

אנליזה נומרית - עבודה 1

אילנה פרבוי – 318271640, פן אייל - 208722058

שאלה 1

א. נחשב המרה מעשרוני לבינארי עבור 0.1 :

מספר	X2	ערך שלם
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

$$(0.1)_{10} = (0.\overline{00011})_2$$

עבור $0.\tilde{1}$ הייצוג הוא של 23 ספרות ולכן הייצוג יהיה :

$$(0.\tilde{1})_{10} = \left(0.\underbrace{000110011001100110011001100}_{23 \text{ ספרות}} \right)_2$$

ב. חישוב השגיאה המוחלטת :

$$\begin{aligned} \Delta(0.\tilde{1})_{10} &= (|0.1 - 0.\tilde{1}|)_{10} = (|0.\overline{00011} - 0.000110011001100110011001100|)_2 \\ &= \left(0.\underbrace{000000000000000000000000}_{23 \text{ ספרות}} \overline{1100} \right)_2 = (0.1 * 2^{-20})_{10} \end{aligned}$$

חישוב השגיאה היחסית :

$$\delta(0.\tilde{1})_{10} = \left(\frac{\Delta(0.\tilde{1})}{0.1} \right)_{10} = \left(\frac{0.1 * 2^{-20}}{0.1} \right)_{10} = (2^{-20})_{10}$$

ג. נחשב את מספר הספרות הבינאריות המשמעותיות בקירוב $(0.\tilde{1})_{10}$.

נחפש מהו ה d הגדול ביותר כך ש: $\delta(0.\tilde{1})_{10} \leq 2^{1-d}$

$$\delta(0.\tilde{1})_{10} = 2^{-20} \leq 2^{1-d} \leftrightarrow -20 \leq 1 - d \leftrightarrow d \leq 21$$

כלומר $d = 21$.

ד.

$$\begin{aligned} n_1 &= \frac{8_{hours}}{(0.1)_{seconds}} = \frac{(8 * 60 * 60)_{seconds}}{(0.1)_{seconds}} = \frac{(28800)_{seconds}}{(0.1)_{seconds}} = (288000)_{10} \\ n_2 &= \frac{8_{hours} + 2_{seconds}}{(0.1)_{seconds}} = \frac{(8 * 60 * 60)_{seconds} + 2_{seconds}}{(0.1)_{seconds}} = \frac{(28802)_{seconds}}{(0.1)_{seconds}} \\ &= (288020)_{10} \end{aligned}$$

$$\tilde{t}_1 = 0.\tilde{1} * n_1 = 0.\tilde{1} * 288000$$

$$\tilde{t}_2 = 0.\tilde{1} * n_2 = 0.\tilde{1} * 288020$$

$$\tilde{t} = \tilde{t}_2 - \tilde{t}_1 = 0.1 * 288020 - 0.1 * 288000 = 0.1 * (288020 - 288000) = 0.1 * 20$$

$$t = t_2 - t_1 = 0.1 * 288020 - 0.1 * 288000 = 0.1 * (288020 - 288000) = 0.1 * 20$$

$$\begin{aligned}\Delta\tilde{t} &= |\tilde{t} - t| = |0.1 * 20 - 0.1 * 20| = |20 * (0.1 - 0.1)| = 20 * \Delta(0.1) \\ &= 20 * 0.1 * 2^{-20} = 2^{-19}\end{aligned}$$

$$\delta\tilde{t} = \frac{\Delta\tilde{t}}{|t|} = \frac{2^{-19}}{0.1 * 20} = 2^{-20}$$

ה. מכיוון שהיחס $\tilde{t} = \tilde{t}_2 - \tilde{t}_1$ נשמר, אנו נקבל את אותה התוצאה, כפי שניתן לראות:

$$\begin{aligned}n_1 &= \frac{100_{hours}}{(0.1)_{seconds}} = \frac{(100 * 60 * 60)_{seconds}}{(0.1)_{seconds}} = \frac{(360000)_{seconds}}{(0.1)_{seconds}} = (3600000)_{10} \\ n_2 &= \frac{100 + 2_{seconds}}{(0.1)_{seconds}} = \frac{(100 * 60 * 60)_{seconds} + 2_{seconds}}{(0.1)_{seconds}} = \frac{(360002)_{seconds}}{(0.1)_{seconds}} \\ &= (3600020)_{10}\end{aligned}$$

$$\tilde{t}_1 = 0.1 * n_1 = 0.1 * 3600000$$

$$\tilde{t}_2 = 0.1 * n_2 = 0.1 * 3600020$$

$$\tilde{t} = \tilde{t}_2 - \tilde{t}_1 = 0.1 * 3600020 - 0.1 * 3600000 = 0.1 * (3600020 - 3600000) = 0.1 * 20$$

$$t = t_2 - t_1 = 0.1 * 3600020 - 0.1 * 3600000 = 0.1 * (3600020 - 3600000) = 0.1 * 20$$

$$\begin{aligned}\Delta\tilde{t} &= |\tilde{t} - t| = |0.1 * 20 - 0.1 * 20| = |20 * (0.1 - 0.1)| = 20 * \Delta(0.1) \\ &= 20 * 0.1 * 2^{-20} = 2^{-19}\end{aligned}$$

$$\delta\tilde{t} = \frac{\Delta\tilde{t}}{|t|} = \frac{2^{-19}}{0.1 * 20} = 2^{-20}$$

ו. נשים לב כי רק החישוב של \tilde{t}_2 והחישובים התלויים בו ישתנו:

$$n_1 = (288000)_{10}$$

$$n_2 = (288020)_{10}$$

$$\tilde{t}_1 = 0.1 * n_1 = 0.1 * 288000$$

$$\tilde{t}_2 = 0.1 * n_2 = 0.1 * 288020 = 28802$$

נשים לב:

$$(0.1)_{10} = (0.1)_{10} - \Delta(0.1)_{10} = 0.1 - 0.1 * 2^{-20} = 0.1 * (1 - 2^{-20})$$

$$\tilde{t} = \tilde{t}_2 - \tilde{t}_1 = 28802 - 0.1 * 288000 = 28802 - 0.1 * (1 - 2^{-20}) * 288000 = 2.02746582$$

$$t = t_2 - t_1 = 0.1 * 288020 - 0.1 * 288000 = 0.1 * (288020 - 288000) = 0.1 * 20 = 2$$

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

$$\delta \tilde{t} = \frac{\Delta \tilde{t}}{|\tilde{t}|} = \frac{0.02746582}{2} = \mathbf{0.0137329102}$$

ז. בדומה לסעיף ה' נקבל כי n_1, n_2 יהיו שווים ל:

$$n_1 = (3600000)_{10}$$

$$n_2 = (3600020)_{10}$$

$$\tilde{t}_1 = 0. \tilde{1} * n_1 = 0. \tilde{1} * 3600000$$

$$\tilde{t}_2 = 0.1 * n_2 = 0.1 * 3600020 = 360002$$

$$\tilde{t} = \tilde{t}_2 - \tilde{t}_1 = 360002 - 0. \tilde{1} * 3600000 = 360002 - 0.1 * (1 - 2^{-20}) * 3600000 = 2.343322754$$

$$t = t_2 - t_1 = 0.1 * 3600020 - 0.1 * 3600000 = 0.1 * (3600020 - 3600000) = 0.1 * 20 = 2$$

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

$$\delta \tilde{t} = \frac{\Delta \tilde{t}}{|\tilde{t}|} = \frac{0.343322754}{2} = \mathbf{0.171661377}$$

שאלה 2

א. עבור *single* בייצוג 32 ביט:

<i>Normal\Un</i>	<i>Bias</i>	סימן	חזקה	מנטיסה
<i>Normal</i>	$2^7 - 1 = 127$	0	$\underbrace{0 \dots 0}_7 1$ 7 ביטים	$\underbrace{0 \dots 0}_{23}$ 23 ביטים

זהו המספר החיובי הקטן ביותר הניתן להציג בדיוק Single Normal, כלומר:

$$x = +1 * \left(1. \underbrace{0 \dots 0}_{23 \text{ ביטים}} \right)_2 * 2^{-126} = +1 * 1 * 2^{-126} = \mathbf{2^{-126}}$$

<i>Normal\Un</i>	<i>Bias</i>	סימן	חזקה	מנטיסה
<i>Non-Normal</i>	$2^7 - 1 = 127$	0	$\underbrace{0 \dots 0}_7 1$ 7 ביטים	$\underbrace{0 \dots 0}_{22} 1$ 22 ביטים

זהו המספר החיובי הקטן ביותר הניתן להציג בדיוק Single Non-Normal, כלומר:

$$x = +1 * \left(0. \underbrace{0 \dots 0}_{22 \text{ ביטים}} 1 \right)_2 * 2^{-126} = +1 * 2^{-23} * 2^{-126} = +1 * 2^{-149} = \mathbf{2^{-149}}$$

ב. עבור *double* בייצוג 64 ביט:

<i>Normal\Un</i>	<i>Bias</i>	סימן	חזקה	מנטיסה
<i>Normal</i>	$2^{10} - 1 = 1023$	0	$\underbrace{0 \dots 0}_{10 \text{ ביטים}} 1$	$\underbrace{0 \dots 0}_{52 \text{ ביטים}}$

זהו המספר החיובי הקטן ביותר הניתן להציג בדיוק Double Normal, כלומר:

$$x = +1 * \left(1. \underbrace{0 \dots 0}_{52 \text{ ביטים}} \right)_2 * 2^{-1022} = +1 * 1 * 2^{-1022} = 2^{-1022}$$

<i>Normal\Un</i>	<i>Bias</i>	סימן	חזקה	מנטיסה
<i>Non-Normal</i>	$2^{10} - 1 = 1023$	0	$\underbrace{0 \dots 0}_{10 \text{ ביטים}} 1$	$\underbrace{0 \dots 0}_{51 \text{ ביטים}} 1$

זהו המספר החיובי הקטן ביותר הניתן להציג בדיוק Double Non-Normal, כלומר:

$$x = +1 * \left(0. \underbrace{0 \dots 0}_{51 \text{ ביטים}} 1 \right)_2 * 2^{-1022} = +1 * 2^{-52} * 2^{-1022} = +1 * 2^{-1074} = 2^{-1074}$$

שאלה 3

א.

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

נקבל מאי שוויון המשולש כי:

$$\leq |x - \tilde{x}| + |y - \tilde{y}| = \Delta\tilde{x} - \Delta\tilde{y} \blacksquare$$

ב. נראה תחילה כי $\delta(\tilde{x} * \tilde{y}) \lesssim \delta\tilde{x} + \delta\tilde{y}$

נסמן: $\tilde{y} = y - \Delta y$, $\tilde{x} = x - \Delta x$

$$\begin{aligned} \Delta(\tilde{x} * \tilde{y}) &= |\tilde{x} * \tilde{y} - x * y| = |(x - \Delta x) * (y - \Delta y) - x * y| \\ &= |x * y - x\Delta y - y\Delta x + \Delta x\Delta y - x * y| = |-x\Delta y - y\Delta x + \Delta x\Delta y| \end{aligned}$$

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

מאי שוויון המשולש נקבל:

$$\leq \frac{|x\Delta y|}{|x * y|} + \frac{|y\Delta x|}{|x * y|} + \frac{|\Delta x\Delta y|}{|x * y|} = \frac{|\Delta y|}{|y|} + \frac{|\Delta x|}{|x|} + \frac{|\Delta x|}{|x|} * \frac{|\Delta y|}{|y|} = \delta\tilde{x} + \delta\tilde{y} + \delta(\tilde{x} * \tilde{y})$$

נשים לב שכאשר $\delta\tilde{x}$, $\delta\tilde{y}$ קטנים אז $\delta(\tilde{x} * \tilde{y})$ כמעט מתאפס ולכן: $\delta(\tilde{x} * \tilde{y}) \lesssim \delta\tilde{x} + \delta\tilde{y}$.

$$\text{כעת נראה כי } \delta\left(\frac{\tilde{x}^2}{\tilde{y}^2}\right) \lesssim 2(\delta\tilde{x} + \delta\tilde{y})$$

נשים לב: $\delta\left(\tilde{x}^2 * \frac{1}{\tilde{y}^2}\right) \lesssim \delta\tilde{x}^2 + \delta\frac{1}{\tilde{y}^2}$: ומהטענה הקודמת נקבל: $\delta\left(\frac{\tilde{x}^2}{\tilde{y}^2}\right) = \delta\left(\tilde{x}^2 * \frac{1}{\tilde{y}^2}\right)$
מכיוון ש: $\delta(\tilde{x}^2) = \delta(\tilde{x} * \tilde{x}) \lesssim \delta\tilde{x} + \delta\tilde{x}$ ומכיוון ש: $\delta\left(\frac{1}{\tilde{y}^2}\right) = \delta\left(\frac{1}{\tilde{y}} * \frac{1}{\tilde{y}}\right) \lesssim \delta\frac{1}{\tilde{y}} + \delta\frac{1}{\tilde{y}}$
נקבל $\delta\left(\tilde{x}^2 * \frac{1}{\tilde{y}^2}\right) \lesssim \delta\tilde{x} + \delta\tilde{x} + \delta\frac{1}{\tilde{y}} + \delta\frac{1}{\tilde{y}}$

כעת נראה כי $\delta\left(\frac{1}{\tilde{y}}\right) = \delta(\tilde{y})$
לפי הנוסחה:

$$\delta\left(\frac{1}{\tilde{y}}\right) = \frac{\left|\frac{1}{\tilde{y}} - \frac{1}{y}\right|}{\left|\frac{1}{y}\right|} = \frac{\left|\frac{y - \tilde{y}}{\tilde{y} * y}\right|}{\left|\frac{1}{y}\right|} = \left|\frac{y - \tilde{y}}{\tilde{y}}\right| = \left|\frac{\pm\Delta y}{y \mp \Delta y}\right| = \left|\frac{\pm\Delta y}{y \mp \Delta y} * \frac{y \pm \Delta y}{y \pm \Delta y}\right|$$

$$= \left|\frac{\pm\Delta y * y \pm \Delta y^2}{y^2 - \Delta y^2}\right|$$

בהנחה כי השגיאה המוחלטת קטנה, Δy^2 אפסית ולכן:

$$\delta\left(\frac{1}{\tilde{y}}\right) \approx \left|\frac{\pm\Delta y * y}{y^2}\right| = \left|\frac{\pm\Delta y}{y}\right| = \delta(\tilde{y})$$

לבסוף, נקבל כי:

$$\delta\left(\tilde{x}^2 * \frac{1}{\tilde{y}^2}\right) \lesssim \delta\tilde{x} + \delta\tilde{x} + \delta\tilde{y} + \delta\tilde{y} = 2(\delta\tilde{x} + \delta\tilde{y}) \blacksquare$$

שאלה 4

א. נפתח את השגיאה המוחלטת ע"פ כלל הגרדיאנט:

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

$$\Delta\tilde{Z} \approx |\nabla Z(x, y)| * (\Delta\tilde{x}, \Delta\tilde{y}) = \left|\begin{pmatrix} \alpha e^{\alpha(x-y)} \\ -\alpha e^{\alpha(x-y)} \end{pmatrix}\right| * (\Delta x, \Delta y) = \left|\begin{pmatrix} \alpha Z \\ -\alpha Z \end{pmatrix}\right| * (\Delta\tilde{x}, \Delta\tilde{y})$$

$$= |\alpha Z| * \Delta\tilde{x} + |\alpha Z| * \Delta\tilde{y} = |\alpha Z| * (\Delta\tilde{x} + \Delta\tilde{y})$$

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

ב.

$$\Delta\tilde{x} = \Delta\tilde{y} = 2$$

$$\rightarrow \delta\tilde{Z} = |\alpha| * (\Delta\tilde{x} + \Delta\tilde{y}) = 4 * |\alpha|$$

נרצה לבדוק מתי $\delta\tilde{Z} < 5\%$:

$$4 * |\alpha| < 0.05$$

$$\Leftrightarrow |\alpha| < 0.0125$$

$$\Leftrightarrow -0.0125 < \alpha < 0.0125$$

שאלה 5

הקוד פייתון מצורף כקובץ נפרד, נדרש רק להריץ את הפונקציה *main*.
בנוסף מצורפת תמונה של הקוד בעמוד הבא.

א. לאחר הרצת הקוד נקבל:

```
approx for 70 itr: 0.281
```

```
real for 70 itr: 0.281
```

```
absolute error at 70 itr = 0.0
```

```
approx for 7000 itr: 1.0
```

```
real for 7000 itr: 28.001
```

```
absolute error at 7000 itr = 27.001
```

שגיאת הצובר גדולה יותר ככל שעובר הזמן מכיוון שלאחר וסכום הצובר מגיע ל-1 כל פעולת הוספה של מספר הקטן מ: 0.01 ל-1 תגרור שהתוצאה תהיה גם כן 1 ללא שינוי.
בדוגמא שלנו מכיוון ש: $1 = 1 * 10^0$ וש: $10^{-3} = 4 * 0.004$ שנחבר את המספרים, נקבל שההפרש בחזקות הוא מגודל 4, ולכן בעת החישוב של חיבור המנטיסות: $\text{floor}\left(1000 + \frac{4}{10}\right)$ יהיה שווה ל-100 עם חזקה -3 ועל כן נקבל שהתוצאה לאחר החיבור עדיין שווה ל- $100 * 10^{-3} = 1$

ב. לאחר הרצת הקוד נקבל:

```
the error difference between 70 itr and 72 itr = 0.0
```

```
the error difference between 8000 itr and 8002 itr = 0.008
```

שני ההפרשים שונים זה מזה בעקבות ההסבר מסעיף א'.
מכיוון ובאיטרציות 70 ו-72 עדיין לא התחילו להיווצר שגיאות אנו נקבל כי אין שגיאות ביניהן.
אך לאחר 250 איטרציות ($1 = 250 * 0.004$) בעקבות ההסבר בסעיף א', נקבל כי יתחילו להופיע שגיאות, וכפי המצופה השגיאה של האיטרציה ה-8002 תהיה גדולה בפעמיים ערך החיבור המקורי במנטיסה מאשר האיטרציה ה-8000 (מכיוון שהיו 2 יותר חיבורי מנטיסה בגודל 0.004 שלא נלקחו בחשבון)

```

import math

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 (zeroedout) digits) """

    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)
    return (num // (10 ** non_significant)), non_significant

def special_adder(small_mantissa: int, small_exp: int, large_mantissa: int, large_exp: int):
    if small_exp > large_exp:
        # calculate c + acc
        # better to add larger (= acc) to smaller (= c)
        special_adder(large_mantissa, large_exp, small_mantissa, small_exp)

    # bring to the same exp
    diff_exp = large_exp - small_exp
    large_mantissa *= (10 ** diff_exp)
    new_mantissa, new_exp = most_significant(small_mantissa + large_mantissa, 3)
    new_exp += small_exp
    return new_mantissa, new_exp

def accumulator(n: int):
    acc_mantissa = 1
    acc_exp = -3
    c_mantissa = 4
    c_exp = -3
    for i in range(n):
        acc_mantissa, acc_exp = special_adder(acc_mantissa, acc_exp, c_mantissa, c_exp)

    return acc_mantissa * (10 ** acc_exp)

def print_num_and_error(n: int):
    acc = accumulator(n)
    print(f"approx for {n} itr: {acc}")
    real = round(0.001 + 0.004 * n, 5)
    print(f"real for {n} itr: {real}")
    err = round(abs(real - acc), 5)
    print(f"absolute error at {n} itr = {err}")
    print("")
    return err

# Press the green button in the gutter to run the script.
if __name__ == '__main__':
    print("-----a-----")
    err_70 = print_num_and_error(70)
    err_7000 = print_num_and_error(7000)

    print("-----b-----")
    err_72 = print_num_and_error(72)
    err_8000 = print_num_and_error(8000)
    err_8002 = print_num_and_error(8002)

    diff_70_72 = round(abs(err_70 - err_72), 5)
    print(f"the error difference between 70 itr and 72 itr = {diff_70_72}")
    diff_8000_8002 = round(abs(err_8000 - err_8002), 5)
    print(f"the error difference between 8000 itr and 8002 itr = {diff_8000_8002}")

```