

# Numerical Methods 2023 — Midterm 1

## Solutions

**Problem 1 (20 pts).** In our course, we have already learned double and single precision floating-point formats. Nevertheless, there exists another useful floating-point format that uses only 2 Bytes on the storage:

sign	exponent	significant
1 bit	5 bits	10 bits

which is called half precision floating-point, and this format has the rules:

$$\left\{ \begin{array}{ll} \pm 0, & \text{exponent} = 00000_2 \text{ and significant} = 0 \\ \pm 2^{-14} \times 0.\text{significant}_2, & \text{exponent} = 00000_2 \text{ and significant} \neq 0 \\ \pm 2^{\text{exponent}-15} \times 1.\text{significant}_2, & 11111_2 > \text{exponent} > 00000_2 \\ \pm \infty, & \text{exponent} = 11111_2 \text{ and significant} = 0 \\ \text{NaN}, & \text{exponent} = 11111_2 \text{ and significant} \neq 0 \end{array} \right. \quad (1)$$

We further consider rounding even in this problem.

(a) (5 pts) What is the binary representation of the value

528.625

in half precision floating-point? Please show your the calculation in the answer.

(b) (5 pts) What is the rounding error of encoding

528.625

in half precision floating-point?

(c) (10 pts) What is the largest relative error in rounding a half precision floating-point number that belongs to the interval  $[-10000, 10000]$ ? Note that the de-normalized number should be considered, so you need to consider the second and the third cases in (1).

*Solution.*

(a)

$$\begin{aligned} & 528.625 \\ &= 512 + 16 + 0.5 + 0.125 \\ &= 2^9 + 2^4 + 2^{-1} + 2^{-3} \\ &= (10\ 0001\ 0000.101)_2 \\ &= 2^9 \times (1.0000\ 1000\ 0101)_2 \end{aligned} \quad (2)$$

Therefore, the significand part is rounding to

$$0000\ 1000\ 01 \tag{3}$$

and the exponent part is

$$9 + 15 = 24 = 11000_2.$$

The binary representation of 528.625 in half precision floating-point is

$$0\ 11000\ 0000\ 1000\ 01.$$

(b) Let  $v$  be the value of encoding 528.625 in half precision floating-point. By (2) and (3), we can calculate the error

$$\begin{aligned} & 528.625 - v \\ &= 2^9 \times ((1.0000\ 1000\ 0101)_2 - (1.0000\ 1000\ 01)_2) \\ &= 2^9 \times (0.0000\ 0000\ 0001)_2 \\ &= 2^{-3} \times (1.0)_2 \end{aligned}$$

Hence, the rounding error is

$$2^{-3} = 0.125.$$

(c) By the formula of relative error

$$\left| \frac{\text{real value} - \text{rounding value}}{\text{real value}} \right|,$$

let us discuss the relative error in two cases:

Case 1:  $11111_2 > \text{exponent} > 00000_2$ . We have relative error as

$$\begin{aligned} & \left| \frac{2^k \times (1.\text{significand}_2 + \text{rounding error}) - 2^k \times 1.\text{significand}_2}{2^k \times (1.\text{significand}_2 + \text{rounding error})} \right| \\ &= \left| \frac{\text{rounding error}}{(1.\text{significand}_2 + \text{rounding error})} \right|. \end{aligned}$$

Note that

$$\text{rounding error} < 2^{-10}$$

because we have 10 digits for

$$\text{significand}_2.$$

To get the largest relative error in this case, the significand must be as small as possible, which is zero. Therefore, we can use the following optimization problem

$$\max_{x < 2^{-10}} \frac{x}{1+x} \equiv \max_{x < 2^{-10}} 1 - \frac{1}{1+x}$$

to get the largest relative error, and the solution  $x^*$  is

$$2^{-10} \times (0.1111\dots)_2,$$

which implies the largest relative error is bounded by

$$1 - \frac{1}{1+2^{-10}}$$

in this case.

Case 2: exponent = 00000<sub>2</sub>. Similarly, we have relative error as

$$\begin{aligned} & \left| \frac{2^{-14} \times (0.\text{significand}_2 + \text{rounding error}) - 2^{-14} \times 0.\text{significand}_2}{2^{-14} \times (0.\text{significand}_2 + \text{rounding error})} \right| \\ &= \left| \frac{\text{rounding error}}{(0.\text{significand}_2 + \text{rounding error})} \right|. \end{aligned}$$

When significand is equal to zero, the largest relative error is

$$\frac{\text{rounding error}}{(0.0 + \text{rounding error})} = 1$$

in this case.

By combining Case 1 and Case 2, we have the largest relative error as

$$\max\left(1 - \frac{1}{1 + 2^{-10}}, 1\right) = 1.$$

**Problem 2 (10 pts).** Consider the following matrix

$$A = \begin{bmatrix} 9 & 3 & 18 & 27 \\ 3 & 17 & 10 & 21 \\ 18 & 10 & 62 & 72 \\ 27 & 21 & 72 & 100 \end{bmatrix}.$$

Show every step of performing Cholesky factorization on  $A$  with the outer product form. Hint: the resulting  $L$  contains only integer values.

*Solution.*

Step 1: Calculate  $\sqrt{\alpha} = \sqrt{9} = 3$ , so that

$$L^{(1)} = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 6 & 0 & 1 & 0 \\ 9 & 0 & 0 & 1 \end{bmatrix}$$

in this step. Also, we calculate

$$B - \frac{\mathbf{v} \cdot \mathbf{v}^T}{\sqrt{\alpha}} = \begin{bmatrix} 17 & 10 & 21 \\ 10 & 62 & 72 \\ 21 & 72 & 100 \end{bmatrix} - \begin{bmatrix} 1 & 6 & 9 \\ 6 & 36 & 54 \\ 9 & 54 & 81 \end{bmatrix} = \begin{bmatrix} 16 & 4 & 12 \\ 4 & 26 & 18 \\ 12 & 18 & 19 \end{bmatrix}.$$

Thus, we have

$$A^{(1)} = \begin{bmatrix} 3 & 1 & 6 & 9 \\ 1 & 16 & 4 & 12 \\ 6 & 4 & 26 & 18 \\ 9 & 12 & 18 & 19 \end{bmatrix}.$$

Step 2: Now we re-do the calculation on  $A^{(1)}$ , we have  $\sqrt{\alpha^{(1)}} = \sqrt{16} = 4$ , and

$$L^{(2)} = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 1 & 4 & 0 & 0 \\ 6 & 1 & 1 & 0 \\ 9 & 3 & 0 & 1 \end{bmatrix}$$

in this step. Then, calculate

$$B^{(1)} - \frac{\mathbf{v}^{(1)} \cdot (\mathbf{v}^{(1)})^T}{\sqrt{\alpha^{(1)}}} = \begin{bmatrix} 26 & 18 \\ 18 & 19 \end{bmatrix} - \begin{bmatrix} 1 & 3 \\ 3 & 9 \end{bmatrix} = \begin{bmatrix} 25 & 15 \\ 15 & 10 \end{bmatrix}.$$

Therefore,

$$A^{(2)} = \begin{bmatrix} 3 & 1 & 6 & 9 \\ 1 & 4 & 1 & 3 \\ 6 & 1 & 25 & 15 \\ 9 & 3 & 15 & 10 \end{bmatrix}.$$

Step 3: Similarly, we have  $\sqrt{\alpha^{(2)}} = \sqrt{25} = 5$ , and

$$L^{(3)} = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 1 & 4 & 0 & 0 \\ 6 & 1 & 5 & 0 \\ 9 & 3 & 3 & 1 \end{bmatrix}$$

in this step. Furthermore,

$$B^{(2)} - \frac{\mathbf{v}^{(2)} \cdot (\mathbf{v}^{(2)})^T}{\sqrt{\alpha^{(2)}}} = [10] - [9] = [1].$$

Therefore,

$$A^{(3)} = \begin{bmatrix} 3 & 1 & 6 & 9 \\ 1 & 4 & 1 & 3 \\ 6 & 1 & 5 & 3 \\ 9 & 3 & 3 & 1 \end{bmatrix}.$$

Step 4: In the final,  $\sqrt{\alpha^{(3)}} = \sqrt{1} = 1$ , so that

$$L^{(4)} = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 1 & 4 & 0 & 0 \\ 6 & 1 & 5 & 0 \\ 9 & 3 & 3 & 1 \end{bmatrix}$$

is the same as  $L^{(3)}$ . Moreover,  $A^{(4)}$  is also equal to  $A^{(3)}$ .

**Problem 3 (30 pts).** Consider the following matrix:

$$A = \begin{bmatrix} 15 & 18 & 60 \\ 5 & 42 & 95 \\ -35 & -21 & -119 \end{bmatrix}$$

(a) (5 pts) In pivoted LU factorization, please give

$$P_1, M_1, P_2, M_2, U$$

such that

$$M_2 P_2 M_1 P_1 A = U.$$

Note that you should choose the pivot which has the largest absolute value. Hint:  $U$  has only integer values.

(b) (5 pts) Following (a), what are the  $P$  and  $L$  such that  $PA = LU$ ?

(c) (10 pts) In pivoted LU factorization, we have a general method that considers the pivot in both rows and columns

$$PAQ = LU,$$

where  $Q$  is also a permutation matrix that can interchange the columns. To simplify the calculations, the number of row and column interchange is limited to **ONE** when a pivot is determined, i.e., you cannot pick both row interchange and column interchange. For example, the pivot of the  $(1, 1)$  position is the largest absolute value element of the first row and the first column. If it is from the first column, then

$$Q_1 = I,$$

otherwise

$$P_1 = I.$$

Follow the similar idea of (a), please give

$$P_1, Q_1, M_1, P_2, Q_2, M_2$$

and  $U$  so that

$$M_2 P_2 M_1 P_1 \cdot A \cdot Q_1 Q_2 = U.$$

Note that you should choose the pivot which has the largest absolute value. If there are more than two candidates of pivot, we choose the one that is the closest to the diagonal part. Note that  $U$  has fraction numbers, so you want to be careful in doing the calculation.

(d) (5 pts) Following (c), please calculate  $P$ ,  $Q$  and  $L$  such that

$$PAQ = LU.$$

(e) (5 pts) Following (d), solving the linear system

$$Ax = \mathbf{b}$$

is equivalent to solve

$$P^{-1}LUQ^{-1}\mathbf{x} = \mathbf{b}.$$

Therefore, we can get the solution  $\mathbf{x}^*$  by the following steps.

- (i) Solve  $L\mathbf{y} = P\mathbf{b}$  to get  $\mathbf{y}^*$ , where  $\mathbf{y}^* = UQ^{-1}\mathbf{x}$ .
- (ii) Solve  $U\mathbf{z} = \mathbf{y}^*$  to get  $\mathbf{z}^*$ , where  $\mathbf{z}^* = Q^{-1}\mathbf{x}$ .
- (iii) Solve  $Q^{-1}\mathbf{x} = \mathbf{z}^*$  to get  $\mathbf{x}^*$ .

Now, we have

$$\mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

Please solve

$$A\mathbf{x} = \mathbf{b}$$

by the aforementioned steps.

*Solution.*

(a) Step 1: The pivot is  $-35$ , so we have

$$P_1 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

such that

$$P_1 A = \begin{bmatrix} -35 & -21 & -119 \\ 5 & 42 & 95 \\ 15 & 18 & 60 \end{bmatrix}.$$

Then, we can calculate

$$M_1 = \begin{bmatrix} 1 & 0 & 0 \\ 1/7 & 1 & 0 \\ 3/7 & 0 & 1 \end{bmatrix},$$

so that

$$M_1 P_1 A = \begin{bmatrix} -35 & -21 & -119 \\ 0 & 39 & 78 \\ 0 & 9 & 9 \end{bmatrix}.$$

Step 2: The pivot is  $39$ , so

$$P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

and  $P_2 M_1 P_1 A$  is still

$$\begin{bmatrix} -35 & -21 & -119 \\ 0 & 39 & 78 \\ 0 & 9 & 9 \end{bmatrix}.$$

Thus,  $M_2$  is calculated by

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -3/13 & 1 \end{bmatrix},$$

and

$$M_2 P_2 M_1 P_1 A = \begin{bmatrix} -35 & -21 & -119 \\ 0 & 39 & 78 \\ 0 & 0 & -9 \end{bmatrix}.$$

(b) We have

$$P = P_2 P_1 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

and

$$\begin{aligned} L &= M_1^{-1} M_2^{-1} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ -1/7 & 1 & 0 \\ -3/7 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 3/13 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ -1/7 & 1 & 0 \\ -3/7 & 3/13 & 1 \end{bmatrix} \end{aligned}$$

(c) Step 1: The pivot is 60, so

$$P_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

and

$$Q_1 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

such that

$$P_1 \cdot A \cdot Q_1 = \begin{bmatrix} 60 & 18 & 15 \\ 95 & 42 & 5 \\ -119 & -21 & -35 \end{bmatrix}.$$

Thereby, we can calculate

$$M_1 = \begin{bmatrix} 1 & 0 & 0 \\ -95/60 & 1 & 0 \\ 119/60 & 0 & 1 \end{bmatrix}$$

and

$$M_1 P_1 \cdot A \cdot Q_1 = \begin{bmatrix} 60 & 18 & 15 \\ 0 & 27/2 & -75/4 \\ 0 & 147/10 & -21/4 \end{bmatrix}.$$

Step 2: Now the pivot is  $-75/4$ , so

$$P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

and

$$Q_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

such that

$$P_2 M_1 P_1 \cdot A \cdot Q_1 Q_2 = \begin{bmatrix} 60 & 15 & 18 \\ 0 & -75/4 & 27/2 \\ 0 & -21/4 & 147/10 \end{bmatrix}.$$

Hence, we have

$$M_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -21/75 & 1 \end{bmatrix}$$

and then

$$M_2 P_2 M_1 P_1 \cdot A \cdot Q_1 Q_2 = \begin{bmatrix} 60 & 15 & 18 \\ 0 & -75/4 & 27/2 \\ 0 & 0 & 273/25 \end{bmatrix}.$$

(d) We have

$$P = P_2 P_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

$$Q = Q_1 Q_2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

and

$$L = M_1^{-1} M_2^{-1}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 95/60 & 1 & 0 \\ -119/60 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 21/75 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 95/60 & 1 & 0 \\ -119/60 & 21/75 & 1 \end{bmatrix}.$$

(e) In step (i), we have to solve

$$L\mathbf{y} = P\mathbf{b}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 95/60 & 1 & 0 \\ -119/60 & 21/75 & 1 \end{bmatrix} \mathbf{y} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{b}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 95/60 & 1 & 0 \\ -119/60 & 21/75 & 1 \end{bmatrix} \mathbf{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

and the solution  $\mathbf{y}^*$  can be calculated as

$$\mathbf{y}^* = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

In step (ii), we have to solve

$$U\mathbf{z} = \mathbf{y}^*$$

$$\Rightarrow \begin{bmatrix} 60 & 15 & 18 \\ 0 & -75/4 & 27/2 \\ 0 & 0 & 273/25 \end{bmatrix} \mathbf{z} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$



The solution  $\mathbf{z}^*$  can be calculated as

$$\mathbf{z}^* = \begin{bmatrix} -12/273 \\ 18/273 \\ 25/273 \end{bmatrix}.$$

In step (iii), we have to solve

$$Q^{-1}\mathbf{x} = \mathbf{z}^*,$$

which implies

$$\mathbf{x} = Q\mathbf{z}^* = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \mathbf{z}^*.$$

Thus,

$$\mathbf{x}^* = \begin{bmatrix} 18/273 \\ 25/273 \\ -12/273 \end{bmatrix}.$$

**Problem 4 (25 pts).** In our lecture slide “linear\_matrixcondition2.pdf”, we have defined the condition of a matrix  $A$  to be  $\|A\|\|A^{-1}\|$ . In this problem, we explore the relation between the condition number and the singular value decomposition.

Assume that we have a real, invertible,  $n \times n$  matrix  $A$ . We say that the matrix  $A$  has a singular value decomposition if there exists matrices  $U, \Sigma, V$  such that

$$A = U\Sigma V^T$$

where  $U, V$  are  $n \times n$  orthogonal matrices satisfying  $U^T = U^{-1}, V^T = V^{-1}$ ,  
and  $\Sigma$  is an  $n \times n$  diagonal matrix with  $\Sigma_{ii} = \sigma_i > 0$ .

In this problem, let us use  $u_i$  and  $v_i$  to denote the  $i$ th column in  $U$  and  $V$ , respectively.

(a) (5 pts) Use the definition of matrix norm given in our lecture slide “linear\_matrixcondition1.pdf” to show that for  $A$  we have

$$\|A\|_2 = \max_{\|x\|_2=1} \|\Sigma V^T x\|_2.$$

(b) (10 pts) Following subproblem (a), show that

$$\max_{\|x\|_2=1} \|\Sigma V^T x\|_2 = \max_i \sigma_i$$

to conclude that

$$\|A\|_2 = \max_i \sigma_i.$$

Hint: One way to do the proof is by first showing that

$$\|\Sigma V^T x\|_2 \leq \max_i \sigma_i \text{ for all } x \text{ satisfying } \|x\|_2 = 1,$$

and then show that

$$\|\Sigma V^T x\|_2 = \max_i \sigma_i \text{ for some } x \text{ satisfying } \|x\|_2 = 1.$$

(c) (5 pts) Show that for the inverse of  $A$  we also have

$$\|A^{-1}\|_2 = \max_i \frac{1}{\sigma_i}.$$

Then, conclude that the condition number of  $A$  induced by the norm  $\|\cdot\|_2$  is equivalent to

$$\frac{\max_i \sigma_i}{\min_i \sigma_i}.$$

(d) (5 pts) Recall from lecture slide “linear\_matrixcondition2.pdf” that when solving a linear system

$$Ax = b,$$

we have the inequality

$$\frac{\|\delta x\|_2}{\|x\|_2} \leq \|A\|_2 \|A^{-1}\|_2 \frac{\|\delta b\|_2}{\|b\|_2} \quad (4)$$

to quantify the relative change in the solution  $x$  when the input  $b$  is slightly perturbed.

Combining equation (4) with the conclusion in subproblem (c), we can obtain this inequality:

$$\frac{\|\delta x\|_2 / \|x\|_2}{\|\delta b\|_2 / \|b\|_2} \leq \frac{\max_i \sigma_i}{\min_i \sigma_i} \quad (5)$$

In this problem, we discuss the situation when the left hand side of (5) becomes equal to the right hand side (i.e., when  $x$  becomes the most sensitive to perturbation in  $b$ ). Let

$$i_{\max} = \arg \max_i \sigma_i \text{ and } i_{\min} = \arg \min_i \sigma_i.$$

You need to show that

$$\text{if } b = c_1 u_{i_{\max}} \text{ and } \delta b = c_2 u_{i_{\min}} \text{ for some } c_1, c_2 \in \mathbb{R} \\ \text{then the equality in (5) holds.}$$

(This result means that, when  $b$  is parallel to the column vector in  $U$  corresponding to the largest singular value and  $\delta b$  is parallel to that of the smallest singular value, the relative change in  $x$  is the largest.)

*Solution.*

(a) By the definition of matrix norm, we have

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 = \max_{\|x\|_2=1} \|U\Sigma V^T x\|_2. \quad (6)$$

Because  $U$  is orthogonal, we have

$$\|Uw\|_2 = \sqrt{(Uw)^T(Uw)} = \sqrt{w^T U^T U w} = \sqrt{w^T w} = \|w\|_2 \quad (7)$$

for any  $w \in \mathbb{R}^n$ . Therefore, (6) can be simplified to

$$\max_{\|x\|_2=1} \|\Sigma V^T x\|_2.$$

(b) We show that  $\|\Sigma V^T x\|_2$  have an upperbound  $\max_i \sigma_i$  and the upperbound is attained by some vector. To derive the upperbound, let

$$y = V^T x = V^{-1} x.$$

Then, we can see that

$$\begin{aligned} \max_{\|x\|_2=1} \|\Sigma V^T x\|_2 &= \max_{\|Vy\|_2=1} \|\Sigma y\|_2 \\ &= \max_{\|y\|_2=1} \|\Sigma y\|_2 && \text{(for the same reason in (7))} \\ &= \max_{\|y\|_2=1} \sqrt{\sigma_1^2 y_1^2 + \dots + \sigma_n^2 y_n^2} \\ &\leq \max_{\|y\|_2=1} \left[ (\max_i \sigma_i) \sqrt{y_1^2 + \dots + y_n^2} \right] \\ &= \max_{\|y\|_2=1} \left[ (\max_i \sigma_i) \|y\|_2 \right] \\ &= \max_{\|y\|_2=1} \left[ (\max_i \sigma_i) \cdot 1 \right] \\ &= \max_i \sigma_i \end{aligned} \tag{8}$$

Next, we need to show that the upperbound is attained. Let

$$i' = \arg \max_i \sigma_i.$$

Then, the upperbound (8) is attained when

$$x = v_{i'}.$$

This is because

$$\|\Sigma V^T x\|_2 = \|\Sigma e_{i'}\|_2 = \sigma_{i'} = \max_i \sigma_i.$$

In the equation above,  $e_{i'}$  is the vector that has 1 in the  $i'$ th position and zero at other positions.

(c) The singular value decomposition of  $A^{-1}$  can be calculated as

$$\begin{aligned} A^{-1} &= (U \Sigma V^T)^{-1} = (V^T)^{-1} \Sigma^{-1} U^{-1} = V \Sigma^{-1} U^T, \\ &\text{where } \Sigma_{ii}^{-1} = \frac{1}{\sigma_i}. \end{aligned}$$

Then, using the result from subproblem (a) and (b), we can know that

$$\|A^{-1}\|_2 = \max_i \frac{1}{\sigma_i}.$$

Therefore, the condition number of  $A$  is thus

$$\|A\|_2 \|A^{-1}\|_2 = (\max_i \sigma_i) (\max_i \frac{1}{\sigma_i}) = \frac{\max_i \sigma_i}{\min_i \sigma_i}.$$

(d) We can reorganize (5) into this form:

$$\frac{\|\delta x\|_2}{\|\delta b\|_2} \cdot \frac{\|b\|_2}{\|x\|_2} \leq \frac{\max_i \sigma_i}{\min_i \sigma_i} \quad (9)$$

Because  $b = c_1 u_{i_{\max}}$ , we have

$$x = A^{-1}b = V\Sigma^{-1}U^T c_1 u_{i_{\max}} = c_1 V\Sigma^{-1} e_{i_{\max}} = \frac{c_1}{\max_i \sigma_i} V e_{i_{\max}} = \frac{c_1}{\max_i \sigma_i} v_{i_{\max}}.$$

Similarly, since  $\delta b = c_2 u_{i_{\min}}$ , we have

$$\delta x = A^{-1}\delta b = \frac{c_2}{\min_i \sigma_i} v_{i_{\min}}.$$

Plugging all the information we have into (9), we get

$$\frac{\|\delta x\|_2}{\|\delta b\|_2} \cdot \frac{\|b\|_2}{\|x\|_2} = \frac{\frac{c_2}{\min_i \sigma_i}}{c_2} \cdot \frac{c_1}{\frac{c_1}{\max_i \sigma_i}} = \frac{\max_i \sigma_i}{\min_i \sigma_i},$$

so the equality indeed holds.

**Problem 5 (15 pts).** On page 5 of lecture slide “fp\_guarddigit1.pdf”, we showed that subtracting near values of the same sign can lead to large relative errors. Then, in slide “fp\_guarddigit2.pdf”, we showed that with an guard digit, the relative error becomes small. However, the proof only applies to subtracting numbers of the same sign.

In this problem, we discuss the error of adding two floating-point numbers of the same sign when **using  $p$  digits with no guard digit**. As in the lecture slides, we will assume that

$$x > y \geq 0$$

and  $x = x_0.x_1x_2 \cdots x_{p-1} \times \beta^0$  with  $x_0 \geq 1$  and  $\beta \geq 2$ .

We assume that before addition,  $y$  is shifted and then **truncated** to  $p$  digits. We denote the truncated value of  $y$  after shifting as  $\bar{y}$ . Then, after the addition is done, we assume the value is rounded to the closest number (using **rounding even**) before being stored. The rounding error is denoted  $\delta$ . To discuss the final error generated in this process, we separate our discussion into two cases. You need to analyze the error in these two cases using  $y, \bar{y}$  and  $\delta$ , similar to the proof in lecture slide “fp\_guarddigit2.pdf”.

- (a) (10 pts)  $x + y < \beta$ : In this case, we must also have  $x + \bar{y} < \beta$  because truncation can only decrease the result. Since there is no carry out to the  $\beta^1$  bit, no rounding have to performed after addition. The only error arises from the discarded bits when shifting  $y$ . Show that the relative error is less than or equal to  $2\epsilon$ .
- (b) (5 pts)  $x + y \geq \beta$ : In this case, we must also have  $x + \bar{y} \geq \beta$ . This is because the addition in the truncated position does not carry over to the untruncated bits, and thus removing it does not affect the most significant bit. Therefore, there is a carry out to the  $\beta^1$  bit and you also have to consider the rounding error after computing the sum. Show that the relative error is still less than or equal to  $2\epsilon$ .

*Solution.*

(a) Because no rounding have to performed after computation, the error is

$$|(x + y) - (x + \bar{y})| = |y - \bar{y}| = y - \bar{y}.$$

Suppose  $l$  bits is shifted out, then the error is at most

$$y - \bar{y} \leq (\beta - 1)(\beta^{-p} + \beta^{-p-1} + \dots + \beta^{-p-l+1}) < \beta^{-p+1}$$

no matter how large  $l$  is. Then, the relative error is

$$\frac{y - \bar{y}}{x + y} \leq \frac{\beta^{-p+1}}{1} = 2\left(\frac{1}{2}\beta^{-p+1}\right) = 2\epsilon$$

because  $x + y \geq 1$ .

(b) Because there is a carry out, we have to round off the rightmost bit before storing the result. The rounded result can be written as

$$x + \bar{y} + \delta \text{ where } |\delta| \leq \frac{\beta^{-p+2}}{2}.$$

Therefore, the error is

$$|(x + y) - (x + \bar{y} + \delta)| = |(y - \bar{y}) - \delta| \leq |y - \bar{y}| + |\delta| \leq \beta^{-p+1} + \frac{\beta^{-p+2}}{2}.$$

The relative error is then

$$\frac{|(x + y) - (x + \bar{y} + \delta)|}{x + y} \leq \frac{\frac{1}{2}\beta^{-p+2}(1 + \frac{2}{\beta})}{\beta} = \left(1 + \frac{2}{\beta}\right)\frac{1}{2}\beta^{-p+1} \leq 2\epsilon.$$

The first inequality is due to  $x + y \geq \beta$  and the last inequality is due to  $\beta \geq 2$ .