

Floating Point

1. a) # elements of $F = (1023+1+1022) \cdot 2^{51} \cdot 2$
 $= 2046 \cdot 2^{52}$

b) # elements where $1 \leq x \leq 2$
 $= 2^{51}$

$$\text{Fraction} = \frac{2^{51}}{2046 \cdot 2^{52}}$$

$$= \frac{1}{4092}$$

c) $(\frac{1}{4})_{10} = (.1 \times 2^{-1})_2$

$$(\frac{1}{2})_{10} = (.1 \times 2^0)_2$$

\therefore # elements where $\frac{1}{4} \leq x < \frac{1}{2}$
 $= 2^{51}$

$$\text{Fraction} = \frac{2^{51}}{2046 \cdot 2^{52}}$$

$$= \frac{1}{4092}$$

d) refer to page 7.

$$2. \frac{|f_l(a) \oplus f_l(b) - (a+b)|}{|a+b|}$$

$$= \frac{|(a(1+\delta_1) \oplus b(1+\delta_2)) - (a+b)|}{|a+b|} \quad \text{where } |\delta_1|, |\delta_2| \leq \epsilon$$

$$= \frac{|((a(1+\delta_1) + b(1+\delta_2))(1+\delta_3) - (a+b))|}{|a+b|} \quad \text{where } |\delta_3| \leq \epsilon$$

$$= \frac{|(a + a\delta_1 + b + b\delta_2)(1 + \delta_3) - (a+b)|}{|a+b|}$$

$$= \frac{|(\cancel{a+b}) + a\delta_1 + b\delta_2 + (a+b)\delta_3 + (a\delta_1 + b\delta_2)\delta_3 - \cancel{(a+b)}|}{|a+b|}$$

$$= \frac{|a(\delta_1 + \delta_3 + \delta_1\delta_3) + b(\delta_2 + \delta_3 + \delta_2\delta_3)|}{|a+b|}$$

$$\leq \frac{|a| |\delta_1 + \delta_3 + \delta_1\delta_3| + |b| |\delta_2 + \delta_3 + \delta_2\delta_3|}{|a+b|} \quad \begin{array}{l} \text{Since } |\delta_i| \leq \epsilon \\ \text{for } i=1,2,3 \end{array}$$

$$\leq \frac{|a|(2\epsilon + \epsilon^2) + |b|(2\epsilon + \epsilon^2)}{|a+b|} \leq \frac{|a| + |b|}{|a+b|} \epsilon(2 + \epsilon)$$

$$\begin{array}{r}
 3. a) \quad 149,659,299,043,739,794 \\
 - 149,659,299,043,739,691 \\
 \hline
 0.103 \text{ mm}
 \end{array}$$

\therefore Thickness is 0.103 mm

b) From MATLAB: thickness = 0 mm

$$\text{Rel. Err.} = \frac{|0 - 0.103|}{0.103} = 1 = 100\%$$

c) From MATLAB: thickness = 0.0938 mm

$$\begin{aligned}
 \text{Rel. Err.} &= \frac{|0.0938 - 0.103|}{0.103} \approx 0.0873204 \\
 &= 8.73204\%
 \end{aligned}$$

d) Upper Bound

$$\text{Rel. Err} \leq \frac{|149659299043739.794| + |-149659299043739.691|}{|149659299043739.794 - 149659299043739.691|} E(2+E)$$

$$\begin{aligned}
 &\leq \frac{299318458087499.485}{0.103} E(2+E) \quad \left| \begin{array}{l} E = \frac{1}{2} B^{1-z} \\ = 2^{-1} 2^{1-52} \\ = 2^{-52} \end{array} \right. \\
 &\leq \frac{299318458087499.485}{0.103} 2^{-52} (2 + 2^{-52})
 \end{aligned}$$

$$\begin{aligned}
 &\leq 1.29052522 \\
 &\leq 129\%
 \end{aligned}$$

4. refer to $\text{page } 8$

$$5. a) x^2 + 100.01x + 1.2121 = 0$$

using calculator:

$$r_1 = -0.01212125712$$

$$r_2 = -99.9978787429$$

$$b) r_1 = \frac{-100.01 \oplus \sqrt{100.01^2 \ominus 4 \otimes 1 \otimes 1.2121}}{2 \otimes 1}$$

$$= \frac{-100.01 \oplus \sqrt{\text{fl}(10002.0001) \ominus \text{fl}(4.8484)}}{\text{fl}(2)}$$

$$= \frac{-100.01 \oplus \sqrt{\underline{10002} \ominus 4.8484}}{2}$$

$$= \frac{-100.01 \oplus \sqrt{\text{fl}(9997.1516)}}{2}$$

$$= \frac{-100.01 \oplus \sqrt{\underline{9997.2}}}{2}$$

$$= \frac{-100.01 \oplus \text{fl}(99.9859990199)}{2}$$

$$r_1 = \frac{-100.01 \oplus \underline{99.986}}{2}$$

$$= \text{fl}(-0.0240) \oplus 2$$

$$= -0.0240 \oplus 2$$

$$= \text{fl}(-0.0120)$$

$$r_1 = -0.0120$$

using the same discriminant

$$r_2 = \frac{-100.01 \ominus \underline{99.986}}{2}$$

$$= \text{fl}(-199.996) \oplus 2$$

$$= \underline{-200.00} \oplus 2$$

$$= \text{fl}(-100.00)$$

$$r_2 = -100.00$$

Red underline represents loss of precision

$$r_1 \text{ Rel. Err.} = \frac{|(-0.0120) - (-0.01212125712)|}{0.01212125712}$$

$$= 0.01000367526 \approx 1\%$$

$$r_2 \text{ Rel. Err.} = \frac{|(-100.00) - (-99.9978787429)|}{99.9978787429}$$

$$= 2.121302 \times 10^{-5} \approx 0.002\%$$

$$\begin{aligned} c) \quad r_1 &= \frac{-b + \sqrt{b^2 - 4ac}}{2a} \cdot \frac{-b - \sqrt{b^2 - 4ac}}{-b - \sqrt{b^2 - 4ac}} \\ &= \frac{\cancel{b^2} - (\cancel{b^2} - 4ac)}{2a(-b - \sqrt{b^2 - 4ac})} \\ &= \frac{\cancel{4ac}^2}{\cancel{2a}(-b - \sqrt{b^2 - 4ac})} \\ &= \frac{2c}{-b - \sqrt{b^2 - 4ac}} \end{aligned} \quad \left| \quad \begin{aligned} r_2 &= \frac{-b - \sqrt{b^2 - 4ac}}{2a} \cdot \frac{-b + \sqrt{b^2 - 4ac}}{-b + \sqrt{b^2 - 4ac}} \\ &= \frac{\cancel{b^2} - (\cancel{b^2} - 4ac)}{2a(-b + \sqrt{b^2 - 4ac})} \\ &= \frac{\cancel{4ac}^2}{\cancel{2a}(-b + \sqrt{b^2 - 4ac})} \\ &= \frac{2c}{-b + \sqrt{b^2 - 4ac}} \end{aligned} \right.$$

$$d) \quad r_1 = \frac{2 \otimes 1.2121}{-100.01 \ominus \sqrt{100.01^2 - 4 \otimes 1 \otimes 1.2121}}$$

using the same discriminant calculated in b)

$$\begin{aligned} r_1 &= \frac{f1(2.4242)}{-100.01 \ominus \underline{99.986}} \\ &= \frac{2.4242}{f1(-199.996)} \end{aligned}$$

$$\begin{aligned} r_2 &= \frac{f1(2.4242)}{-100.01 \oplus \underline{99.986}} \\ &= \frac{2.4242}{f1(-0.0240)} \end{aligned}$$

$$r_1 = 2.4242 \oplus \underline{-200.00}$$

$$= fl(-0.012121)$$

$$= \underline{-0.012121}$$

$$r_2 = 2.4242 \oplus -0.0240$$

$$= fl(-101.008333333)$$

$$= \underline{-101.01}$$

$$r_1 \text{ Rel. Err.} = \frac{|(-0.012121) - (-0.01212125712)|}{0.01212125712}$$

$$= 2.121232 \times 10^{-5} \approx 0.002\%$$

$$r_2 \text{ Rel. Err.} = \frac{|(-101.01) - (-99.9978787429)|}{99.9978787429}$$

$$= 0.01012142727 \approx 1\%$$