

Instructor: Rich Little

Question #1 - 8 marks.

Consider a base 5 normalized, floating-point number system. Assume that a hypothetical computer using this system has the following floating-point representation: $s m f_1 f_2 f_3 f_4 s e_1 e_2$ where $s m$ is the sign of the mantissa, $s e$ is the sign of the exponent (1 for negative, 0 for positive), f_i are the digits of the mantissa, and e_j are the digits of the exponent.

- (a) Consider the base 5 number, given using the above representation, 02003004. What exact decimal value does it represent?

Handwritten solution for part (a):

Leoza Kabir
V00840048
CSC 349A
Date: Sep 27, 2018
Assignment #2

1 a) $02003004 = (0.2003 \times 5^4)_5 = (2003)_5$
 $= 2 \times 5^3 + 0 \times 5^2 + 0 \times 5^1 + 3 \times 5^0$
 $= 250 + 0 + 0 + 3 = (253)_{10}$

- (b) What decimal value does 11004003 represent?

Handwritten solution for part (b):

1 b) $11004003 = (-0.11004 \times 5^3)_5 = (-100.4)_5$
 $= -1 \times 5^2 + 0 \times 5^1 + 0 \times 5^0 + 4 \times 5^{-1}$
 $= -25 + 0 + 0 + \frac{4}{5}$
 $= (-24.8)_{10}$

- (c) What is the smallest positive, non-zero, number that can be represented in this system? Give the answer in the above form (i.e. as 8 base-5 digits.) and in decimal.

Handwritten solution for part (c):

1 c) $s m f_1 f_2 f_3 f_4 s e_1 e_2$
 $0 1 0 0 0 1 1 4 4$
 $(0.1000144)_5 = 1 \times 5^{-4} = (3.51844 \times 10^{-32})_{10}$

- (d) What is the size of the gap between any two consecutive numbers in the interval $25(10)$ and $125(10)$ in this floating-point representation system? Your answer should be in decimal.

Instructor: Rich Little

d) $[5^2, 5^3]$ size of interval $= 5^3 - 5^2 = 100$
 $\#$ of sub intervals $= (b-1)b^{k-1} = (5-1)(5)^{4-1} = 500$
 size of gap $= b^{-k} = 5^{3-4} = 5^{-1} = 5^{-1} = 1/5 = 0.2$

Question #2 - 6 Marks.

The polynomial $P(x) = x^2 - 83.12x + 3.123$ has two roots, at approximately 0.0375892 and 83.0824. The roots of a quadratic polynomial $ax^2 + bx + c$ can be computed by (i) $-b \pm \sqrt{b^2 - 4ac}$ or equivalently (ii) $-2c \pm \sqrt{b^2 - 4ac}$. Using floating-point arithmetic, one of these formulas is often much more accurate than the other. For example, if $(-b + \sqrt{b^2 - 4ac})/(2a)$ is used to compute one of the roots of $P(x) = x^2 - 83.12x + 3.123 = 0$ with base $b = 10$, precision $k = 4$, idealized chopping arithmetic, the results are as follows: $\text{fl}(b^2) = \text{fl}((6908.9344)) = 6908$ or 0.6908×10^4 . $\text{fl}(4a) = 4$ or 0.4000×10^1 . $\text{fl}(4ac) = \text{fl}(4 \times 3.123) = \text{fl}(12.492) = 12.49$. $\text{fl}(b^2 - 4ac) = \text{fl}((6908 - 12.49)) = \text{fl}((6895.51)) = 6895$. $\text{fl}(\sqrt{b^2 - 4ac}) = \text{fl}(\sqrt{6895}) = \text{fl}(83.036136 \dots) = 83.03$. $\text{fl}(-b + \sqrt{b^2 - 4ac}) = \text{fl}(83.12 + 83.03) = \text{fl}(166.15) = 166.1$. $\text{fl}(2a) = 2$. $\text{fl}(-b + \sqrt{b^2 - 4ac}) / 2a = \text{fl}(166.1 / 2) = 83.05$ or 0.8305×10^2 which is very accurate. The relative error is about 0.00039 or 0.039%. On the other hand, it can be shown (similar to the above) that $\text{fl}(-2c \pm \sqrt{b^2 - 4ac}) / \sqrt{b^2 - 4ac} = 69.40$ or 0.6940×10^2 which (using the exact value of 83.0824 ...) has a large relative error of 0.165 or 16.5%.

- (a) Use base $b = 10$, precision $k = 4$, idealized chopping arithmetic and each of the mathematically equivalent formulas $-2c \pm \sqrt{b^2 - 4ac}$ and $-b \pm \sqrt{b^2 - 4ac}$ to compute an approximation to one root of $P(x) = 1.2x^2 - 78.99x + 1.234 = 0$. As above, specify each step of the computation. Note that many of the computations for the two formulas are identical, and need only be done once. Use your calculator to do this, not MATLAB.

2a) $\text{fl}(b^2) = \text{fl}(78.99^2) = \text{fl}(6239.4201) = 6239$ or 0.6239×10^4
 $\text{fl}(4a) = \text{fl}(4 \times 1.2) = \text{fl}(4.8) = 4.8$ or 0.4800×10^1
 $\text{fl}(4ac) = \text{fl}(4 \times 1.234) = \text{fl}(4.936) = 4.936$
 $\text{fl}(b^2 - 4ac) = \text{fl}(6239 - 4.936) = \text{fl}(6234.064) = 6234$
 $\text{fl}(\sqrt{b^2 - 4ac}) = \text{fl}(\sqrt{6234}) = \text{fl}(78.94983124) = 78.94$
 $\text{fl}(-b - \sqrt{b^2 - 4ac}) = \text{fl}(-78.99 - 78.94) = -157.9$
 $\text{fl}(-2c) = \text{fl}(-2 \times 1.234) = \text{fl}(-2.468) = -2.468$
 $\text{fl}\left(\frac{-2c}{-b - \sqrt{b^2 - 4ac}}\right) = \text{fl}\left(\frac{-2.468}{-157.9}\right) = 0.01563$
 $\text{fl}(b^2 - 4ac) = 6234$
 $\text{fl}(-b - \sqrt{b^2 - 4ac}) = \text{fl}(-78.99 - 78.94) = -157.9$
 $\text{fl}(2a) = \text{fl}(2 \times 1.2) = \text{fl}(2.4) = 2.4$ or 0.2400×10^1
 $\text{fl}\left(\frac{-b - \sqrt{b^2 - 4ac}}{2a}\right) = \text{fl}\left(\frac{-157.9}{2.4}\right) = \text{fl}(-65.79166666) = -65.79$
 $\text{fl}\left(\frac{-b + \sqrt{b^2 - 4ac}}{2a}\right) = \text{fl}\left(\frac{0.05}{2.4}\right) = \text{fl}(0.020833333) = 0.02083$

- (b) Compute the relative errors of each of the approximations in (a) using the fact that the exact value of the root is $0.01562594 \dots$. Give at least 2 significant digits.

b) Relative error for $\frac{-b - \sqrt{b^2 - 4ac}}{2a}$ method

$$|E_r| = \left| \frac{0.01562594 - 0.01563}{0.01562594} \right| = 0.000259829 = 0.026\%$$

Relative error for $\frac{-b + \sqrt{b^2 - 4ac}}{2a}$ method

$$|E_r| = \left| \frac{0.01562594 - 0.02083}{0.01562594} \right| = 0.3330398 \approx 33.30\%$$

- (c) One of the two zeros of a quadratic polynomial $ax^2 + bx + c$ can be computed using either the formula (i) $\frac{-b + \sqrt{b^2 - 4ac}}{2a}$ or (ii) $\frac{-2c}{b + \sqrt{b^2 - 4ac}}$. For each of the specified polynomials in the table below, place an X in the appropriate box to indicate which of these formulas is more accurate in precision $k = 4$ floatingpoint arithmetic. Put exactly one X in each row of the table. (No justification for your answers is required. It is NOT necessary to do any floating-point computation to answer this question.)

c)

Polynomial	(i) is more accurate	(ii) is more accurate
$0.01x^2 - 125x + 0.05$		X
$-0.3x^2 + 125x + 0.005$	X	

Instructor: Rich Little

Question #3 - 6 Marks

- (a) Determine the second order ($n = 2$) Taylor series expansion for $f(x) = \sqrt{x+3}$ expanded about $a = 1$ including the remainder term. Leave your answer in terms of factors $(x - 1)$ (that is, do not simplify). Show all your work.

Handwritten solution for part (a):

$$\begin{aligned}
 &3) a) \quad f(x) = \sqrt{x+3} \quad f(1) = \sqrt{1+3} = \sqrt{4} = 2 \\
 &f'(x) = (x+3)^{-1/2} = \frac{1}{2(x+3)} \quad f'(1) = \frac{1}{2\sqrt{1+3}} = \frac{1}{2\sqrt{4}} = \frac{1}{4} \\
 &f''(x) = -\frac{1}{2}(x+3)^{-3/2} = -\frac{1}{4(x+3)^{3/2}} \quad f''(1) = -\frac{1}{4(1+3)^{3/2}} = -\frac{1}{4 \times 8} = -\frac{1}{32} \\
 &f'''(x) = -\frac{1}{4}(x+3)^{-5/2} = -\frac{3}{8(x+3)^{5/2}} \quad f'''(1) = -\frac{3}{8(1+3)^{5/2}} = -\frac{3}{256} \\
 &f(x) \approx f(1) + f'(1)(x-1) + \frac{f''(1)}{2!}(x-1)^2 + \frac{f'''(1)}{3!}(x-1)^3 \\
 &= 2 + \frac{1}{4}(x-1) - \frac{1}{64}(x-1)^2 - \frac{3}{16384}(x-1)^3 \\
 &= 2 + \frac{1}{4}(x-1) - \frac{1}{64}(x-1)^2 + \frac{1}{16384}(x-1)^3
 \end{aligned}$$

- (b) Use the polynomial approximation in (a) (without the remainder term) to approximate $f(1.12) = \sqrt{4.12}$. Use either hand computation, your calculator or MATLAB. Give an exact answer. 3

Handwritten solution for part (b):

$$\begin{aligned}
 &b) \quad \sqrt{4.12} = f(1.12) \\
 &\approx 2 + \frac{1}{4}(1.12-1) - \frac{1}{64}(1.12-1)^2 = 2.029975
 \end{aligned}$$

Instructor: Rich Little

c) Determine a good upper bound for the truncation error of the Taylor polynomial approximation in (a) for all values of x such that $1 \leq x \leq 1.2$ by bounding the remainder term.

The image shows a handwritten solution on lined paper. It starts with the formula for the remainder term R_2 and the interval E_2 . Then it substitutes the values to calculate the upper bound, and finally states the inequality for the truncation error $|E_t|$.

$$\begin{aligned} \textcircled{1} \quad R_2 &= \frac{1}{16(E_2+3)^{5/2}} (x-1)^3 & E_2 &= [1, 1.2] \\ &= \frac{(1.2-1)^3}{16(1+3)^{5/2}} = 0.000015625 \\ |E_t| &\leq |R_3| = 0.000015625 \end{aligned}$$