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

**a)** $f'(x) = e^{-x}$. Thus,

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

on $[0,1]$, $e^x \geq 1$ and $e^x \leq e$

As $x \to 0$ $\quad \dfrac{x}{e^x - 1} \to \dfrac{1}{1} \leq 1$ and as $x \to 1$,

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

It seems that, as $(\text{cond} f)(x)$ has no local minima or maxima, it goes monotonically from $1$ to $\frac{1}{e-1}$ on this interval.

**b)** The condition of $A$ is:

$$(\text{cond } A)(\bar{x}) = \frac{|\bar{x}_A - \bar{x}|}{|\bar{x}|}$$

find first what $f_A(x)$ yields. Then find

$$x_A = f^{-1}(f_A(x))$$

$f_A(x)$: first, $x \longrightarrow -x(1 + \varepsilon_x)$

then, $e^{-x(1+\varepsilon_x)}(1 + \varepsilon) = e^{-x}(1 - x\varepsilon_x + \varepsilon)$

and, $1.0 - e^{-x}(1 - x\varepsilon_x + \varepsilon) = f_A(x)$

Assume the $1.0$ has no error. So any error only comes from the $e^{-x}$ let's try to move the error to the exponent for ease of finding $x_A$.

$$e^{-x}(1 + \varepsilon - x\varepsilon_x) \approx e^{-x} e^{\varepsilon - x\varepsilon_x} \approx e^{-x + \varepsilon - x\varepsilon_x}$$

$$\Rightarrow \quad x_A = x + \varepsilon + x\varepsilon_x$$

$$\Rightarrow (\text{cond } A)(x) = \frac{|x + x\varepsilon_x + \varepsilon - x|}{|x|} \varepsilon^{-1}$$

$$= \left| \varepsilon_x + \varepsilon/x \right| \varepsilon^{-1}$$

Assume that $\varepsilon = \varepsilon_x$. Then

$$(\text{cond } A)(x) = \left| 1 + 1/x \right| > 1 \quad \text{on } 0 \text{ to } 1.$$

In fact as $x \to 0$, $(\text{cond } A)(x)$ diverges.

(c) See plot attached.
The root of the problem is that the error from exponentiation, $\varepsilon$, gets comparable to the magnitude of $x$ as $x$ gets small.

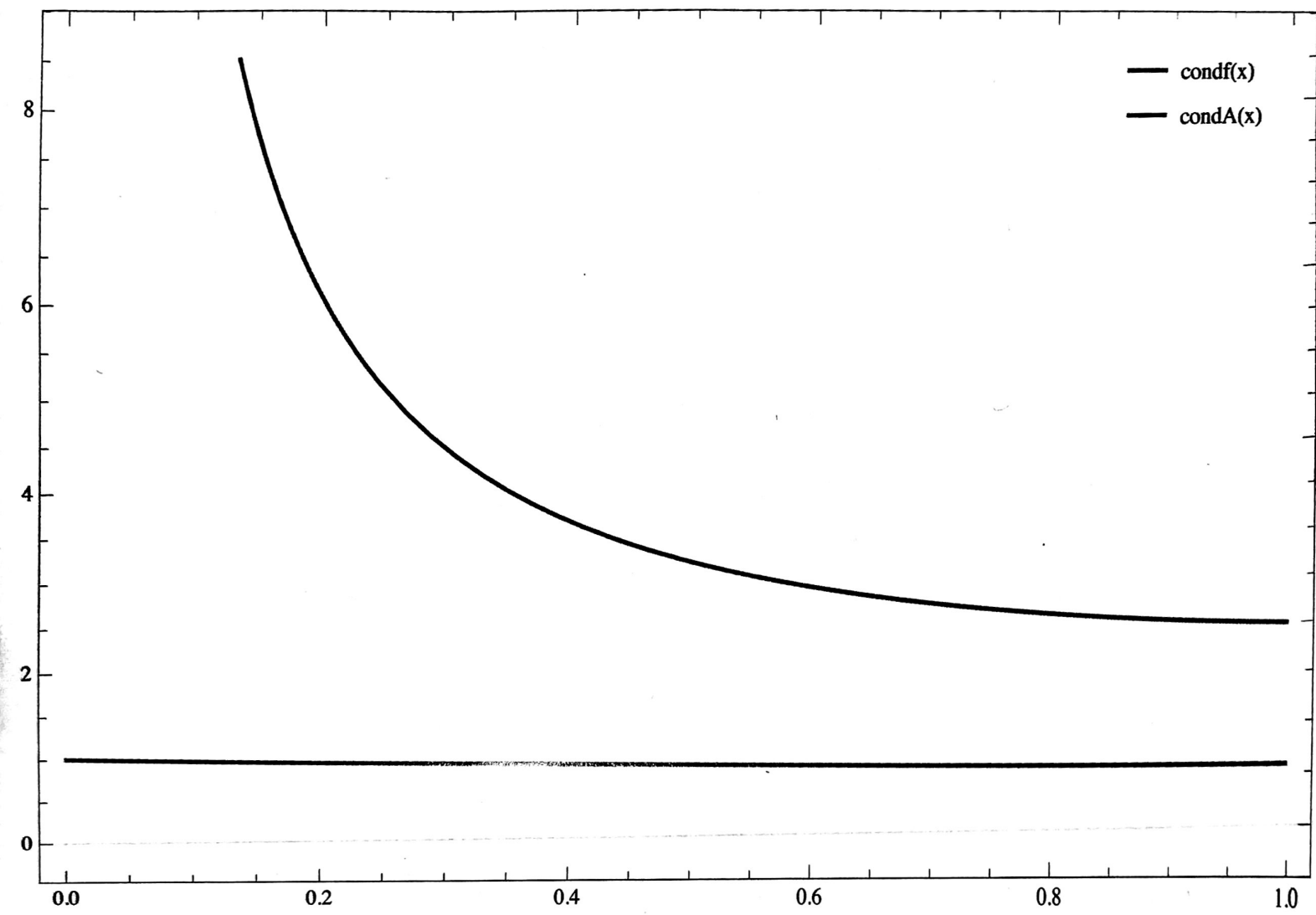d) S'pose we have a $p$ bit mantissa.
Any number on $[0,1]$ can be written as

$$x = \sum_{n=1}^{p} b_n 2^{-n} \quad (\text{rounded})$$

The last bit of significance is $b_p 2^{-p}$.
This is true even when $x$ is small. Let's at least write $x = 2^e \left( \sum_{n=1}^{p} b_n 2^{-n} \right)$

It's not clear to me if it's meant $n$ bits in $x_A$ or in $f$. $f$ isn't to sensitive to errors, so I'll assume it's in $x$.

If $\varepsilon = 2^{-p-1}$, then $2\varepsilon$ is losing one digit of significance, and $2^n \varepsilon$ is losing $n$ bits of significance.

when does $1 + \frac{1}{x} = 2^n$?

$$(2^n - 1)^{-1} = x_{min}$$

We require $x > x_{min}$ to lose at most $n$ bits of significance. When $n=1$,

$x > x_{min} = 1$. We always lose one bit of precision.

What if $n = 2, 3, 4$?          $\frac{1}{3}$

| $n$ | $x_{min}$ |
|---|---|
| 2 | $\frac{1}{3}$ |
| 3 | $\frac{1}{7}$ |
| 4 | $\frac{1}{15}$ |

e) We knew the relative error is bounded by

$$\text{cond} f(x) \left( \varepsilon + \varepsilon (\text{cond } A)x \right) = \text{err}(x)$$

So :

| $x$ | $\text{err}(x)/\varepsilon$ |
|---|---|
| 1 | 1.746 |
| $\frac{1}{3}$ | 1.966 |
| $\frac{1}{7}$ | 1.993 |
| $\frac{1}{15}$ | 1.999 |

f) When $x$ is small, we could use an alternative method. Namely,

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

When $x$ is small, this is well behaved. we see that

$$f(x) \to \frac{1 + 2x}{1 - x} \quad \text{for small } x.$$

This is fine.