

Q1

9 / 10

Give a binary floating point representation with a 4 bit mantissa ($b_1.b_2b_3b_4 \times 2^e$) of the following numbers (use rounding to truncate). Calculate the relative rounding error in each case.

- (a) 3.9
- (b) 0.19
- (c) 3.375
- (d) 8.5

Mahmoud Moustafa
3648228
CS3113 AI

1) d) 8.5
 $8 = (1000)_2$
 $8 - 8 = 0$
 $0.5 \times 2 = 1.0$ | $0.5 = (1)_2$
 $8.5 = (1000.1)_2$
 floating rounding & truncating
 $1.0001 \times 2^3 \Rightarrow 1.0001 \Rightarrow 1.001 \times 2^3$

RRE = $\frac{|x_1 - x|}{|x|}$ $1.001 \times 2^3 \rightarrow 10$
 $\frac{|9 - 8.5|}{8.5} = \frac{|0.5|}{8.5} = 0$

c) 3.375
 $3 = (11)_2$
 $3 - 2 = 1$
 $1 - 1 = 0$
 $0.375 = (0.11)_2$
 $0.375 \times 2 = 0.75$
 $0.75 \times 2 = 1.5$
 $0.50 \times 2 = 1.00$
 $3.375 = (11.011)_2$
 $1.1011 \times 2^1 \Rightarrow 1.1011 \Rightarrow 1.110 \times 2^1$

rounding -1 error.
Should round to $(1.000)_2 \times 2^3$

$1.110 \times 2^1 \rightarrow 11.10 = 2^1 + 2^0 \cdot 2^{-1} = 3.5$

$RRE = \frac{|3.5 - 3.375|}{13.375} = \frac{|0.125|}{13.375} = 0.03704$ ✓

a) 3.9 2^1
2 1

$3 = (11)_2$
 $3 - 2 = 1$
 $1 = 1$
 $0.9 = (0.11100)_2$

$0.9 \times 2 = 1.8$	1
$0.8 \times 2 = 1.6$	1
$0.6 \times 2 = 1.2$	1
$0.2 \times 2 = 0.4$	0
$0.4 \times 2 = 0.8$	0
$0.8 \times 2 = 1.6$	1

$3.9 = (11.11100)_2$

$1.111100 \times 2^1 \Rightarrow 1.11111 \Rightarrow 10.00 \Rightarrow 1.000 \times 2^2$

$1.000 \times 2^2 \Rightarrow 100.0 = 2^2 = 4$

$RRE = \frac{|4 - 3.9|}{13.9} = \frac{|0.1|}{13.9} = 0.02564$ ✓

b) ~~0.9
 $0.9 = (0.11100)_2$
 $0.9 \times 2 = 1.8$
 $0.8 \times 2 = 1.6$
 $0.6 \times 2 = 1.2$
 $0.2 \times 2 = 0.4$
 $0.4 \times 2 = 0.8$
 $0.8 \times 2 = 1.6$~~

1. 1100×2^{-1}

$1100 \times 2^{-1} \rightarrow 0.1100 \times 2^{-1+2+2-3} = 0.875$

$RRE = \frac{1100 \times 2^{-1} - 0.91}{0.91} = 0.02778$

1) b) 0.19

$0.19 = (0.0010000101111111)_2$

$0.19 \times 2 = 0.38$ 0
 $0.38 \times 2 = 0.76$ 0
 $0.76 \times 2 = 1.52$ 1
 $0.52 \times 2 = 1.04$ 1
 $0.04 \times 2 = 0.08$ 0
 $0.08 \times 2 = 0.16$ 0
 $0.16 \times 2 = 0.32$ 0
 $0.32 \times 2 = 0.64$ 0
 $0.64 \times 2 = 1.28$ 1
 $0.28 \times 2 = 0.56$ 0
 $0.56 \times 2 = 1.12$ 1

$1.1000101 \times 2^{-3} \rightarrow 1.100 \times 2^{-3}$

$1100 \times 2^{-3} \rightarrow 0.0011 = 2^{-3} + 2^{-4} = 0.1875$

$RRE = \frac{1100 \times 2^{-3} - 0.19}{0.19} = \frac{0.0025}{0.19} = 0.01316$

Q2

5 / 5

Perform the multiplication

$42.00 \times 1.2344 \times 1600$ in two different orders, using 4 digits and rounding, **base 10**. In both cases, compute the relative error of the result.

a) $(42.00 \times 1.2344) \times 1600$

b) $42.00 \times (1.2344 \times 1600)$

2) a) $(42.00 \times 1.2344) \times 1600$
 $= (42.00 \times 1.2344) \times 1600$
 $= 51.8288 \times 1600$
 $= 51.83 \times 1600$
 $= 82928$
 $= 8.293 \times 10^4 \Rightarrow x_c$ $82951.68 \rightarrow x$
 $RRE = \frac{|x_c - x|}{|x|}$
 $= \frac{|82930 - 82951.68|}{|82951.68|}$
 $= 0.0002614 \quad \checkmark$

b) $42.00 \times (1.2344 \times 1600)$
 $= 42.00 \times (1.2344 \times 1600)$
 $= 42.00 \times 1974.4$
 $= 42.00 \times 1974$
 $= 82908$
 $= 8.291 \times 10^4 \Rightarrow x_c$ $82951.68 \rightarrow x$
 $RRE = \frac{|82910 - 82951.68|}{|82951.68|}$
 $= 0.0005025 \quad \checkmark$

good! 5

Q3

3 / 4

Subtract $\sqrt{8280}$ from 91 using 4 digits base 10 and rounding. Compute the relative error and comment on its magnitude.

3) 4 digit base 10 & rounding

$$91 - \sqrt{8280} \quad \sqrt{8280} = 90.99 \text{ rounded } \downarrow$$

$$91 - 90.99$$

$$= 0.01$$

$$91 - \sqrt{8280} = 0.00549467138$$

$$RRE = \frac{0.01 - 0.00549467138}{0.00549467138}$$

$$= 0.8199$$

RRE is too high: 81.99 %

Ans: large
error due to
catastrophic
cancellation

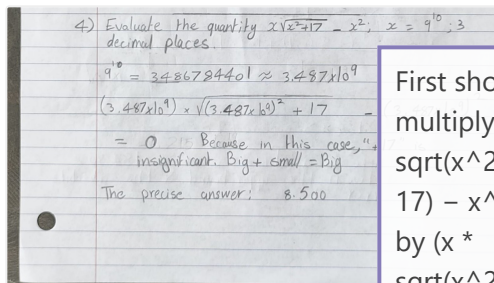
-1

Q4

0 / 5

(textbook, section 0.4, exercise 4)

Evaluate the quantity

 $x\sqrt{x^2 + 17} - x^2$ where $x = 9^{10}$,
correct to at least 3 decimal places.


First should **-5**
multiply $x \cdot$
 $\sqrt{x^2 +$
 $17} - x^2$
by $(x \cdot$
 $\sqrt{x^2 +$
 $+17) +$
 $x^2)/(x \cdot$
 $\sqrt{x^2 +$
 $+17) +$
 $x^2)$. After
simplifying,
then
substituting
 $x = 9^{10}$.

Q5

9 / 10

Consider a hypothetical 8 bit floating point machine representation with a sign bit, a 3 bit exponent, and a 4 bit mantissa ($se_1e_2e_3b_1b_2b_3b_4$), where the exponent bias is 3 (add 3 to exponent of number to form machine representation). Recall that actual mantissa has 5 bits, since the leading 1 is not stored on the machine.

(a) What is the number $e \approx 2.718$ in this 8-bit format?

(b) What is the number that $(10100111)_2$ represents in this 8-bit format?

(c) What is the upper bound of the relative error when representing a real number in this 8-bit format?

5)

a) 2.718

$2 = (10)_2$

0.718 =		
0.718 x 2 = 1.436	1	
0.436 x 2 = 0.872	0	
0.872 x 2 = 1.744	1	
0.744 x 2 = 1.488	1	
0.488 x 2 = 0.976	0	
0.976 x 2 = 1.952	1	
0.952 x 2 = 1.904	1	
0.904 x 2 = 1.808	1	
0.808 x 2 = 1.616	1	
0.616 x 2 = 1.232	1	
0.232 x 2 = 0.464	0	
0.464 x 2 = 0.928	0	
0.928 x 2 = 1.856	1	
0.856 x 2 = 1.712	1	
	!	
	!	
	!	

2.718 = (10.101101111...)

= 1.0101101111 x 2¹

= 1.0110 x 2¹

sign = 0

bias = 3

exponent = 3 + 1 = 4

mantissa = 0110

0100110 ⇒ 8-bit representation (format)

b) (10100111) to decimal

1 010 0111
 sign exponent mantissa

exponent = $2 - 3 = -1$ (010)₂ = bias
 mantissa = 1.0111

$1.0111 \times 2^{-1} = 0.10111 = 2^{-1} + 2^{-3} + 2^{-4} + 2^{-5}$
 $= 0.71875$

c)

$E_{\text{mach}} = 2^{-6} \rightarrow 2^{-4} = \frac{1}{16}$
 upper bound of the relative error = $\frac{1}{2}$
 $= \frac{1}{2}$
 $= \frac{1}{32}$

sign
 not
 correct

Q6

6 / 6

(textbook, section 0.5, exercise 2)

Find c satisfying the Mean Value Theorem for $f(x)$ on the interval $[0, 1]$. (a) $f(x) = e^x$ (b) $f(x) = x^2$ (c) $f(x) = 1/(x + 1)$

Handwritten solution for the Mean Value Theorem problem:

6) $[0, 1]$

a) $f(x) = e^x \Rightarrow f'(x) = e^x$

$$f'(c) = \frac{f(b) - f(a)}{b - a}$$

$$= \frac{e^1 - e^0}{1 - 0} = \frac{e - 1}{1} = e - 1$$

$f'(x) = e^x$

Solve $e^c = e - 1 \Rightarrow \ln(e^c) = \ln(e - 1)$

$$c = \ln(e - 1) \approx 0.5413$$

b) $f(x) = x^2 \Rightarrow f'(x) = 2x$

$$f'(c) = \frac{f(b) - f(a)}{b - a} = \frac{1 - 0}{1} = 1$$

$f'(x) = 2x$

Solve $2c = 1$

$$c = \frac{1}{2} = 0.5$$

c) $f(x) = \frac{1}{x+1} \Rightarrow f'(x) = \frac{-1}{(x+1)^2}$

$f'(c) = \frac{f(b) - f(a)}{b - a}$ $f(b) = 0.5$
 $f(a) = 1$

$= \frac{0.5 - 1}{1 - 0} = \frac{-0.5}{1} = -0.5$

Solve $\frac{-1}{(c+1)^2} = -0.5 =$

$\frac{1}{(c+1)^2} = 0.5 =$

$1 = 0.5 \times (c+1)^2$

$2 = (c+1)^2$

↓ ↓

$\sqrt{2} = c+1$ $-\sqrt{2} = c+1$

$\sqrt{2} - 1 = c$ $-\sqrt{2} - 1 = c$

$0.4142 = c$ $-2.4142 = c$ X

↳ ✓ ✓ ✓ in interval less than 0 ←

Q7

7 / 8

Determine the second-degree Taylor polynomial and associated remainder term for the function $f(x) = e^{-x^2}$, expanding about 0, and use it to estimate $e^{-0.1^2}$.

Compare the upper bound of the remainder to the exact absolute error.

$$f(x) = e^{-x^2}$$

$$f'(x) = (e^{-x^2})(-2x)$$

$$f''(x) = (e^{-x^2})(-2x) + (-2)(e^{-x^2}) = e^{-x^2}(-2x^2 - 2)$$

$$f'''(x) = (e^{-x^2})(-4x) + (-2)(e^{-x^2})(-2x) = e^{-x^2}(-4x^2 + 4x)$$

$$f(0) = 1 \quad f'(0) = 0 \quad f''(0) = -2 \quad f'''(0) = 0$$

$$p_2(x) = 1 + 0(x-0) + \frac{-2}{2!}(x-0)^2 = 1 - \frac{2}{2}x^2 = 1 - x^2$$

$$p_2(0.1) = 1 - (0.1)^2 = 1 - 0.01 = 0.99$$

$$\text{absolute error} = 0.99 - 0.9900498 = 0.0000498$$

$$f'''(c) < f'''(x) \quad \frac{f'''(c)(x-x_0)^3}{3!}$$

$$\boxed{0 < c < 0.1}$$

$$\frac{e^{-0.01}(-4(0.1)^2 + 4(0.1))}{3!} = \frac{0.01848}{3} = 0.00616$$

$$0.00616 > 0.0000498$$

There is a difference of 0.0184302 between the upper bound & the absolute error

remainder not
correct

