## Question 1

[15 marks]

Consider the normalized floating-point system $\mathbb{R}_3(3, 1)$ with *limited* exponent range $-1 \leq e \leq 1$.

a. What is the smallest positive (nonzero) number representable? Give your answer in both base-3 and base-10.

$$\text{Base-3}: \quad (0.100)_3 \times 3^{-1}$$

$$\text{Base-10}: \quad 1 \times 3^{-2} = \frac{1}{9} = 0.\overline{1}$$

b. What is the largest positive number representable? Give your answer in both base-3 and base-10.

$$\text{Base-3}: \quad (0.222)_3 \times 3^{1}$$

$$\text{Base-10}: \quad 2 \cdot 3^0 + 2 \cdot 3^{-1} + 2 \cdot 3^{-2} = 2 + \frac{2}{3} + \frac{2}{9} = 2 + \frac{6+2}{9} = 2 + \frac{8}{9} = 2.\overline{8}$$

c. Assuming round-to-nearest, what is the tightest upper bound on the relative error $|fl(x) - x|/|x|$ when $x \in \mathbb{R}$ is stored as $fl(x) \in \mathbb{R}_3(3, 1)$ in this floating-point system? Give your answer in base-10.

$$\frac{|fl(x) - x|}{|x|} = \varepsilon < \frac{1}{2} b^{1-t} = \frac{1}{2} \cdot 3^{1-3} = \frac{1}{2} \cdot \frac{1}{9} = \frac{1}{18}$$

d. What is the floating-point representation of $(407)_{10}$ in this system? (*Hint*: Does the representation exist?)

There is no representation of $(407)_{10}$ in this system since the largest positive number representable is $2.\overline{8}$

e. What is the floating-point representation of $(0.567)_{10}$ in this system? Give your answer in base-3. Recall that there are only *three* base-3 digits in the mantissa.

| | | | | | |
|---|---|---|---|---|---|
| 0.567 | 3 | 1.701 | 1 | 0.701 | |
| 0.701 | 3 | 2.103 | 2 | 0.103 | |
| 0.103 | 3 | 0.309 | 0 | 0.309 | |
| 0.309 | 3 | 0.927 | 0 | 0.927 | |
| 0.927 | 3 | 2.781 | 2 | 0.781 | |
| 0.781 | 3 | 2.343 | 2 | 1 | |

So $(0.567)_{10} \approx (0.120)_3$ in this system

f. List all possible normalized, non-zero mantissas in this system. In total, how many floating-point numbers are representable? Recall that the exponent range is *limited*.

Normalized so: combinations excluding 0.000 are

$0. \underline{2} \times \underline{3} \times \underline{3} = 18$ possible numbers

cannot be 0,
so 1 or 2

$\times 2$    because of sign

$\times 3$    because of $e^{-1}, e^0$ or $e^1$ multiplied

For a total of $18 \times 2 \times 3 + 1 = 109$ representable numbers in this system

$\downarrow$
$(0.000)_3$

# Question 2

To convert to base 10, we can just do:

$$d_0 \times b^0 + d_1 \times b^1 + \ldots + d_k \times b^k$$

$$= d_0 + b(d_1 + b(d_2 + \ldots + b(d_{k-1} + bd_k))\ldots))$$

1 flop

2 flops

k

k-1

Thus, we can clearly see that to convert a $(k+1)$-digit in base $b$, we would just need $k$ flops

# Question 3

[15 marks]

Consider the linear system $Ax = b$ where

$$A = \begin{bmatrix} 2 & 5 & 10 \\ 8 & 32 & 8 \\ 1 & 8 & 13 \end{bmatrix}, \quad b = \begin{bmatrix} 7 \\ -16 \\ 6 \end{bmatrix}.$$

a. Compute the $PA = LU$ factorization of $A$. Use exact arithmetic. Show all intermediate calculations, including Gauss transforms and permutation matrices.

$$P_{12} A = \begin{bmatrix} 8 & 32 & 8 \\ 2 & 5 & 10 \\ 1 & 8 & 13 \end{bmatrix} \qquad L_1 = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{2}{8} & 1 & 0 \\ -\frac{1}{8} & 0 & 1 \end{bmatrix} \qquad L_1 P_{12} A = \begin{bmatrix} 8 & 32 & 8 \\ 0 & -3 & 8 \\ 0 & 0 & 11 \end{bmatrix}$$

So Let $P = P_{12}$ $\quad L = L_1^{-1}$

$$L_1 P_{12} A = U \implies P_{12} A = L_1^{-1} U \implies PA = LU$$

**b.** Use the factorization computed in (a) to solve the system.

$$P A x = P b \implies L U x = P b \qquad \text{Let } U x = d \qquad L d = P b$$

$$\begin{bmatrix} 1 & 0 & 0 \\ \frac{2}{8} & 1 & 0 \\ \frac{1}{8} & 0 & 1 \end{bmatrix} d = \begin{bmatrix} -16 \\ 7 \\ 6 \end{bmatrix} \qquad 1 d_1 = -16 \qquad \frac{2}{8} d_1 + d_2 = 7 \implies \boxed{d_2 = 11}$$

$$\frac{1}{8} d_1 + d_3 = 6 \implies \boxed{d_3 = 8}$$

So,

$$\begin{bmatrix} 8 & 32 & 8 \\ 0 & -3 & 8 \\ 0 & 0 & 11 \end{bmatrix} x = \begin{bmatrix} -16 \\ 11 \\ 8 \end{bmatrix} \qquad 11 x_3 = 8 \implies \boxed{x_3 = \frac{8}{11}}$$

$$-3 x_2 + 8 x_3 = 11 \implies \boxed{x_2 = \frac{-19}{11}}$$

$$8 x_1 + 32 x_2 + 8 x_3 = -16 \implies \boxed{x_1 = \frac{46}{11}}$$

**c.** Instead of first computing the $PA = LU$ factorization, we could have solved the system above by processing the left and right-hand sides simultaneously with Gauss transforms and permutations. Would this alternate approach incur any extra cost? **Explain.** We are solving one system only.

$$\implies L_1 P_{12} A x = L_1 P_{12} b$$

Solving with $PA = LU$ factorization cost $\frac{n^3}{3} + O(n^2)$ as shown in class whereas here

## Question 7

[15 marks]

Consider the data points $\{(0, 3), (1, 7), (2, 37), (3, 141)\}$.

a. Set up the Vandermonde system for determining the monomial-basis form of the polynomial which interpolates these data points. Do **not** solve the system.

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 8 \\ 1 & 3 & 9 & 27 \end{bmatrix} a = \begin{bmatrix} 3 \\ 7 \\ 37 \\ 141 \end{bmatrix}$$

b. Derive the Newton form of the interpolating polynomial. Show all of your work, including the divided-difference table.

$$X_i \quad y[X_i] \quad y[X_{i+1}, X_i] \quad y[X_{i+2}...X_i] \quad y[X_{i+3}...X_i]$$

| $X_i$ | $y[X_i]$ | $y[X_{i+1}, X_i]$ | $y[X_{i+2}...X_i]$ | $y[X_{i+3}...X_i]$ |
|-------|----------|-------------------|--------------------|--------------------|
| 0 | 3 | | | |
| 1 | 7 | $\frac{7-3}{1-0} = 4$ | $\frac{30-4}{2-0} = 13$ | $\frac{37-13}{3-0} = 8$ |
| 2 | 37 | $\frac{37-7}{2-1} = 30$ | $\frac{104-30}{3-1} = 37$ | |
| 3 | 141 | $\frac{141-37}{3-2} = 104$ | | |

So Polynomial is $P(x) = 3 + 4(x) + 13(x)(x-1) + 8(x)(x-1)(x-2)$

c. Derive the Lagrange form of the interpolating polynomial. Verify it is the same polynomial as in (b).

$$P(x) = \sum_{i=0}^{n} \ell_i(x) \, y_i = \ell_0(x) \, y_0 + \ell_1(x) \, y_1 + \ell_2(x) \, y_2 + \ell_3(x) \, y_3$$

$$\ell_0(x) = \prod_{\substack{j=0 \\ j \neq 0}}^{n} \frac{x - x_j}{x_i - x_j} = \left(\frac{x-1}{0-1}\right)\left(\frac{x-2}{0-2}\right)\left(\frac{x-3}{0-3}\right) = \frac{1}{-6}(x-1)(x-2)(x-3)$$

$$\ell_1(x) = \left(\frac{x-0}{1-0}\right)\left(\frac{x-2}{1-2}\right)\left(\frac{x-3}{1-3}\right) = \frac{1}{2}(x)(x-2)(x-3)$$

$$\ell_2(x) = \left(\frac{x-0}{2-0}\right)\left(\frac{x-1}{2-1}\right)\left(\frac{x-3}{2-3}\right) = -\frac{1}{2}(x)(x-1)(x-3)$$

$$\ell_3(x) = \left(\frac{x-0}{3-0}\right)\left(\frac{x-1}{3-1}\right)\left(\frac{x-2}{3-2}\right) = \frac{1}{6}(x)(x-1)(x-2)$$

$$P(x) = -\frac{1}{2}(x-1)(x-2)(x-3) + \frac{7}{2}(x)(x-2)(x-3) - \frac{37}{2}(x)(x-1)(x-3) + \frac{141}{6}(x)(x-1)(x-2)$$

d. Briefly discuss the relative efficiency of the methods in (a), (b), and (c). Which method is best if we need to include additional data points? Explain.

Divided difference since if we add additional data points we do not need To erase anything From our divided-difference table but just require to add the corresponding additional calculation unlike the other methods in which we would need to recalculate everything From scratch

**e.** Construct the linear spline (i.e., the piecewise linear interpolant) which interpolates all four data points $\{(0,3),(1,7),(2,37),(3,141)\}$.

First line: $y = 3 + \left(\frac{7-3}{1-0}\right)(x-0) = 3 + 4x$

Second line: $y = 7 + \left(\frac{37-7}{2-1}\right)(x-1) = 7 + 30(x-1) = -23 + 30x$

Third line: $y = 37 + \left(\frac{141-37}{3-2}\right)(x-2) = 37 + (104)(x-2)$

$$= -171 + 104x$$

So spline is

$$f(x) = \begin{cases} 3+4x & , \text{ if } 0 \le x \le 1 \\ -23+30x & , \text{ if } 1 < x \le 2 \\ -171+104x & , \text{ if } 2 < x \le 3 \end{cases}$$

we could extend this to cover all of $\mathbb{R}$

**END OF EXAM**