

Numerical analysis Assignment 1

Q1:

a.

Initial number	Times 2	Rest
0.1	0.2	0
0.2	0.4	0
0.4	0.8	0
0.8	1.6	1
0.6	1.2	1
0.2	0.4	0 ...

So:

$$(0.1)_2 = 0.0\overline{0011}$$

$$(\widetilde{0.1})_2 = (0.000110011001100110011001100) = 0.1_{10} - 2^{-20} \cdot 0.1_{10} = 0.1_{10} \cdot (1 - 2^{-20})$$

$$b. \Delta(\widetilde{0.1}) = |\widetilde{0.1} - 0.1| = \left(0.\overbrace{0 \dots 0}^{23} \overline{1100}\right)_2 = 2^{-20} \cdot 0.1_{10}$$

$$\delta(\widetilde{0.1}) = \frac{\Delta(\widetilde{0.1})}{0.1} = \left(\frac{0.\overbrace{0 \dots 0}^{23} \overline{1100}}{0.00011}\right)_2 = \frac{0.0\overline{0011} \cdot 2^{-20}}{0.0\overline{0011} \cdot 2^0} = 2^{-20}$$

$$c. \delta(\widetilde{0.1}) \leq 2^{1-d} \rightarrow 2^{-20} \leq 2^{1-d} \rightarrow -20 \leq 1-d \rightarrow d \leq 21 \rightarrow d = 21$$

$$d. n_1 = 8_h = (8 \cdot 60 \cdot 60 \cdot 10)_{10nth} = (288,000)_{10nth}$$

$$n_2 = 8_h + 2_s = n_1 + (2 \cdot 10)_{10nth} = (288,020)_{10nth}$$

$$\tilde{t}_2 - \tilde{t}_1 = \widetilde{0.1} \cdot n_2 - \widetilde{0.1} \cdot n_1 = \widetilde{0.1} \cdot 20$$

$$t = t_2 - t_1, \tilde{t} = \tilde{t}_2 - \tilde{t}_1 \rightarrow \Delta\tilde{t} = |t - \tilde{t}| = |0.1 \cdot 20 - \widetilde{0.1} \cdot 20| = |20 \cdot \Delta(\widetilde{0.1})| = 20 \cdot 2^{-20} \cdot 0.1 = 2^{-19}$$

$$\delta(\tilde{t}) = \frac{\Delta\tilde{t}}{t} = \frac{2^{-19}}{t_2 - t_1} = \frac{2^{-19}}{0.1 \cdot 20} = 2^{-20}$$

e. Notice that at section d we didn't depend on the absolute number of hours, but on the difference between n_2 and n_1 , which doesn't change in this section. So we get the same answer.

$$f. \tilde{t} = \tilde{t}_2 - \tilde{t}_1 = 0.1 \cdot n_2 - \widetilde{0.1} \cdot n_1 = 0.1 \cdot n_2 - 0.1_{10} \cdot (1 - 2^{-20}) \cdot n_1 = 0.1 \cdot (288,020 - (1 - 2^{-20}) \cdot 288,000) = 2.02746582$$

$$\Delta(\tilde{t}) = |t - \tilde{t}| = |2 - 2.02746582| = 0.02746582$$

$$\delta(\tilde{t}) = \frac{\Delta(\tilde{t})}{t} = \frac{0.027465}{2} = 0.0137329102$$

$$g. \tilde{t} = \tilde{t}_2 - \tilde{t}_1 = 0.1 \cdot n_2 - \widetilde{0.1} \cdot n_1 = 0.1 \cdot n_2 - 0.1_{10} \cdot (1 - 2^{-20}) \cdot n_1 = 0.1 \cdot (3,600,020 - (1 - 2^{-20}) \cdot 3,600,000) = 2.343322754$$

$$\Delta(\tilde{t}) = |t - \tilde{t}| = |2 - 2.343322754| = 0.343322754$$

$$\delta(\tilde{t}) = \frac{\Delta(\tilde{t})}{t} = \frac{0.343322754}{2} = 0.171661377$$

Q2 a. 32bit:

1. The smallest positive normalized number with 32bit in IEEE754 is when we have the lowest normalized exponent possible which is 1, and then by subtracting the bias we will get exponent of -126 and the mantissa will be all zeros.

$$\text{So the number will be: } \left(1.\overbrace{0 \dots 0}^{23}\right)_2 \cdot 2^{-126} = 2^{-126}$$

2. The smallest positive un-normalized number with 32bit in IEEE754 is when we have the exponent field at 0, which indicates the exponent to be -126 and that the number is un-normalized and starts with 0., and the mantissa will be all zeros and 1 at the last bit.

$$\text{So the number will be: } \left(0.\overbrace{0 \dots 0}^{22}1\right)_2 \cdot 2^{-126} = 2^{-149}$$

b.

1. The smallest positive normalized number with 64bit in IEEE754 is when we have the lowest normalized exponent possible which is 1, and then by subtracting the bias we will get exponent of -1022 and the mantissa will be all zeros.

$$\text{So the number will be: } \left(1.\overbrace{0 \dots 0}^{52}\right)_2 \cdot 2^{-1022} = 2^{-1022}$$

2. The smallest positive un-normalized number with 64bit in IEEE754 is when we have the exponent field at 0, which indicates the exponent to be -1022 and that the number is un-normalized and starts with 0., and the mantissa will be all zeros and 1 at the last bit.

$$\text{So the number will be: } \left(0.\overbrace{0 \dots 0}^{51}1\right)_2 \cdot 2^{-1022} = 2^{-1074}$$

Q3

a. $\Delta(\tilde{x} - \tilde{y}) \leq \Delta(\tilde{x}) + \Delta(\tilde{y})$

$$\Delta(\tilde{x} - \tilde{y}) = |(x - y) - (\tilde{x} - \tilde{y})|$$

משולש

$$= |(x - \tilde{x}) + (\tilde{y} - y)| \stackrel{\text{משולש}}{\leq} |x - \tilde{x}| + |\tilde{y} - y|$$

$$= |x - \tilde{x}| + |y - \tilde{y}| = \Delta(\tilde{x}) + \Delta(\tilde{y})$$

b. $\delta\left(\frac{\tilde{x}^2}{\tilde{y}^2}\right) \lesssim? 2(\delta\tilde{x} + \delta\tilde{y})$

We saw in class that for x, y we get $\delta(\tilde{x} \cdot \tilde{y}) \lesssim \delta(\tilde{x}) + \delta(\tilde{y})$.

We saw in practical session that for x, y we get $\delta\left(\frac{\tilde{x}}{\tilde{y}}\right) \lesssim \delta(\tilde{x}) + \delta(\tilde{y})$

We will use both of those inequalities to prove that $\delta\left(\frac{\tilde{x}^2}{\tilde{y}^2}\right) \lesssim? 2(\delta\tilde{x} + \delta\tilde{y})$

$$\delta\left(\frac{\tilde{x}^2}{\tilde{y}^2}\right) = \delta\left(\frac{\tilde{x}}{\tilde{y}} \cdot \frac{\tilde{x}}{\tilde{y}}\right) \lesssim \delta\left(\frac{\tilde{x}}{\tilde{y}}\right) + \delta\left(\frac{\tilde{x}}{\tilde{y}}\right) \lesssim \delta(\tilde{x}) + \delta(\tilde{y}) + \delta(\tilde{x}) + \delta(\tilde{y}) = 2(\delta\tilde{x} + \delta\tilde{y})$$

Q4 $f(x, y) = Z = e^{\alpha(x-y)}$

$$f_x = \alpha \cdot Z, f_y = -\alpha \cdot Z$$

a. $\Delta\tilde{Z} \approx |\nabla f(x, y)| \cdot |(\Delta\tilde{x}, \Delta\tilde{y})| = |\alpha \cdot Z| \cdot \Delta\tilde{x} + |\alpha \cdot Z| \cdot \Delta\tilde{y} = |\alpha \cdot Z| \cdot (\Delta\tilde{x} + \Delta\tilde{y})$

$$\delta(\tilde{Z}) = \frac{\Delta(\tilde{Z})}{|Z|} = \frac{|\alpha \cdot Z| \cdot \Delta\tilde{x} + |\alpha \cdot Z| \cdot \Delta\tilde{y}}{|Z|} = \frac{|\alpha Z| \cdot (\Delta\tilde{x} + \Delta\tilde{y})}{|Z|} = |\alpha| \cdot (\Delta\tilde{x} + \Delta\tilde{y})$$

b. Let $\Delta X = \Delta Y = 2$:

We want $\delta(\tilde{Z}) \leq 0.05$

From a:

$$\delta(\tilde{Z}) = |\alpha| \cdot (\Delta\tilde{x} + \Delta\tilde{y}) = |\alpha|(2 + 2) = 4|\alpha| \leq 0.05 \rightarrow -0.05 \leq 4\alpha \leq 0.05 \\ \rightarrow -\mathbf{0.0125} \leq \alpha \leq \mathbf{0.0125}$$

c. שתי השיטות זהות לחלוטין מבחינה נומרית. אין יתרון ל"סילוק" החיסור מכיוון שאיבוד משמעות היה מתקיים רק אם $|Z| \sim 0$, מצב שלא יתכן בפונקציית האקספוננט. מבחינת כמות החישובים עדיפה השיטה הראשונה הדורשת פחות הפעלות של פונקציית ה $\exp()$.

Q5.

a. $\text{err } 70 = 0.0$

$$\text{err } 7000 = 27.001$$

Omri Bar-oz 313325961

Nevo Strauss 311177638

- שגיאת הצובר גדולה יותר ככל שעובר הזמן מכיוון שאנחנו תומכים רק ב-3 ספרות עשרוניות, אז כאשר מגיעים למצב שיש באקומולטור 1 ואנחנו מנסים להוסיף לו $0.00x$ עבור $1 \leq x \leq 9$ אז לאחר החיתוך המנטיסה לא משתנה, ולכן השגיאה גדלה.

b. $\text{err}72 - \text{err}70 = 5.551115123125783e-17$

$$\text{err}8002 - \text{err}8000 = 0.00800000000000002672$$

- ההפרשים הנ"ל כל כך שונים זה מזה מכיוון שלאחר 70 איטרציות השגיאה היא 0, ולאחר 72 איטרציות השגיאה היא בערך 10^{-17} , כלומר השגיאה מתחילה להצטבר החל מהצעד ה-72.

בתיאוריה, ההפרש בין err_72 ל- err_70 הוא 0, כיוון שבצעדים אלו כמעט ואין שגיאה, ובין err_8002 ל- err_8000 הוא בערך 0.008 כי זה ההפרש האבסולוטי בין המספרים האמיתיים שיחושבו בצעדים ה-8002 וה-8000.

קוד התכנית מצורף בעמוד הבא.

Omri Bar-oz 313325961

Nevo Strauss 311177638

```
import math

# [i] python has a simple typing mechanism.
# ': <type>' signifies the type of the parameter (integer in the following)
def most_significant(num: int, digits_to_keep: int):
    """ return the digits_to_keep most significant digits of num
    :return (the kept digits, the number of non-significant (zeroed- out) digits)
    """
    # [i] 'assert' is a special keyword in python.
    # it verifies if the condition is True
    assert digits_to_keep > 0, 'digits_to_keep should be positive'
    num_digits = math.floor(math.log10(num)) + 1
    non_significant = max(0, num_digits - digits_to_keep)
    # [i] '/' is division without remainder in python 3+
    return (num // 10 ** non_significant), non_significant

def adder(man_a: int, exp_a: int, man_b: int, exp_b: int):
    if exp_a > exp_b:
        return adder(man_b, exp_b, man_a, exp_a)
    exp_diff = exp_b - exp_a
    man_b *= (10 ** exp_diff)
    man_output, exp_output = most_significant(man_a + man_b, 3)
    exp_output += exp_a
    return man_output, exp_output

def accumulate(n: int):
    man_acc = 100
    exp_acc = -5
    man_c = 400
    exp_c = -5
    for i in range(n):
        man_acc, exp_acc = adder(man_a=man_acc, exp_a=exp_acc, man_b=man_c,
exp_b=exp_c)

    real_out = 0.001 + 0.004 * n
    approximate_out = man_acc * (10 ** exp_acc)
    absolute_err = math.fabs(real_out - approximate_out)
    return absolute_err

def main():
    err_70 = accumulate(70)
    err_72 = accumulate(72)
    err_7000 = accumulate(7000)
    err_8000 = accumulate(8000)
    err_8002 = accumulate(8002)
    print(f"err 70 = {err_70}")
    print(f"err 7000 = {err_7000}")
    print(f"err72 - err70 = {err_72 - err_70}")
    print(f"err8002 - err8000 = {err_8002 - err_8000}")

if __name__ == "__main__":
    main()
```