

BLG 202E HOMEWORK #1

- 1) In general floating point system which is consist of β, t, L, U

$$\text{Rounding unit} = \frac{1}{2}\beta^{1-t}$$

For $\beta = 2$ and $t = 52$:

$$\text{Rounding unit} = \frac{1}{2}\beta^{1-t} \rightarrow \frac{1}{2}2^{-51} \rightarrow 2^{-52}$$

- 2) Real answer of $x - y = 0,00016$

- a. With rounding:

$$x = 0,8532 \quad y = 0,8530 \rightarrow x - y = 0,0002$$

$$\text{Absolute Error} = |\text{real} - \text{found}| = |0,00016 - 0,0002| = \mathbf{0,00004}$$

$$\text{Relative Error} = \frac{|\text{real} - \text{found}|}{\text{real}} = \frac{|0,00016 - 0,0002|}{0,00016} = \mathbf{0.25}$$

- b. With chopping:

$$x = 0,8531 \quad y = 0,8530 \rightarrow x - y = 0,0001$$

$$\text{Absolute Error} = |\text{real} - \text{found}| = |0,00016 - 0,0001| = \mathbf{0,00006}$$

$$\text{Relative Error} = \frac{|\text{real} - \text{found}|}{\text{real}} = \frac{|0,00016 - 0,0001|}{0,00016} = \mathbf{0.375}$$

- 3)

- a. When we do matrix multiplication on the left side, we get:

$$ax + by = 1 \quad bx + ay = 0$$

If we add second equation to the first one and do some proper calculations:

$$ax + ay + bx + by = 1$$

$$a(x + y) + b(x + y) = 1$$

$$(x + y)(a + b) = 1$$

$$(x + y) = \frac{\mathbf{1}}{\mathbf{a + b}}$$

- b.

- i. While determining a condition of linear system solving, we should obtain expressions for the variables which are given,

$$\begin{bmatrix} a & b \\ b & a \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \rightarrow ax + by = 1, bx + ay = 0$$

If we solve the equation with two unknowns, we find:

$$x = \frac{a}{a^2 - b^2} \text{ and } y = \frac{b}{b^2 - a^2}$$

We should determine effects of small changes in a and b on x and y. To illustrate,

$$\text{If } a = 1.00001 \text{ and } b = 1 \rightarrow x \approx 50000, y \approx -49999$$

$$\text{If } a = 1.00002 \text{ and } b = 1 \rightarrow x = 0.00004, y \approx -0.00004$$

As we saw, small changes in values causes too big differences. Because of that solving linear system is ill-conditioned.

ii. As we found in part a:

$$(x + y) = \frac{1}{a + b}$$

To reach a result in condition checking, we should change values very little.

$$\text{If } a = 1.00001 \text{ and } b = 1 \rightarrow x + y \approx 0.49999 \approx \frac{1}{2}$$

$$\text{If } a = 1.00002 \text{ and } b = 1 \rightarrow x + y \approx 0.49999 \approx \frac{1}{2}$$

Proximity of values does not affect our results significantly. Because of that solving linear system is well-conditioned.

c. In general floating systems, β, t, L, U determine limits of the systems. With the help of these values, we are able to calculate truth accuracy of our results, biggest and smallest numbers that we can represent and our calculation sensitivity.

d. Standard formula for roots of the quadratic equation is: $\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

$$\text{i. } x_1 = \frac{10^5 + \sqrt{(-10^5)^2 - 4 \cdot 1 \cdot 1}}{2 \cdot 1} \quad x_2 = \frac{10^5 - \sqrt{(10^5)^2 - 4 \cdot 1 \cdot 1}}{2 \cdot 1}$$

While calculating x_2 , value of the square root approximately equal to 10^5 . This means share of the fraction value is calculated as 0 and cancellation error occurs.

Solution: To avoid cancellation error, we can use $x_1 * x_2 = c/a$ formula. Firstly, we should calculate x_1 and using formula we can obtain x_2 .

$$x_1 = \frac{10^5 + \sqrt{(-10^5)^2 - 4 \cdot 1 \cdot 1}}{2 \cdot 1}, \quad x_2 = \frac{c}{a \cdot x_1} = \frac{1}{x_1}$$

ii.

$$x_1 = \frac{-(5 * 10^{30}) + \sqrt{(5 * 10^{30})^2 - 4 * (6 * 10^{30}) * (-4 * 10^{30})}}{2 * (6 * 10^{30})}$$

$$x_2 = \frac{-(5 * 10^{30}) - \sqrt{(5 * 10^{30})^2 - 4 * (6 * 10^{30}) * (-4 * 10^{30})}}{2 * (6 * 10^{30})}$$

While calculating inside of the square root, numbers exceed 10^{60} . But our upper bound is 10^{50} . Therefore, we cannot calculate and represent results of the multiplications.

Solution: To avoid overflow, we should **rescale our equation with any $s \neq 0$** ,

$$c = \sqrt{b^2 - 4ac} \rightarrow c = s \sqrt{\left(\frac{b}{s}\right)^2 - 4 \frac{a}{s} \frac{c}{s}}$$

If we apply this method with $s = 10^{30}$,

$$x_1 = \frac{-(5 * 10^{30}) + 10^{30} * \sqrt{(5)^2 - 4 * (6) * (-4)}}{2 * (6 * 10^{30})}$$

$$x_2 = \frac{-(5 * 10^{30}) - 10^{30} \sqrt{(5)^2 - 4 * (6) * (-4)}}{2 * (6 * 10^{30})}$$