit becomes one, (which is repeated for 10^{17} (i.e. n_{stop} is 10^{17})) which is because at 10^{16} $\frac{1}{10^{16}}$ becomes too small to store (underflow). Thus, $\frac{1}{10^{16}}$ is stored as 0, so the limit becomes $(1+0)^{10^{16}}$, which is just 1.

6.) Charles Maher
(plot is in notebook)

remaining values for S2 are

S3

S4

S0

e^0, e^1, e^2

e^0, e^2

$e^0, e^{3/4}, e^{5/4}, e^{1/4}, e^{3/2}$

$e^{1/4}, e^{3/4}, e^{7/4}, e^{9/4}$

$e^{1/4}, e^2, e^{9/4}$

from IEEE double
precision
53-n bits

→ powers are those representable with
of precision, where n is the number of
square roots/squares taken

This is because each \sqrt{x} removes a bit of
precision for each iteration, and so any #
above something representable with 53-n bits
will get floored because the less significant
bits wind up getting lost.

We get e to these powers of this because
when we square these functions back up to their
"original value" we can think of the following:
we now have numbers written like:

$(1 + \text{something tiny}^{5/4} \text{ with } 53\text{-n bits of mantissa})$

after the sqrts we have something like:
 $e^{1/4 + \frac{1}{4}\text{sqrting}} e^x$ but now x has very little
precision, i.e. can only take
the values listed above \square .

8

Charles Maher

7.1

a) (coeff in attached notebook)

b.) I) Newton largest root: 19.9999999571418 II) Numpy "roots": 20.00054209

$$\epsilon = 10^{-8}$$

I) 9.585389646516598 II) $20.6475355 \pm 1.186912412i$

$$\epsilon = 10^{-6}$$

I) 7.752713003402644 II) $23.1496169 \pm 2.7409845i$

$$\epsilon = 10^{-4}$$

I) 5.969334849605957 II) $28.400212411 \pm 6.51043422i$

$$\epsilon = 10^{-2}$$

I) 5.469592915093453 II) $38.47818362 + 20.93432359i$ (all roots also
in attached
notebook)(Newton based method
only detects real roots,
so largest real root is
taken)

d.) roots 16,17 have collapsed to the complex pair:

$$16.73096403 \pm 2.81265595i$$

$$e.) (\text{cond } \Omega_{k+1})(\bar{a}) = \sum_{l=0}^{n-1} (\Gamma_{kl})(\bar{a})$$

$$= \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{a_k \frac{\partial P_k(\bar{a})}{\partial a_k}}{P_k(\bar{a})} \right| = \left| \frac{a_k \sum_{i=0}^{k-1}}{P'(P_k)} \right|$$

ii.) (calculations in mathematica notebook attached)

$$14: 1.33297 \times 10^9$$

$$16: 2.40751 \times 10^9$$

$$17: 1.9044 \times 10^9$$

$$20: 4.30999 \times 10^7$$

extremely poorly conditioned, seems to
be worse conditioned for intermediate
roots (<20), so the problem itself is
prone to significant error

iii.) unlikely because the problem itself is poorly conditioned, so
the framework itself is flawed

Charles Maher

8.)

a.) $y_{n+1} = e^{-(n+1)}y_n \rightarrow y_n = \frac{e^{-y_{n+1}}}{n+1} \rightarrow \boxed{y_{n+1} = \frac{e^{-y_n}}{n}}$

$$y_{n-1} = \frac{e^{-y_n}}{n}; y_{n-2} = \frac{e^{-y_{n-1}}}{n-1} = \frac{e}{n-1} - \frac{e}{n(n-1)} - \frac{y_n}{n(n-1)},$$

$$y_{n-3} = \frac{e}{n-2} - \frac{e}{(n-2)(n-1)} + \frac{e}{n(n-1)(n-2)} + \frac{y_n}{n(n-1)(n-2)}$$

$$y_{N-k} = \dots = \frac{y_N}{N(N-1)\dots(N-k)} \rightarrow = \frac{y_N k!}{N!}$$

$$(\text{cond } g_k) = \left| \frac{y_N}{y_k} g'(y_N) \right| = \left| \frac{y_N}{y_k} \frac{k!}{2!} \right|$$

$\frac{y_N}{y_k}$ is at greatest ~ 2 if $N > k$

$$\rightarrow (\text{cond } g_k) < \left| \frac{k!}{N!} \right|$$

b.) $\varepsilon_y = (\text{cond } l) \varepsilon_x$

$$\varepsilon_y = 2; \varepsilon_x = 1$$

$$\varepsilon = (\text{cond } l) \rightarrow \varepsilon \geq \frac{k!}{N!} \quad \boxed{\frac{k!}{\varepsilon} \leq N!}$$

c.) let $k=20, \varepsilon = 2^{-53}$

$$20! 2^{-53} \approx 2.2 \times 10^{-34}$$

$$\underbrace{2.2 \times 10^{-34}}_{\text{closest two, but}} \leq N!$$

less than
821

$$\boxed{N \geq 32}$$

d.) (notebook w/ code attached)

↳ result = 0.1238038

Wolfram Alpha: 0.1238038
(to some amount of digits)

$$\text{error: } \frac{1.1238038 - .1238038}{.1238038} = 0$$