

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

note: $\text{rd}(x) = x(1 + \varepsilon)$ where $\varepsilon \leq \text{eps}$

$$x = \pm \left(\sum_{l=1}^{\infty} b_l 2^{-l} \right) 2^e$$

$$x - \text{rd}(x) = ? \Rightarrow 2 \text{ cases}$$

Case 1: First discarded bit is 0

$$x - \text{rd}(x) = \pm \left(\sum_{l=p+2}^{\infty} b_{-l} 2^{-l} \right) 2^e$$

$p+1$ digit is correct

$$|(x - \text{rd}(x))|_{\max} = \left(\sum_{l=1}^{\infty} 2^{-l} \right) 2^{e-p-1} \rightarrow b_{-l} = 1 \text{ if } l=p+2$$

$$\left(\lim_{l \rightarrow \infty} \sum_{l=1}^{\infty} 2^{-l} = 1 \right)$$

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

$$(x)_{\min} = 2^{e-1} \rightarrow \text{smallest mantissa} = 0.1 = 2^{-1}$$

$$\therefore \text{for Case 1: } \left| \frac{x - \text{rd}(x)}{x} \right| \leq \left| \frac{2^{e-p-1}}{2^{e-1}} \right| = 2^{-p} \quad \checkmark$$

Case 2: First discarded bit is 1

$$x - \text{rd}(x) = (2^{-p} - \sum_{l=p+1}^{\infty} b_{-l} 2^{-l}) 2^e \rightarrow b_{-l} = 1 \forall p+1 \leq l < \infty \wedge b_{-l} \leq 0 \forall l > p+1$$

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

$$\therefore \left| \frac{x - \text{rd}(x)}{x} \right| \leq \frac{2^{e-p-1}}{2^{e-1}} = 2^{-p}$$

x_{\min} doesn't change

Qii) \rightarrow X is a machine # \therefore error in xy comes from the $f1(xy)$ operation

$$f1(x*x) = (1 + \varepsilon_f)(xx)$$

{ floating error

$$\therefore f1((xx)x) = (xx)x(1 + \varepsilon_f)(1 + \varepsilon_f) = xx(1 + 2\varepsilon_f)$$

$$\therefore f1(x^n) = x^n(1 + (n-1)\varepsilon_f)$$

\Rightarrow max error is $x^n \underbrace{(n-1)\varepsilon_f}_{\varepsilon_{\text{tot}}}$

i) ① $f1(\ln x)$

② $f1(n \cdot f1(\ln(x)))$ or $f1(\ln(x)) + f1(\ln(x)) + \dots$

③ $f1(\exp(\text{step} \theta))$

① $f1(\ln x) = \ln(x)(1 + \varepsilon_e)$

② $f1(n \cdot f1(\ln(x))) = n \ln(x)(1 + \varepsilon_e)(1 + \varepsilon_n) = n \ln(x)(1 + \varepsilon_e + \varepsilon_n)$

or

$$f1(f1(\ln(x)) + f1(\ln(x))) = (\ln x + \ln x)(1 + \varepsilon_e + \varepsilon_n)$$

float of \ln
float of ε

$$\underbrace{z + f1(\ln(x))}_{z} = (\ln x + \ln x + \ln x)(1 + \frac{\ln x}{3\ln x}\varepsilon_e) + \frac{2\ln x}{3\ln x}(\varepsilon_e + \varepsilon_n)$$

$$w = 3\ln x \left(1 + \frac{3\varepsilon_e + 2\varepsilon_n}{3}\right) \rightarrow f1(w) = 3\ln x \left(1 + \frac{3\varepsilon_e + 5\varepsilon_n}{3}\right)$$

$$\frac{\ln x}{4\ln x} \varepsilon_2 + \frac{3\ln x}{4\ln x} \left(\frac{3\varepsilon_2 + 5\varepsilon_n}{3} \right) = \varepsilon_2 + \frac{5}{12} \varepsilon_n$$

$$\frac{4\varepsilon_2 + 5\varepsilon_n}{4} = \varepsilon_2 + \frac{5}{12} \varepsilon_n$$

$$f(1) \Rightarrow \varepsilon_2 + \frac{5}{12} \varepsilon_n$$

A1
Q3
P2

3iii) using adding technique:
control

$$f(1)(f(1(\ln x)) + f(1(\ln x)) + \dots) = n \ln x \left(1 + \varepsilon_2 + K\varepsilon_n \right)$$

where $K =$

③ using $f(1(nf(1(\ln x)))$ we get:

$$x^n = \exp(n \ln x (1 + \varepsilon_2 + \varepsilon_n)) = e^{n \ln x} * e^{n \ln x (\varepsilon_2 + \varepsilon_n)}$$

↓ taylor

$$\approx e^{n \ln x} (1 + n \ln x (\varepsilon_2 + \varepsilon_n))$$

$$\text{now, } f(1(x^n)) = e^{n \ln x} (1 + n \ln x (\varepsilon_2 + \varepsilon_n) + \varepsilon_f)$$

$$\Rightarrow \text{max error is } e^{n \ln x} \underbrace{(2n \ln x + 1)}_{\varepsilon_{\text{totii}}} \varepsilon_{\text{ps}}$$

note:

\rightarrow if x is small (< 1), $\ln x$ blows up $\therefore e^{n \ln x}$ is a bad choice as $\varepsilon_{\text{totii}}$ will blow up

\rightarrow if $x \sim 1$, $\varepsilon_{\text{totii}} \sim \varepsilon_{\text{ps}}$, $e^{n \ln x}$ is a good choice

b) $\rightarrow x \neq a$ not machine numbers

\rightarrow no floating point error

i) $a = a + \varepsilon_a$, $x = x$

$$x^{a+\varepsilon_a} = x^a x^{\varepsilon_a} = x^a (1 + \varepsilon_a \ln(x))$$

\downarrow taylor expand

$$\varepsilon_{\text{tot}} = \varepsilon_a \ln(x)$$

ii) $(x + \varepsilon_x)^a = x^a + a \varepsilon_x x^{a-1} = x^a \left(1 + \frac{a \varepsilon_x}{x}\right)$

\downarrow taylor around $\varepsilon_x \approx 0$

\Rightarrow for both cases error gets large when $x \rightarrow 0$

\Rightarrow Case 1: error will also get large as $x \rightarrow \infty$, but asymptotically will be very small when $x \sim 1$

\Rightarrow Case 2: error \sim large when $x < 1$, diminishes as $x \rightarrow \infty$

Also grows with a

Large $x \rightarrow$ use case 2

$x \sim 1 \rightarrow$ use case 1

$x < 1 \rightarrow$ bad all around

4) $f(x) = 1 - e^{-x} = y, [0, 1]$

a) $(\text{cond } f)(x) = ?$

$$\Delta y = f(x + \Delta x) - f(x) = f'(x) \Delta x$$

$$f'(x) = e^{-x}, \Delta y = e^{-x} \Delta x$$

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

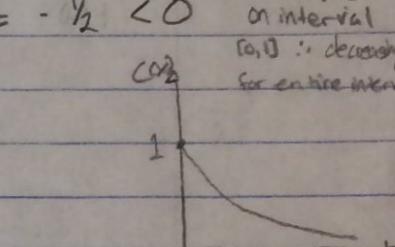
cond # = C

$$C(x) = \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{x e^{-x}}{1 - e^{-x}} \right| = \frac{x e^{-x}}{1 - e^{-x}} \quad (0, 1]$$

\hookrightarrow for $x \neq 0$ $\frac{dc}{dx} = \frac{e^x(x-1)+1}{(e^x-1)^2}, C'(0) = -\frac{1}{2} < 0$

decreasing
 $+C' \neq 0$
on interval
 $[0, 1] \therefore$ decreasing
for entire interval

$$C(1) = \frac{e^{-1}}{1 - e^{-1}} = \frac{1}{e-1} < 1$$



$$C(0+\varepsilon) = \left| \frac{\varepsilon e^\varepsilon}{1 - e^{-\varepsilon}} \right| = \left| \frac{\varepsilon}{e^\varepsilon - 1} \right| < 1$$

$\therefore C(x) < 0 \quad \forall x \in (0, 1]$

for $x=0$: $C(x) = \left| \frac{e^{-x}}{1 - e^{-x}} \right| = \left| \frac{1}{e^x - 1} \right| < 1 \quad \forall$

$$C(0) = \left| \frac{1}{-1} \right| = 1$$

4b)

$$A(i) : x \times (-1) \rightarrow \varepsilon = 0$$

$$A(ii) : \exp(A(i)) \rightarrow \exp(-x)(1 + \varepsilon_e)$$

$$A(iii) : 1 - A(ii)$$

$$\hookrightarrow (1 - \exp(-x)(1 + \varepsilon_e))(1 + \varepsilon_r) = A$$

$$A = (1 - \exp(-x)) \left(1 + \left(\frac{e^{-x}}{1-e^{-x}} \right) \varepsilon_e \right) (1 + \varepsilon_r)$$

$$A = (1 - e^{-x}) \left(1 + \frac{e^{-x}}{1-e^{-x}} (\varepsilon_e + \varepsilon_r) \right) \quad \begin{matrix} \varepsilon_e \neq \varepsilon_r \text{ banded} \\ \text{by } \text{eps} \end{matrix}$$

$$\Rightarrow \varepsilon_A = \text{eps} \left(\frac{e^{-x}}{1-e^{-x}} + 1 \right)$$

$$A = (1 - e^{-x}) (1 + \varepsilon_A) = f_A(x)$$

\Rightarrow now need $f(x_A)$

$$\Rightarrow \text{let } \varepsilon_{x_A} = \frac{x_A - x}{x}, \quad x_A = x(\varepsilon_{x_A} + 1)$$

$$f(x_A) = 1 - \exp(-x(\varepsilon_{x_A} + 1)) = 1 - \exp(-x) \exp(-x\varepsilon_{x_A})$$

$\stackrel{\text{3 taylor}}{\approx}$

$$f(x_A) = 1 - e^{-x} (1 - x\varepsilon_{x_A}) = 1 - e^{-x} + e^{-x} x \varepsilon_{x_A}$$

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

$$\frac{d-f}{f}$$

4b) now $f(x) = f_A(x)$
cond

$$(1 - e^{-x}) \left(1 + \frac{xe^{-x}}{1 - e^{-x}} \varepsilon_{x_A} \right) = (1 - e^{-x})(1 + \varepsilon_A)$$

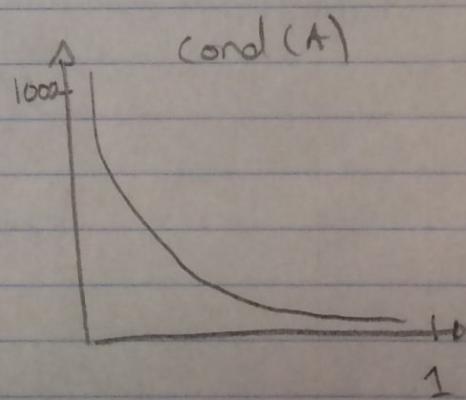
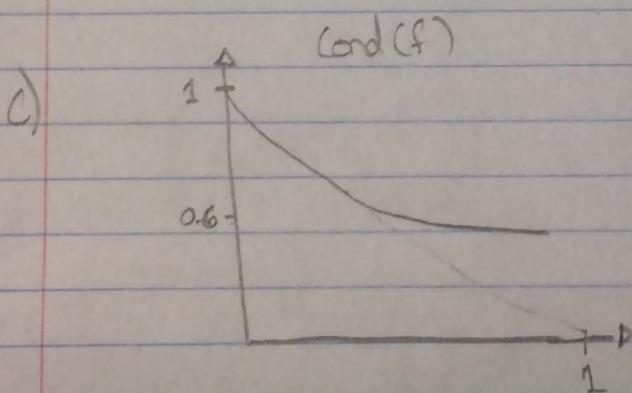
$$\varepsilon_A = \frac{xe^{-x}}{1 - e^{-x}} \varepsilon_{x_A}$$

$$\varepsilon_A = \frac{1 - e^{-x}}{xe^{-x}} \left(\text{eps} \left(\frac{e^{-x}}{1 - e^{-x}} + 1 \right) \right) = \text{eps} \left(\underbrace{\frac{1}{x} + \frac{1 - e^{-x}}{xe^{-x}}}_{\text{cond}(A)} \right)$$

$$\varepsilon_{x_A} = \frac{1}{x} \left(\frac{e^{-x}}{e^{-x}} + \frac{(1 - e^{-x})}{e^{-x}} \right) \text{eps} = \frac{e^x}{x} \text{eps}$$

$$\Rightarrow \text{cond}(A) = \frac{e^x}{x}$$

See code for nice plot



$$\lim_{x \rightarrow 0} \frac{e^x}{x} = \frac{e^0}{0} = \frac{1}{0} = \infty$$

$\rightarrow \gamma_x$ causes $\text{cond}(A)$ to explode on $[0, 1]$

4b) $f = \frac{1 - e^{-x}}{b_x - y}$

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

$$\Rightarrow b=1, 2^{-1} \leq 1 - e^{-x}$$

$$e^{-x} \leq \frac{1}{2}$$

$$-x \geq \ln(1/2)$$

$$x \geq -\ln(1/2) = 0.693$$

$$\Rightarrow b=2, e^{-x} \leq - (2^{-2}-1) = \frac{3}{4}$$

$$x \geq 0.287682 = -\ln(0.75)$$

$$\Rightarrow b=3, e^{-x} \leq - (2^{-3}-1) = \frac{7}{8}$$

$$x \geq -\ln(7/8) = 0.1335$$

$$\Rightarrow b=4, e^{-x} \leq - (2^{-4}-1) = \frac{15}{16}$$

$$x \geq -\ln(15/16) = 0.06454$$

c) $\frac{\|y^* - y\|}{\|y\|} \leq (\text{cond } f)(x) \left\{ \begin{array}{l} \varepsilon + (\text{ord } A)(x^*) \text{eps} \\ b \end{array} \right\}$

$\frac{\|x^* - x\|}{\|x\|} = 0$ since x is a machine #

$$4e) \frac{\|y_i - y\|}{\|y\|} \leq \left(\frac{xe^{-x}}{1-e^{-x}} \right) \left(\frac{e^x}{x} \right) \epsilon_{ps}$$

\downarrow
err

2^{-53}

$$\text{err} \leq \frac{1}{(1-e^{-x})} \epsilon_{ps}$$

ref err

for $b=1$, $e^{-x} = \gamma_2$, $\text{err}_{max} = 2 \epsilon_{ps}$

$b=2$, $e^{-x} = \frac{3}{4}$, $\text{err}_{max} = 4 \epsilon_{ps}$

$b=3$, $e^{-x} = \frac{7}{8}$, $\text{err}_{max} = 8 \epsilon_{ps}$

$b=4$, $e^{-x} = \frac{15}{16}$, $\text{err}_{max} = 16 \epsilon_{ps}$

$$4f) \quad 1 - e^{-x} \left(\frac{1+e^{-x}}{1+e^x} \right) = \frac{1 - (e^{-x})^2}{1+e^{-x}}$$

$e^{-x} = \cosh x - \sinh x$
 $(e^{-x})^2 = c^2 + s^2 - 2sc$
 $1 = \cosh^2 x - \sinh^2 x$
} plug in

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

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

where $\sinh(x) = \sum_{n=0}^{\infty} \frac{x^{2n+1}}{(2n+1)!}$ → no subtractions!

↳ MacLaurin series

a) for all these cases, first we do:

$$x^{\frac{1}{2^n}} \rightarrow \text{can write as } e^{\frac{1}{2^n} \ln(x)}$$

↳ taylor expand

$$x^{\frac{1}{2^n}} = 1 + \ln x \cdot \frac{1}{2^n}$$

now, this is the same as moving the decimal place of $\ln(x)$ n times \therefore
when $n=52$ we only have 1 sigfig left in the mantissa since for double precision it can hold 53 digits.
This means when $n=52$ the only numbers that match $y=x$ are those which are integer powers of e .
Similar logic applies for other values of n .

$\tilde{r}_k \tilde{a}$

7e) $r = r(a_0, a_1, \dots, a_{n-1})$, $r_k = k^{\text{th}}$ root of p_k

i) $\tilde{r}_{k\ell}$, cond # for k^{th} root changes ℓ^{th} coeff

for 1D, $\text{cond}(f) = f(x + \Delta x) - f(x)$

for 2D we get $p(r_k + \delta r_k, (a_\ell + \delta a_\ell)) = p_A$

$$p_A = a_0 + a_1(r_k + \delta r_k) + (a_\ell + \delta a_\ell)(r_k + \delta r_k)^2 + \dots$$

$$\begin{aligned} \text{note: } (r_k + \delta r_k)^\ell &= r_k^\ell \left(1 + \frac{\delta r_k}{r_k}\right)^\ell \\ &\quad \xrightarrow{\text{small, Taylor expand}} \\ &= r_k^\ell \left(1 + \ell \frac{\delta r_k}{r_k}\right) \\ &= r_k^\ell + \underbrace{\ell r_k^{\ell-1} \delta r_k}_{\frac{dp}{dr_k} \Big|_{r_k}} \end{aligned}$$

$$\begin{aligned} p_A &= a_0 + a_1 r_k + \dots + a_\ell r_k^\ell + a_\ell \delta r_k + \dots + a_\ell (\ell r_k^{\ell-1} \delta r_k) \\ &\quad + \delta a_\ell r_k^\ell + \delta a_\ell (\ell r_k^{\ell-1} \delta r_k) \end{aligned}$$

second order, ignore

$$p_A = \left(\frac{dp}{dr_k} \Big|_{r_k}\right) \delta r_k + \delta a_\ell r_k^\ell$$

$\frac{dp}{da_\ell} \Big|_{r_k}$

$$\Rightarrow \text{note } \text{cond}_{r_k}(a) = p[(r_k + \delta r_k), (a_\ell + \delta a_\ell)] - p(r_k, a_\ell)$$

$$\text{but } p(r_k, a_\ell) = 0$$

7e) \therefore in the limit that $\delta\sigma_k \neq \delta\alpha$ are small, $p_A \approx 0$

$$\Rightarrow \frac{\delta\sigma_k}{\delta\alpha} = - \left(\frac{\partial p}{\partial \alpha} \Big|_{\sigma_k} \right) \left(\frac{\partial p}{\partial \sigma_k} \Big|_{\sigma_k} \right)^{-1}$$

$$= \sigma_k^l (p'(\sigma_k))^{-1}$$

$$\therefore \text{cond } \sigma_k(a) = \sum_{l=0}^{n-1} \left| \frac{a \sigma_k^l}{\sigma_k p'(\sigma_k)} \right| = \sum_{l=0}^{n-1} \left| \frac{a \sigma_k^{l-1}}{p'(\sigma_k)} \right|$$

ii) $r = 14, 16, 17, 20$
 $\sigma_k = 15, 17, 18, 21$

$$\left\{ \begin{array}{l} C_{14} = 5.88 \times 10^{13} \\ C_{16} = 3.854 \times 10^3 \\ C_{17} = 1.723 \times 10^{13} \\ C_{20} = 1.3798 \times 10^{11} \end{array} \right.$$

\rightarrow condition number is lower for higher roots
 \rightarrow all condition #s are $\gg 1$
 \rightarrow small perturbations in coeffs lead to large changes in roots

iii) \rightarrow Problem itself is ill-conditioned \therefore
no clever algorithm will help \therefore

$$8) \quad y_n = \int_0^1 x^n e^x dx \quad n \geq 0$$

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

a) g_x is a map from y_N to y_K , $K < N$

(cond g_N)(y_N) wrt K+N = ?

$$y_n = \frac{e - y_{n+1}}{n+1} \Rightarrow y_k = \frac{e - y_N}{N} \quad \text{if } N = k+1 \quad (\text{special case})$$

$$(\text{cond } g_K)(y_N) = \begin{vmatrix} y_N & g'K \\ & y_K \end{vmatrix} \rightarrow \text{need } g'(K)$$

\hookrightarrow figured this out before

$$g^k(y) = \prod_{i=0}^{N-(k+1)} \frac{1}{(N-i)} = \frac{k!}{N!}$$

$$\Rightarrow (x_n g_k)(y_N) = \begin{vmatrix} y_N & k \\ y_N & N! \end{vmatrix}$$

8b) now assume $\varepsilon_N = 1$

$$\frac{\varepsilon_k}{\varepsilon_N} = \left| \frac{y_N g'_k}{y_k} \right|$$

$$\Rightarrow \varepsilon_k = \varepsilon = \frac{y_N g'_k}{y_k} \Rightarrow \text{well conditioned if } \varepsilon \leq 1$$

$$\therefore 1 \geq \frac{y_N g'_k}{y_k} = \frac{y_N}{y_k} \frac{k!}{N!}$$

need max of this

$$y_N = \int_0^1 e^x x^n dx \Rightarrow \text{for } N > k, y_N < y_k$$

$$\therefore \frac{y_N}{y_k} < 1, \text{ let it equal 1}$$

$$\varepsilon \geq \frac{k!}{N!} \Rightarrow \left(N! \geq \frac{k!}{\varepsilon} \right)$$

8c) $\varepsilon = \text{eps}$, $k = 20$, $N = ?$
 $\varepsilon = 2^{-53}$

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

$$30! = 2.6 \times 10^{32}, \quad 31! = 8.22 \times 10^{33}, \quad 32! = 2.63 \times 10^{35} > 2.19 \times 10^{34}$$

$\boxed{N=32}$

d) Using $k = 20$, $N = 32$, $y_N = 0$:

$$y_{20} = 0.1238 \dots \rightarrow \text{python}$$

$$\int_0^1 e^x x^{20} dx = 0.1238 \dots \rightarrow \text{wolfram}$$

$$\text{rel error} = \frac{|y_w - y_{20}|}{y_w} = 2.24 \times 10^{-16} \sim 2^{-53} = 1 \times 10^{-16}$$

wow \therefore