

Part 1

Problem 2

(b) from left to right, the sum is 244.71. If $k \geq 17$, then the sum no longer changes. The relative error is 7.3839×10^{-5}

(c) from right to left, the sum is 244.70. If $k \geq 17$, then the sum no longer changes. The relative error is 3.2971×10^{-5}

(d) from (i) to (iv), the values are 0.0037353, 0.004, 0, 0.01, $k = 24, 19, 19, 19$ respectively. The relative errors are 0.086002, 0.021232, 1, 1.4469

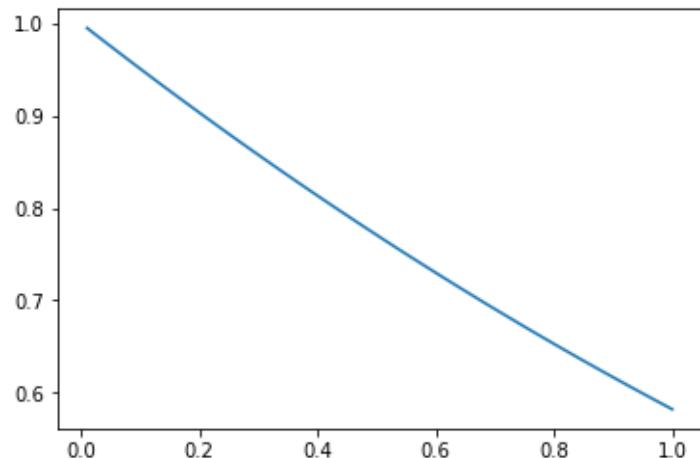
Thus, we can see adding from right to left is always the most accurate, because add small numbers first leads to higher accuracy. And we should avoid add positive and negative part respectively. Adding from right to left also converges more quickly.

(e) compute $e^{-5.5}$ by first compute $e^{5.5}$ the get its inverse $1/e^{5.5}$. The value is 0.0040866, the relative errors is 4.19496×10^{-5} .

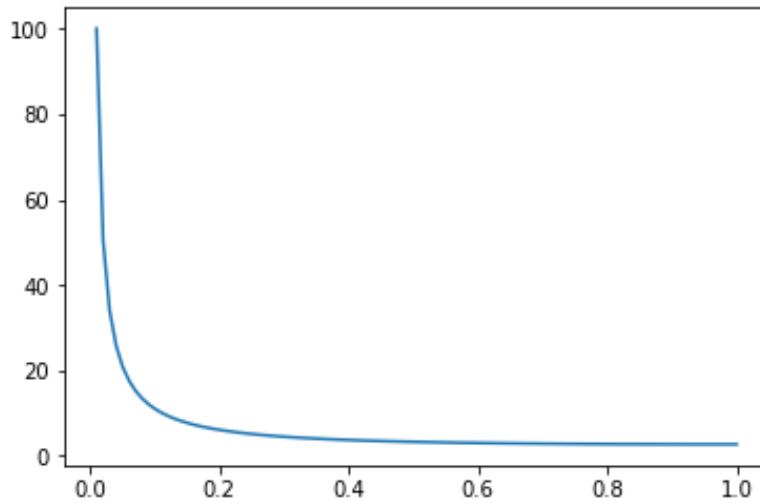
Problem 4

(c)

$$(\text{cond}f)(x) = x/(e^x - 1)$$



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



Problem 5

`n_stop = 10*13`, and the values are

```

2
2.5937424601
2.704813829422
2.716923932236
2.718145926825
2.718268237192
2.718280469096
2.718281694132
2.718281798347
2.718282052012
2.718282053235
2.718282053357
2.718523496037
2.716110034087 (n_stop)
2.716110034087

```

The final value it converges is 1. This is because when n is greater than 2^{53} , that is about 10^{16} , $1/n$ becomes 0 in the machine. Thus $(1+1/n)$ is 1 and $(1+1/n)^n = 1$.

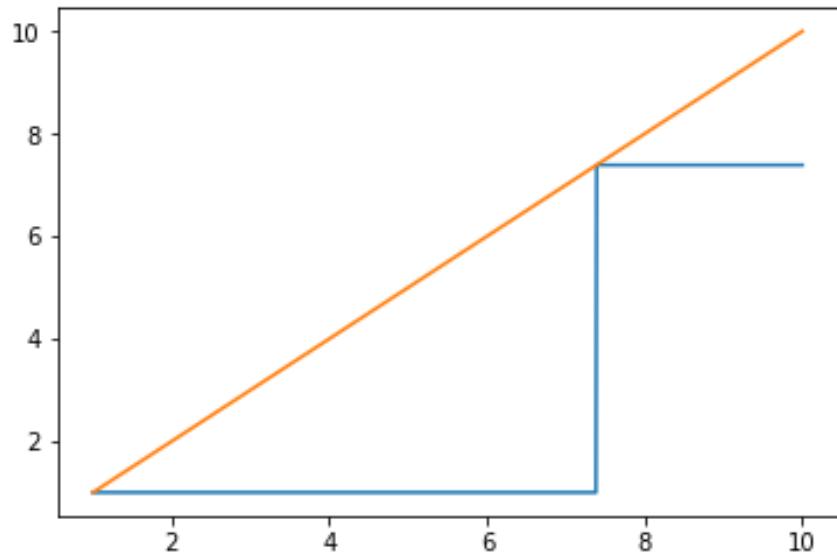
The reason it seemingly converges to 2.716110034087 is

1. The relative error of $1/n$ is about $10^{13}/10^{16} = 1/1000$ (that is e_x), after adding and exponentiating, the relative error is about $(\frac{n}{n+1} e_x + ne_1)$, which is in the order of 10^{-3} , thus only the first 2 or 3 digit is accurate

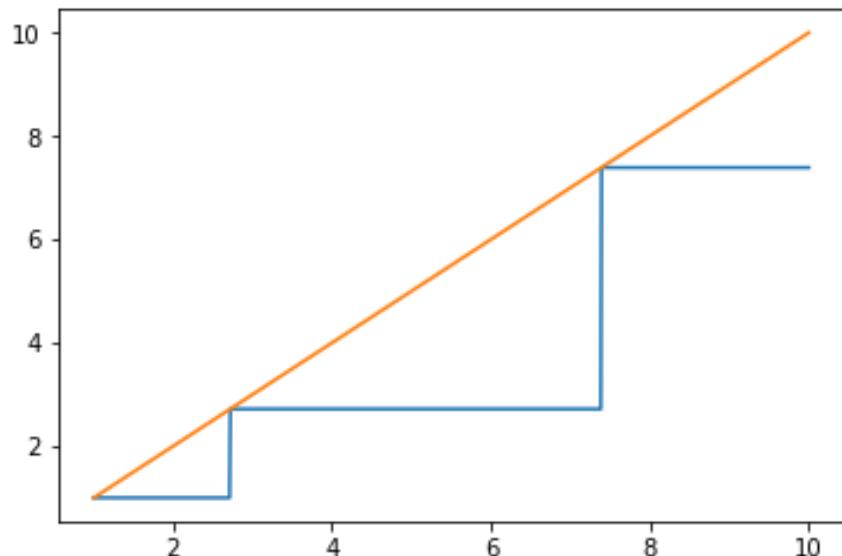
2. Coincidentally, $1+1/10^{13}$ and $(1+1/10^{14})^{10}$ has almost the same binary form, while this is not true for $(1+1/10^{15})^{100}$.

Problem 6

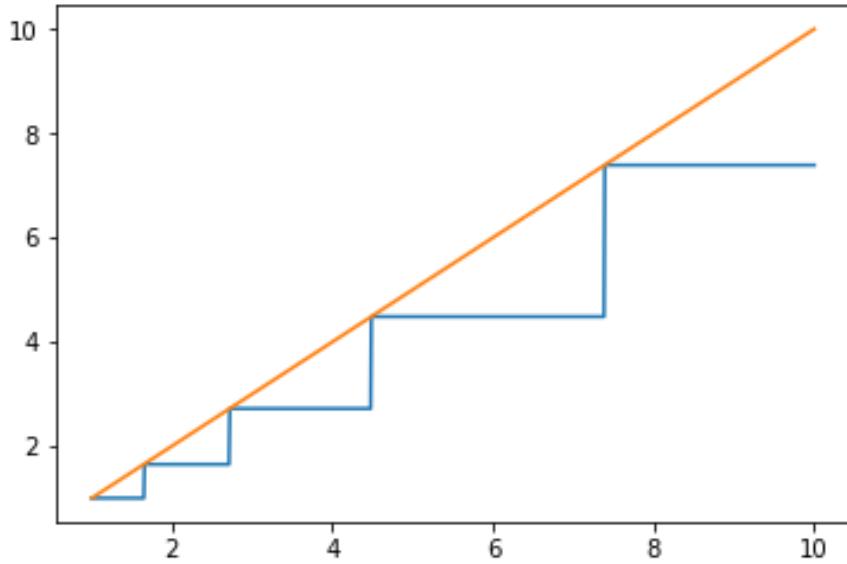
$N = 53$, the figure is as follow:



$N = 52$, the figure is as follow:



$N = 51$, the figure is as follow:



For N=53, all numbers are represented as either $(1/2+1/2^{53})^2$ or $(1/2)^2$ after square root in computer, thus exponentiate them will yield either $[(1/2+1/2^{53})^2]^{(2^{53})} = e^2$ or 1.

For N=52, all numbers represented as $(1/2)^2$, $(1/2+1/2^{53})^2$ or $(1/2+1/2^{52})^2$, thus exponentiate them will yield 1, e or e^2 .

The same thing holds for N=51 and other N

Problem 7

(a) 2432902008176640000,

-8752948036761600000,

13803759753640704000,

-12870931245150988800,

8037811822645051776,

-3599979517947607200,

1206647803780373360,

-311333643161390640,

63030812099294896,

-10142299865511450,

1307535010540395,

-135585182899530,

11310276995381,

-756111184500,

40171771630,

-1672280820,

53327946,

-1256850,

20615,

-210.0,

1

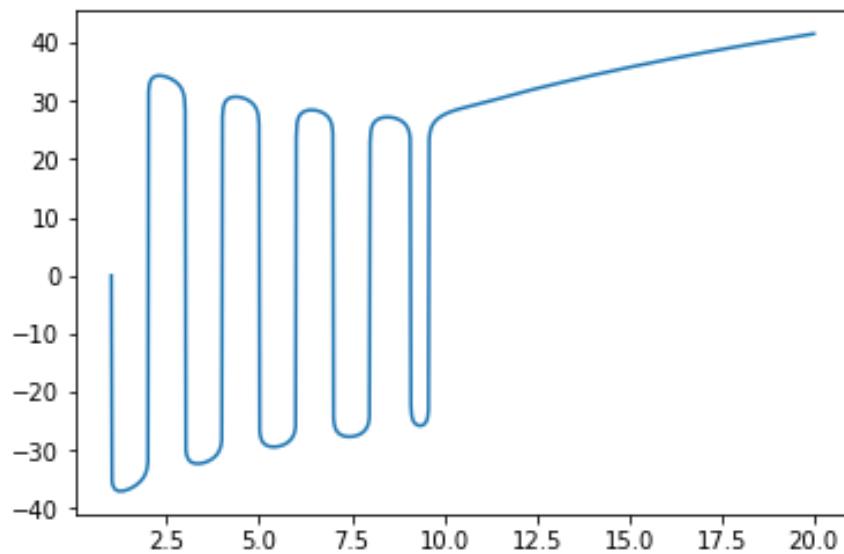
(b) use `scipy.optimize.newton`, the root is 19.9999949571418

(c) $\delta = 10^{-8}$, the root is 9.585389646516598

10^{-6} , the root is 7.752713003402644

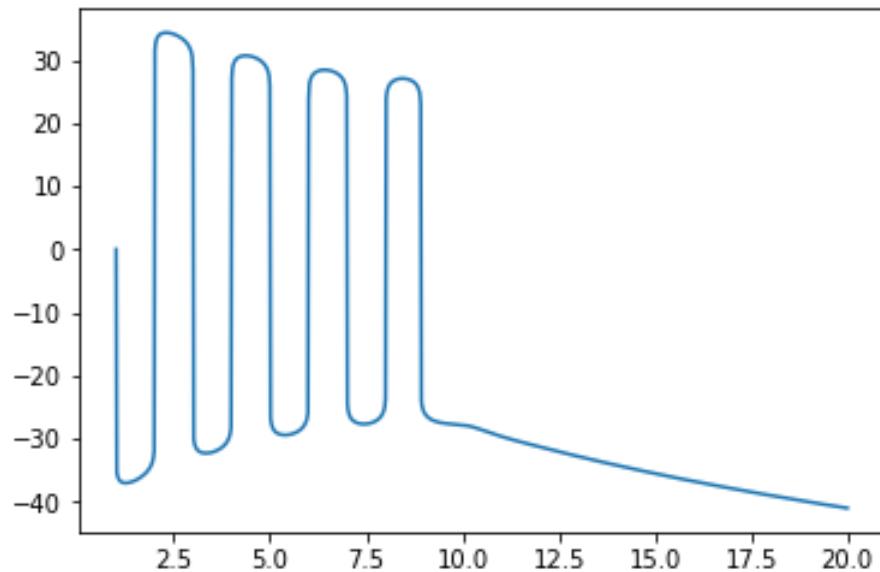
10^{-4} , the root is 5.969334849605957

10^{-2} , the root is 5.469592915093453



plot of $\text{sgn}(w(x)) * \log(1 + |w(x)|)$ for $\delta = 10^{-8}$

(d) for root 16, the value becomes 8.007305890325446
for root 17, the value becomes 8.917114218525562



plot of $\text{sgn}(w(x)) * \log(1 + |w(x)|)$

(e) see other part

Part 2

$$1. \text{ let } x = 2^q \times 0.\underbrace{1 \dots 1}_{p\text{-digits}} \dots = 2^q \times \left(\underbrace{0.1 \dots 1}_{p\text{-digits}} + \underbrace{0.0 \dots 0}_{(all \ 0)} \dots \right) \quad (1)$$

$$\text{then } rd(x) = 2^q \times \left(\underbrace{0.1 \dots 1}_{p\text{-digits}} + 2^{-p} \right)$$

$$\text{thus } |x - rd(x)| = 2^q [2^{-p} - \underbrace{0.0 \dots 0}_{(p+1)\text{ digits}}] \leq 2^{-(p+1)} \times 2^q$$

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

$$\text{let } x = 2^q \times 0.\underbrace{1 \dots 0}_{p\text{-digits}} \dots$$

$$\text{then } rd(x) = 2^q \times \left(\underbrace{0.1 \dots 0}_{p\text{-digits}} \right)$$

$$\text{thus } |x - rd(x)| \leq 0.\underbrace{0000 \dots 0}_{(p+1)\text{ digits}} \leq 2^{-(p+1)} \times 2^q$$

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

$$3. (a) f(\lambda x^k) = \lambda f(x)$$

$$f((x^2)f(\lambda x)) = f((x \circ x) \circ \lambda) = x^2(1 + \varepsilon_1)$$

$$f((x^3)f(\lambda x^3)) = f((x^2 \circ x) \circ \lambda) = x^3(1 + \varepsilon_1) \cdot x(1 + \varepsilon_2) \\ = x^3(1 + \varepsilon_1 + \varepsilon_2)$$

$$x + xe^{-x}$$

$$\frac{1+x e^{-x}}{1-e^{-x}} f(x^n) = x^n (1 + \varepsilon_1 + \dots + \varepsilon_{n-1}) \quad |\varepsilon_i| \leq \varepsilon \quad (1 \leq i \leq n-1)$$

$$f(e^{n \ln x}) = f((e^{n f(\ln x)}) \circ \lambda) = f((e^{n \ln x(1 + \varepsilon_1)}) \circ \lambda) \\ = f((e^{n \ln x(1 + \varepsilon_1 + \varepsilon_2)}) \circ \lambda)$$

$$= e^{n \ln x \cdot (1 + \varepsilon_1 + \varepsilon_2)} (1 + \varepsilon_3) = e^{n \ln x} e^{n \ln x \cdot (\varepsilon_1 + \varepsilon_2)} (1 + \varepsilon_3) \\ = x^{n \ln x \cdot (1 + \varepsilon_1 + \varepsilon_2)} (1 + \varepsilon_3) \approx x^n [\ln x (\varepsilon_1 + \varepsilon_2) + 1] (1 + \varepsilon_3)$$

$$f_l(e^{n \ln x}) \approx x^n [n \ln x (\varepsilon_1 + \varepsilon_2) + 1] (1 + \varepsilon_3) \quad (2)$$

$$\approx x^n [1 + \varepsilon_3 + (n \ln x)(\varepsilon_1 + \varepsilon_2)]$$

$$\leq x^n [1 + (2n \ln x + 1)\varepsilon]$$

$$f_l(x^n) = x^n (1 + \varepsilon_1 + \varepsilon_2) \leq x^n [1 + \cancel{\varepsilon} (n-1)\varepsilon]$$

so if $2n \ln x + 1 \geq n^{-1}$, that means $n \ln x \geq \frac{1}{2}$, that is $\frac{2}{n} \geq 1 - 2 \ln x$

then multiplication is more accurate than log-exponential

$$(b) \text{ (ii)} f_l(x^a) = f_l(e^{a \ln x}) = f_l[e^{a(1+\varepsilon_a) \ln x}] \quad (\varepsilon_a \text{ is the error of } a)$$

$$= f_l[e^{a(1+\varepsilon_a) \ln x \cdot (1+\varepsilon_1)}] \quad (\varepsilon_1 \text{ is the error of } \ln x)$$

$$= f_l[e^{a(1+\varepsilon_a) \ln x \cdot (1+\varepsilon_1 + \varepsilon_2) \ln x}]$$

$$= f_l[e^{(1+\varepsilon_a + \varepsilon_1)(1+\varepsilon_2) a \ln x}] \quad (\varepsilon_2 \text{ is the error of mult.})$$

$$x^a x = f_l(x) dx \quad = f_l[e^{(1+\varepsilon_a + \varepsilon_1 + \varepsilon_2) a \ln x}]$$

$$= e^{(1+\varepsilon_a + \varepsilon_1 + \varepsilon_2) a \ln x} (1 + \varepsilon_3) \quad (\varepsilon_3 \text{ is the error of exp.})$$

$$x(1 + \varepsilon_x) \quad \text{if } a \text{ is very large}$$

$$x + x\varepsilon_x \approx x^n [1 + \varepsilon_3 + (a \ln x)(\varepsilon_a + \varepsilon_1 + \varepsilon_2)] \quad \text{or } x \text{ is very small}$$

(if $A(x)x^a \neq f_l(e^{a \ln x}) / f_l(e^{a \ln(x+\varepsilon_x)})$) ε_x is the error of x
 then the error will be substantial

$$f_l[e^{a \ln(x+\varepsilon_x)(1+\varepsilon_1)}] \quad \varepsilon_1 \text{ is the error of } \ln x$$

$$= f_l[e^{a \ln(x+\varepsilon_x)(1+\varepsilon_1 + \varepsilon_2)}] \quad \varepsilon_2 \text{ is the error of } a \ln x /$$

$$= f_l[e^{a \ln x + \ln \frac{x+\varepsilon_x}{x}}]$$

$$\begin{aligned}
 (\text{iii}) f_l(X^a) &= f_l(e^{alnX}) = f_l\left\{e^{aln[X(1+\varepsilon_x)]}\right\} = f_l\left\{e^{alnX(1+\frac{\ln(1+\varepsilon_x)}{lnX})}\right\} \\
 &= f_l\left\{e^{alnX(1+\varepsilon_1+\varepsilon_2)\left[1+\frac{\ln(1+\varepsilon_x)}{lnX}\right]}\right\} \quad (3) \\
 &= f_l\left\{e^{alnX(1+\varepsilon_1+\varepsilon_2)(1+\frac{\varepsilon_x}{lnX})}\right\} \\
 &= e^{alnX(1+\varepsilon_1+\varepsilon_2+\frac{\varepsilon_x}{lnX})(1+\varepsilon_3)} \\
 &= e^{alnX}\left[1+(\varepsilon_1+\varepsilon_2)\cdot alnX + a\varepsilon_x\right](1+\varepsilon_3) \\
 &= e^{alnX}\left[1+(\varepsilon_1+\varepsilon_2)alnX + a\varepsilon_x + \varepsilon_3\right]
 \end{aligned}$$

thus, the ^{propagated} error of a is $X^a \cdot (alnX) \varepsilon_a$

the propagated error of x is $X^a \cdot a\varepsilon_x$

when a is very large or x is very small or x is very large
 ε_a could be very substantial

when a is very large, ε_x could be substantial

$$4. (a) (\text{cond } f)(x) = \frac{xf'(x)}{f(x)} = \frac{xe^{-x}}{1-e^{-x}} = \frac{x}{e^x-1} \leq 1$$

(4)

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

$$f_A(x) = f(1 - e^{-x(1+\varepsilon_x)}) = f(1 - e^{-x}(1 + \varepsilon_x))$$

$$= f(1 - e^{-x}(1 - x\varepsilon_x + \varepsilon_1)) \quad (\varepsilon_1 \text{ is the error of exp})$$

$$= [1 - e^{-x}(1 - x\varepsilon_x + \varepsilon_1)](1 + \varepsilon_2) \quad (\varepsilon_2 \text{ is the error of multiplication})$$

$$= (1 - e^{-x}) \left[1 + \frac{e^{-x}}{1 - e^{-x}} (-x\varepsilon_x + \varepsilon_1) + \varepsilon_2 \right]$$

D

$$f(x_A)(x - x_A) \neq 0$$

$$(1 - e^{-x}) [f_A(x) - f(x)] = f(x_A) - f(x), \text{ thus}$$

$$(1 - e^{-x}) \left[\frac{e^{-x}}{1 - e^{-x}} (-x\varepsilon_x + \varepsilon_1) + \varepsilon_2 \right] = f'(x)(x - x_A), \text{ thus}$$

$$x - x_A = \left[\frac{1}{1 - e^{-x}} (-x\varepsilon_x + \varepsilon_1) + e^x \varepsilon_2 \right] (1 - e^{-x})$$

$$\therefore (\text{cond } A)x = \left| \frac{x - x_A}{x} \right| / \varepsilon = \left[\frac{1}{1 - e^{-x}} (-x\varepsilon_x + \varepsilon_1) + e^x \varepsilon_2 \right] / (x\varepsilon) \leq$$

$$(1 - e^{-x}) \left[\frac{(1+x)\varepsilon}{1 - e^{-x}} + e^x \varepsilon \right] / x\varepsilon$$

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

if $\varepsilon_x = 0$, $(\text{cond } A)(x) = \frac{e^x}{x}$, is still larger than 1

$$(c) (\text{cond } A)(x) = \left[(-x\epsilon_x + \epsilon_1)/x + \frac{e^x(1-e^{-x})}{x} \epsilon_2 \right] / \epsilon \quad (5)$$

$$= \left[-\epsilon_x + \frac{1}{x} \epsilon_1 + \frac{e^x-1}{x} \epsilon_2 \right] / \epsilon$$

if $\epsilon_x = 0$, then $(\text{cond } A)(x) = \left[\frac{1}{x} \epsilon_1 + \frac{e^x-1}{x} \epsilon_2 \right] / \epsilon$

We can see that the error in exponentiating is expanded by $\frac{1}{x}$
thus it is the cause of poor conditionality

$$(d) \left| \frac{f_A(x) - f(x)}{f(x)} \right| = \left| \frac{e^{-x}}{1-e^{-x}} \left(\cancel{-x\epsilon_x} + \epsilon_1 \right) + \epsilon_2 \right| \quad (\epsilon_x = 0)$$

$$= \left| \frac{e^{-x} \cdot \epsilon_1}{1-e^{-x}} + \epsilon_2 \right| \leq \epsilon \left| \frac{e^{-x}}{1-e^{-x}} + 1 \right| = \left| \frac{e^x}{e^x-1} \right| \cdot \epsilon$$

If we are willing to lose at most 1 bit of significance,

then $\left| \frac{e^x}{e^x-1} \right| \leq 10$, we have $x \geq \cancel{0.1}$

If we are willing to lose at most 2 bits

then $\left| \frac{e^x}{e^x-1} \right| \leq 100$, $\frac{e^x}{e^x-1} \approx \frac{1}{x}$, thus $x \geq 0.01$

for 3 bits, $x \geq 0.001$, for 4 bits, $x \geq 0.0001$

$$(e) \left| \frac{f_A(x) - f(x)}{f(x)} \right| \leq \epsilon \left| \frac{e^x}{e^x-1} \right|$$

(f) we can compute $f(x) = 1 - e^{-x}$ by computer $\frac{e^x-1}{e^x} = \frac{x + \frac{x^2}{2} + \dots}{1 + x + \frac{x^2}{2} + \dots}$
for very small x

7. (2) (i) let $p(\underline{s}_k, \alpha) = 0$

$$\text{then we have } \frac{\partial p}{\partial s_k} \cdot d\underline{s}_k + \sum_{l=0}^{n-1} \frac{\partial p}{\partial a_l} d\underline{a}_l = 0 \quad \textcircled{6}$$

$$\text{thus } \frac{\partial \underline{s}_k}{\partial a_l} = - \frac{\frac{\partial p}{\partial a_l}}{\frac{\partial p}{\partial s_k}} = - \frac{d\underline{s}_k}{p'(s_k)}.$$

$$\text{then } T_{kl} = \left| \frac{a_l \frac{\partial \underline{s}_k}{\partial a_l}}{s_k} \right| = \left| \frac{a_l \underline{s}_k^{l-1}}{p'(s_k)} \right|$$

$$(\text{cond } T_k) (\bar{\alpha}) = \sum_{l=0}^{n-1} \left| \frac{a_l \underline{s}_k^{l-1}}{p'(s_k)} \right|$$

(iii). for $r=14$, the condition number is 5.4×10^{13}

for $r=16$, the condition number is 3.5×10^{13}

for $r=17$, the condition number is 1.8×10^{13}

for $n=20$, the condition number is 1.4×10^{11}

the condition number for this problem is so large that even a small change in the coefficient, the root will change greatly.

(iii) no algorithm can help here, the problem is inherently ill-conditioned, small perturbation of coefficient can cause large changes in the root.

8(a) $\circ y_{n+1} = -(n+1) \circ y_n$, thus

(1) (2) (3)

$$\circ y_N = \frac{N!}{k!} (-1)^{N-k} \circ y_k \quad (k < N)$$

and since $y_{n+1} = \int_0^1 x^{n+1} e^x dx \leq \int_0^1 x^n e^x dx = y_n$

we have $\circ y_N \leq y_k, \quad (k < N)$

that is $\left| \frac{\circ y_N}{y_N} \right| \geq \frac{N!}{k!} \left| \frac{\circ y_k}{y_k} \right|$

so $\varepsilon_{y_k} \leq \frac{k!}{N!} \varepsilon_{y_N}, \quad (\text{cond } g_k)(y_N) \leq \frac{k!}{N!}$

(b) if we want $\varepsilon_{y_k} \leq \varepsilon$, we just need $\frac{k!}{N!} \varepsilon_{y_N} \leq \varepsilon$,

that is $N! \geq \frac{k!}{\varepsilon}$

(c) $\varepsilon = 2^{-53}$, $k = 20$, we have $N! \geq \frac{20!}{2^{-53}} \approx 2.2 \times 10^{34}$

since $31! \approx 8.2 \times 10^{33}$, $32! \approx 2.6 \times 10^{35}$

we have $N = 32$

(d) ~~skip~~ (optional part)

the value of backward recurrence relation is 0.12380383076256993,

the value from direct integration is 0.12380383076256998.

the difference is quite small