① <u>Rounding Error:</u>

For any bases we have to cases: the error for rounding up and the one for rounding down

eg. $1.23\overline{499}...$ vs $1.23\,50..$ to the 2nd decimal are
$\quad\quad\hookrightarrow 1.23 \to \epsilon = 0.005 \quad\quad \hookrightarrow 1.24 \to \epsilon = 0.005$
the worst case scenarios for each operation.

Meaning, for rounding up the error is largest if the last $(p+1)^{th}$ digit is 5 and the rest is zero and for rounding down it's worst when $(p+1)$ digit is 4 and rest is 9

In general, if $\beta$ is the base. the round up error is given by

$$\mathcal{E}_{up} = |x - x_{rd_{up}}| = \left|\frac{\beta}{2}\beta^{-(p+1)}\beta^q\right| = \left|\frac{1}{2}\beta^{-p}\beta^q\right| = \frac{1}{2}\beta^{q-p}$$

Round down:

$$\mathcal{E}_{down} = |x - x_{rd_{down}}| = \left|\left[\left(\frac{\beta}{2}-1\right)\beta^{-(p+1)}_{\underset{\text{digit}}{\uparrow}} + \sum_{\ell=p+2\,in}^{\infty}(\beta-1)\beta^{-\ell}\right]\beta^q\right|$$

$$\text{e.g } 9\text{ base}_{10}, 1\text{ base}_2$$

$$= \left|\left[\frac{1}{2}\beta^{-p} - \beta^{-(p+1)} + \underbrace{\left(\sum_{\ell=p+2}^{\infty}\beta^{-\ell}\right)(\beta-1)}_{\text{geometric series}}\right]\beta^q\right|$$

$$= \left|\left[\frac{1}{2}\beta^{-p} - \beta^{-(p+1)} + \frac{1}{\beta^{p+2}}\frac{\beta^{-1}}{\beta(\beta-1)}\right]\beta^q\right|$$

$$= \left|\left[\frac{1}{2}\beta^{-p} - \beta^{-(p+1)} + \beta^{-(p+1)}\right]\beta^q\right|$$

$$= \frac{1}{2}\beta^{q-p}$$

Since the error for both is the same we can conclude that
the max error for rounding $\mathcal{E}_{rd} = \frac{1}{2}\mathcal{E}_{tr} \iff \frac{1}{2}\beta^{q-p}$

3(a) (i) $fl(x^n) = fl(x \circ fl(x \circ \cdots = x^n(1+\varepsilon)^n$

$$\approx x^n(1+n\varepsilon) + O(\varepsilon^2)$$

Relative error: $\varepsilon_0 = \left| \dfrac{x_n^n - x^n(1+n\varepsilon)}{x_n} \right| \leq |n\varepsilon|$

3(c) (ii)

$$e^{n \ln(x)(1+\epsilon_\ell)(1+\epsilon_N)} (1+\epsilon_e)$$

$$e^{n \ln(x)(1+\epsilon_\ell+\epsilon_N)} (1+\epsilon_e)$$

$$e^{n \ln(x)} \; e^{n \ln(x)\epsilon_\ell} \; e^{n \ln(x)\epsilon_N} (1+\epsilon_e)$$

series expansion
{ omit
  higher
  Order
  terms

$$e^{n \ln(x)} (1 + n \ln(x)\epsilon_\ell)(1 + n \ln(x)\epsilon_N)(1+\epsilon_e)$$

$$e^{n \ln(x)} (1 + n \ln(x)(\epsilon_\ell+\epsilon_n))(1+\epsilon_e)$$

$$e^{n \ln(x)} \big(1 + \underbrace{n \ln(x)(\epsilon_\ell+\epsilon_n) + \epsilon_e}_{\epsilon}\big)$$

Relative Error:

$$\epsilon = \left| \frac{1}{1} - 1 - n\ln(x)(\epsilon_\ell+\epsilon_n) + \epsilon_e \right|$$

$$= \left| n\ln(x)(\epsilon_\ell+\epsilon_n) + \epsilon_e \right|$$

**3 (b)**

(i) 
$$fl(x^a) = x^{a(1+\epsilon_a)}(1+\epsilon_p)$$

$$= x^a x^{a\epsilon_a}(1+\epsilon_p)$$

$$= x^a x^{a\epsilon_a \ln x}(1+\epsilon_p) \qquad \downarrow \text{Taylor}$$

omit higher order eps

$$= x^a [1 + a\epsilon_a \ln(x)](1+\epsilon_p) \qquad \text{omit} - \text{''} -$$

$$= x^a [1 + \underbrace{a\epsilon_a \ln(x) + \epsilon_p}_{\epsilon}]$$

Relative Error is

$$\epsilon = \cancel{x^a} \; a\ln(x)\epsilon_a \; (+\epsilon_p)$$

not given/asked for in question but non negligible?

Problems could occur if either a is large at x is large. Also, x can ~~not to~~ be negative but → Complex and x = 0 is also an issue.

3(b) ii)   $fl(x^a) = [x(1+\epsilon_x)]^a (1+\epsilon_p)$      ? $\rightarrow$ Power Error

$[= [x + x\epsilon_x]^a]$

$= x^a (1+\epsilon_x)^a$

$= x^a e^{a \ln(1+\epsilon_x)}$      (Logarithmic Series $O(\epsilon_x^2)$)

$= x^a e^{a\epsilon_x}$      (Exponential Series $O(\epsilon_x^2)$)

$= x^a [1 + a\epsilon_x]$

$= x^a [1 + \underbrace{a\epsilon_x + \epsilon_p}_{\epsilon}]$      if power error is regarded.

Relative error:

$$\epsilon = a\epsilon_x + \epsilon_p$$

Here, the error could only become substantial if $a$ is large. Otherwise this is fine, I guess.

4 (a)  Find Condition of $f(x) = 1 - e^{-x}$

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

Expand to series.

$$= \left| \frac{x}{e^{x} - 1} \right| = \left| \frac{x}{-1 + 1 + \frac{x}{1} + \cdots + } \right| \qquad \left| \quad \right| : \text{# cancel } x$$

$$\left(\text{cond } f\right)(x) = \left| \frac{1}{1 + \frac{x}{1} + \frac{x^2}{2!} + \cdots} \right| \geqslant 1$$

I.e., on $x \in \{0, 1\}$  $\underline{\underline{(\text{cond } f)(x) \leq 1}}$

4 (b)      $f.(x_A) = f_A(x)$ , where $f(x) = 1-e^{-x_A}$

$$f_A(x) = \left[ 1 - e^{-x}(1+\epsilon_1) \right](1+\epsilon_1)$$

(C)         Find $\epsilon$ error in $f_A$

↳ See Code

$$f_A(x) = \left[ 1 - e^{-x}(1-\epsilon_1) \right](1-\epsilon_2)$$

$$= 1 - e^{-x} - \epsilon_1 e^{-x} + \epsilon_2 - \epsilon_2 e^{-x}$$

$$= (1-e^{-x})\left\{ 1 + \underbrace{\frac{\epsilon_2 - e^{-x}(\epsilon_1+\epsilon_2)}{1-e^{x}}}_{\epsilon_A} \right\}$$

Now, find condition of A

$$f(x_A) = f_A(x)$$

$$1 - e^{-x_A} = \left[1 - e^{-x}\right](1+\epsilon_A)$$

$$\cancel{x} - e^{-x_A} = \cancel{x} - e^{-x} + \epsilon_A - \epsilon_A e^{-x} \quad | \cdot e^{x}$$

$$-e^{x-x_A} = -1 + \epsilon_A e^{+x} - \epsilon \quad | ^{\wedge}(-1)$$

$$e^{x_A - x} = (1 + \epsilon_A - \epsilon_A e^{x})^{-1} \quad | \ln$$

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

$$x_A - x = -\ln(1 + \epsilon_A - e^{x}) \qquad \text{Log series} \cdot \ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} \cdots$$

$$= -\epsilon(1+e^{x}) + O(\epsilon_A^2)$$

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

Then, we need to find the max
                          smallest sensible number

unbounded at

$$\max_{x\in[0,1]}\left|\frac{x_A - x}{x}\right| = \left\{ \max(x_A - x)\Big|_{x \leqslant 1} = \frac{\epsilon_A(e^{x}-1)}{x} \to \begin{array}{c}\text{unbounded at}\\ x=0\end{array}\right.$$

Hence, (cond A)(x) is
always larger than 1
on [0,1]

$$\Rightarrow (\text{cond } A)(x) \qquad \frac{\epsilon_A\left|\frac{x_A - x}{x}\right|}{\epsilon} = \frac{e^{x}-1}{x} > 1$$

$$x \in [0,1]$$

4 (d)     We have $(1)(e^{-y})$
           ↓   ↓
           x   y

$$2^{-b} \le 1 - \frac{y}{x} \le 2^{-a}$$

$$2^{-b} \le 1 - e^{-x} \le 2^{-a}$$

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

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

$$+x \ge -\ln\left(1 - 2^{-b}\right)$$

$$x \ge \ln\left(\left(1 - 2^{-b}\right)^{-1}\right)$$

Willing to lose one bit $\underline{b = 1}$

$$x \ge \ln\left(\left(1 - \frac{1}{2}\right)^{-1}\right) = \underline{\underline{\ln 2}}$$

2 bits  $b = 2$

$$x \ge \ln\left(\left(1 - \frac{1}{4}\right)^{-1}\right) = \underline{\underline{\ln \frac{4}{3}}}$$

~~$x \ge \ln t$~~

~~Although~~ 3 bits  $b = 3$

$$x \ge \ln\left(\left(1 - \frac{1}{8}\right)^{-1}\right) = \underline{\underline{\ln \frac{8}{7}}}$$

4 bits  $b = 4$

$$x \ge \ln\left(\left(1 - \frac{1}{16}\right)^{-1}\right) = \underline{\underline{\ln \frac{16}{15}}}$$

4 (e)   The upper bound for the relative Error is given by

the $(\text{cond } A)(x)\big|_{x=x_0}$ , where $x_0$ is the value/s from

from (d)

1 bits:
$$(\text{cond } A)(\ln 2) = \frac{e-1}{\ln 2} \approx 2.47896.$$

2 bits:
$$(\text{cond } A)\left(\ln \tfrac{4}{3}\right) = \frac{e-1}{\ln \tfrac{4}{3}} \approx 5.97285$$

3 bits:
$$(\text{cond } A)\left(\ln \tfrac{8}{7}\right) = \frac{e-1}{\ln \tfrac{8}{7}} \approx 12.868$$

4 bits:
$$(\text{cond } A)\left(\ln \tfrac{16}{15}\right) = \frac{e-1}{\ln \tfrac{16}{15}} \approx 26.62413$$

4(f)  $f(x) = 1 - e^{-x}$  can be rewritten into

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

This again can be rewritten into Another function, ie, the using Taylor expansion

$$f(x) = \frac{x + \frac{x^2}{2!} + \cdots}{e^x}$$

this way we omit the subtraction and hence the source of largest error.

(6.) Explanation          Mahtissa

$$1.\boxed{0\ 0\ \cdots\ -\qquad\cdots\ \cdots\ 0\ d}$$

Double Precision has 52-bits in mantissa

When we take the square root of a binary number say

$$(4)_{10} = (100)_2 = 2^2 \Rightarrow \sqrt{2^2} = 2^{\frac{2}{2}} \quad \text{Same again 52 times}$$

└→ bit shifts to right in mantissa

Meaning, we will be left with with $1.\ 0\ \cdots\ 0\ d$, i.e.
1 significant bit

Meaning, when taking the square root and the squaring
52 times we are left with:

$$\left(1 + 2^{-52}\right)^{2^{52}}$$

looking back at exercise 5, there is striking similarity and we
find that $\qquad \underset{}{} \left(1 + 2^{-52}\right)^{2^{52}} \approx \lim\limits_{n \to \infty} \left(1 + \frac{1}{n}\right)^{n}$

Inspecting more closely we find that values larger that $2^2$ say
$2^3$ result in $e^2$ e.g.

The reason is simply that their significant digit is ~~one~~ an exponent larger

$$\text{giving} \quad \left(1 + 2^{-51}\right)^{2^{52}} = \left(1 + 2^{-51}\right)^{2^{51} \cdot 2} \approx \lim\limits_{n \to \infty} \left(1 + \frac{1}{n}\right)^{2} = \underline{e^2}$$

~~the~~

For comparison, square-rooting and squaring 51 times leaves one
more significant digit that is $1.\boxed{0\ 0\ \cdots \boxed{d\ d}}$
and the result is $\left(1 + 2^{-50} + 2^{-52}\right)$. Hence, the
plot will output more intermediate values between $e\ e^1\ e^2 \cdots$

7(e)

$$\text{The } (\text{cond } \Omega_k)(\underline{a}) = \sum_{\ell}^{n-1} (\Gamma_{k\ell})\left(\frac{a}{\underline{a}}\right)$$

$$\Gamma_{k\ell} = \left|\frac{a_\ell \frac{\partial \Omega_k}{\partial a_\ell}}{\Omega_k}\right| = \left|\frac{a_\ell \frac{\partial \Omega_k}{\partial \rho(\Omega_k)} \frac{\partial \rho(\Omega_k)}{a_\ell}}{\Omega_k}\right|$$

$$= \left|\frac{a_\ell \frac{1}{\rho'(\Omega_k)} (\Omega_k)^\ell}{\Omega_k}\right|$$

$$= \left|a_\ell \frac{(\Omega_k)^{\ell-1}}{\Omega_k \rho'(\Omega_k)}\right|$$

So, the condition is given by

$$(\text{cond } \Omega_k)(\underline{a}) = \sum_{\ell}^{n-1} \left|a_\ell \frac{(\Omega_k)^{\ell-1}}{\rho'(\Omega_k)}\right|$$

8(a) $\quad y_n = \dfrac{(e - y_{n+1})}{n+1} \qquad \dfrac{\partial y_n}{\partial y_{n+1}} = - \dfrac{1}{n+1}$

$$(\text{cond } g_k)(y_N) = \left| \dfrac{y_N \dfrac{\partial g_k}{\partial y_N}}{g_k} \right|$$

$$= \left| \dfrac{y_N}{y_k} \dfrac{\partial y_k}{\partial y_{k+1}} \dfrac{\partial y_{k+1}}{\partial y_{k+2}} \cdots \dfrac{\partial y_{N-1}}{\partial y_N} \right|$$

$$= \left| \dfrac{\partial y_N}{g_k} \dfrac{(-1)}{k+1} \dfrac{(-1)}{k+2} \cdots \dfrac{(-1)}{N} \right|$$

$$= \left| \dfrac{y_N}{y_k} (-1)^{N-k+1} \dfrac{k!}{N!} \right|$$

$$(\text{cond } g_k)(y_N) = \left| \dfrac{y_N}{y_k} \dfrac{k!}{N!} \right|$$

(b) Now, ~~assuming~~ ~~though~~ ~~$\frac{y_N}{y_k} < g_k$, so the~~
we have that $y_N$ is always smaller than $y_k$, to evaluate
an upper bound we "need" to say ~~and~~ that $\frac{y_N}{y_k}$ is
largest when $\frac{y_N}{y_k} = 1$. Hence, we get the expression

$$\epsilon_x \leq \left| \dfrac{k!}{N!} \right|$$

(c) See in the code!
$$N! \leq \dfrac{k!}{\epsilon}$$

Stirling's approx for guessing $\quad N \ln(N) - N \approx\!\Rightarrow N! = 32$

$\underline{N! \leq 32!}$ $\qquad$ To get a good approximation of $y_k \big|_{k=20}$.

(d) Even though we start at the wrong initial value there is
a recurrence relation which makes up for it an converges to the
value of $0.1238\ldots$ See Code!