

שאלה 1
 א.

סימן	חזקה	שבר
1	10001000	00100010101100110011010

ב. מספר המספרים השלמים הניתנים לייצוג בעזרת IEEE754:
 מתוך כל האפשרויות עבור המנטיסה, בחצי מהספרות ה-1 הכי ימני יהיה הביט האחרון, עבור כל סיפרה כזו נקבל מספר שלם אם נכפיל אותו ב-2 בחזקת X כך X מספר שלם בין 23 ל-127 כולל.
 נחזור על הפעולה עבור כל מספר שה-1 הכי ימני בו הוא הביט האחד לפני האחרון (ברבע מהאפשרויות) כעת ניתן לקבל מספר שלם אם נכפיל ב-2 בחזקת X כך ש-X מספר שלם בין 22 ל-127 כולל.
 נחזור על הפעולה עד שה-1 הכי ימני יהיה הביט הראשון (מקרה בודד) אותו ניתן להכיל ב-2 בחזקת X כך ש-X מספר שלם בין 2 ל-127 כולל.
 נכפול ב-2 עבור כל אפשרות של סיבית הסימן.

$$2 \cdot \sum_{i=1}^{23} 2^{23-i} \cdot (127 - 23 + (i + 1))$$

ג. 6 ספרות בהצגה דסימלית. דבר זה נובע מהפרשים בין חזקות של 2 שניתן ליצג - הרי ההפרש הכי גדול הוא $2^{-23} = \frac{1}{8388604}$. כאשר עוברים את החזקה ה-23, נוצר פער בין מספרים שניתן ליצג, לכן הדיוק יורד ככל שהמספר עולה.

ד.

- i) $3.4028 \cdot 10^{38}$
- ii) $1.18 \cdot 10^{-38}$

ה) 2^{23}

שאלה 2

א.

```
function [est_pi, error, d] = calcPi(n)
    pi_estimate = 2^n;
    last = 0;
    while (n > 1)
        last = sqrt(vpa(2 + last));
        n = n - 1;
    end
    last = sqrt(2-last);
    est_pi = pi_estimate * last;
    error = vpa(pi) - vpa(est_pi);
    d = ceil(log10(error));
    return
```

ב.

n	d
4	2
8	4
12	7
16	9
20	11
24	14

ג. הסטודנט לא צדק, כי מתקיים $\lim_{n \rightarrow \infty} \sqrt{2 + \sqrt{2 + \sqrt{2 + \dots}}} = 2$. כלומר לח גדול מספיק, הערך של

$\lim_{n \rightarrow \infty} \sqrt{2 + \sqrt{2 + \sqrt{2 + \dots}}} - 2$ הוא קטן מאוד ולא מחושב נכון - המחשב יגיע לתוצאה 0 שתוביל

לשגיאה בחישוב.

ד. נוסף דיוק:

```
function [est_pi1, error1, d1] = calcPi1(n)
    digits(1000)
    format LONG
    [est_pi1, error1, d1] = calcPi(n);
    return
```

ה.

n	d
4	2
8	4
12	7
16	9
20	11
24	14

ו. הקוד להצגת הגרפים:

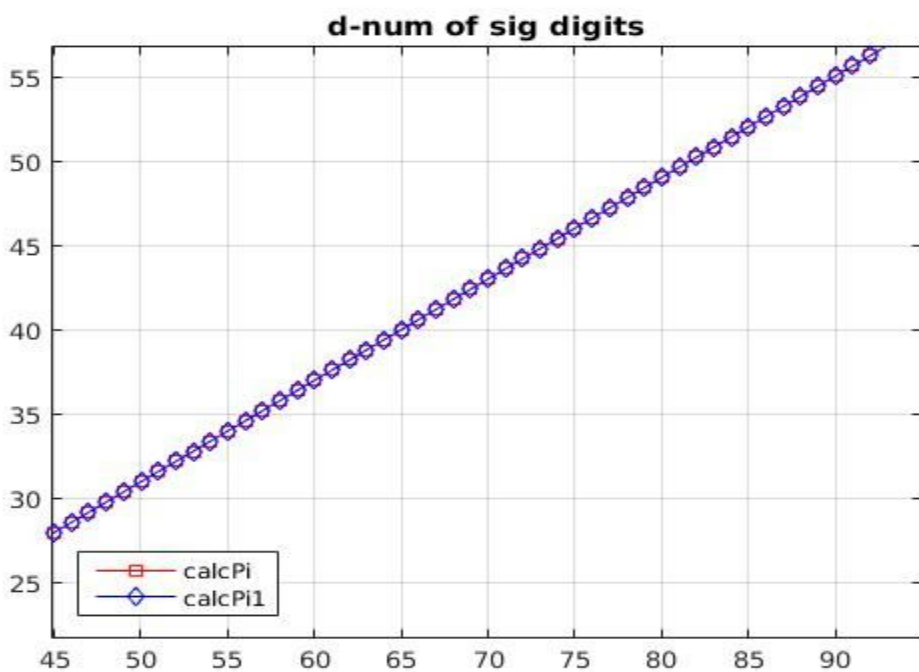
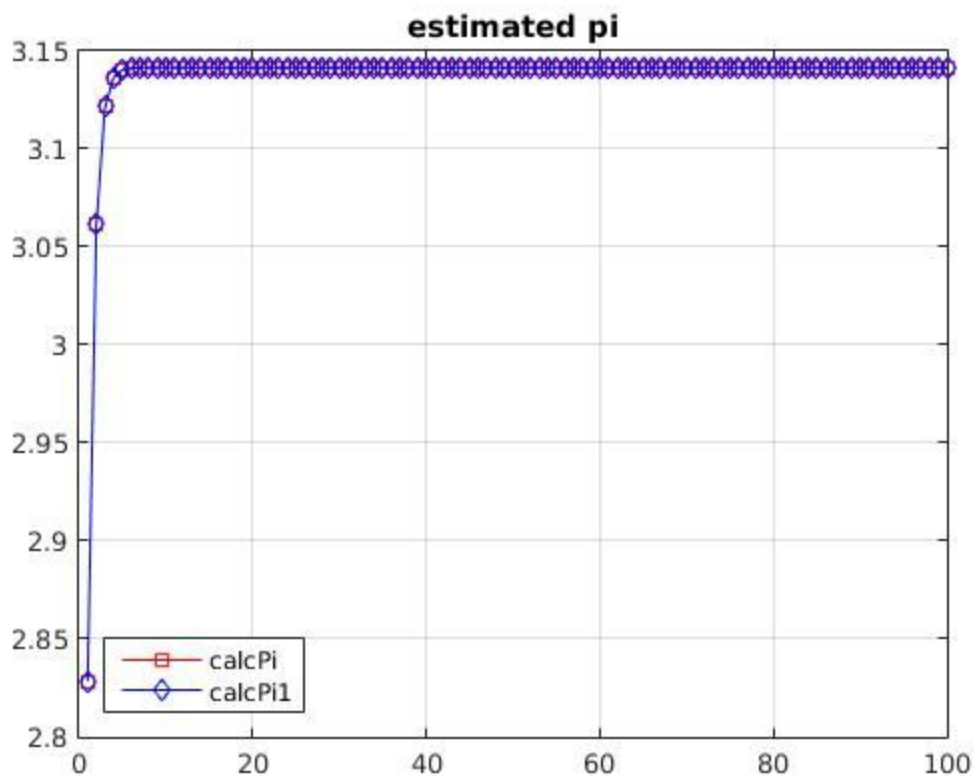
```
t=100;
y1=zeros(1,t);
y2=zeros(1,t);
for n=1:t
    [est_pi,err,d]=calcPi(n);
    e1(1,n)=err;
    d1(1,n)=d;
    p1(1,n)=est_pi;
    [est_pi,err,d]=calcPi1(n);
    e2(1,n)=err;
    d2(1,n)=d;
    p2(1,n)=est_pi;
end

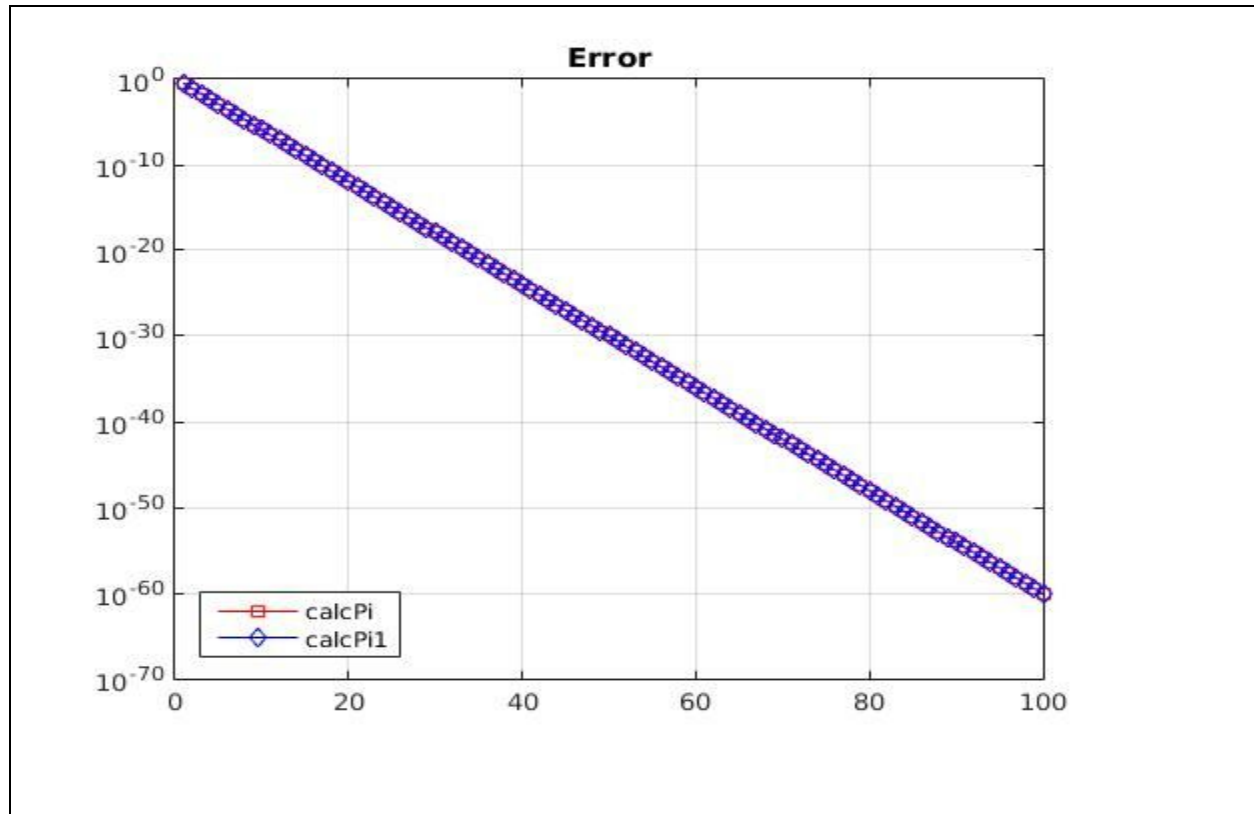
x=1:t;
figure(1)
semilogy(x,e1,'rs-');
hold on

semilogy(x,e2,'bd-');
grid on
hold off
figure(2)
plot(x,d1,'rs-');
hold on
plot(x,d2,'bd-');
grid on
hold off

figure(3)
plot(x,p1,'rs-');
hold on
```

```
plot(x,p2,'bd-');  
hold off
```





3. נסמן:

$$x = y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e$$

$$\tilde{x} = x_1 \cdot x_2 \cdots x_n \cdot \beta^e$$

נשים לב כי:

$$\beta - \text{Base}, n \leq k, m \leq e \leq M$$

אם כן, לפי משפט שלמדנו, מספר הספרות המשמעותיות (בחיתוך) של \tilde{x} הוא d אם מתקיים ש- d הוא המספר השלם האי-שלילי הקטן ביותר המקיים:

$$\delta \tilde{x} = |x - \tilde{x}| / |x| \leq \beta^{1-d}$$

נרצה להראות: $d=n$:

$$\delta \tilde{x} = |x - \tilde{x}| / |x| = \frac{|y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e - x_1 \cdot x_2 \cdots x_n \cdot \beta^e|}{|y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e|}$$

מכיוון ש- \tilde{x} הוא קירוב בחיתוך של x ל- n ספרות, מתקיים $1 \leq x_i = y_i$ אם כן:

$$\frac{|y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e - x_1 \cdot x_2 \cdots x_n \cdot \beta^e|}{|y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e|} = \frac{|0.0 \cdots 0 y_{n+1} \cdots y_k \cdot \beta^e|}{|y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e|} = \frac{|y_{n+1} \cdot y_{n+2} \cdots y_k \cdot \beta^{e-n}|}{|y_1 \cdot y_2 \cdots y_n \cdots y_k \cdot \beta^e|} \leq \beta^{1-n}$$

מכיוון ש- d הוא המספר השלם האי-שלילי הגדול ביותר המקיים זאת, ניתן להניח ש- $d \geq n$. נניח בשלילה כי $d > n$: לדוגמה, $d = n + 1$, ונסתכל על המקרה הבא:

$$\beta = 10, n = 4, x = 1.4321567 \cdot 10^6, \tilde{x} = 1.432 \cdot 10^6$$

אם כן:

$$\delta\tilde{x} = \frac{|1.4321567 \cdot 10^6 - 1.432 \cdot 10^6|}{|1.4321567 \cdot 10^6|} = \frac{|1.567 \cdot 10^2|}{|1.4321567 \cdot 10^6|} \geq 1 \cdot 10^{-4} = 10^{1-(n+1)}$$

בסתירה להנחה.

.4

א.

```
function sum = anss(array)
    sum = 0;
    for ii = 1:numel(array)
        sum = sum + array(ii);
    end
```

נוכיח כי פונקציה זאת תעבוד לדיוק יחסי של 10^{-7} במקרה שנקבל מערך בגודל $2^n, 2 \leq n \leq 20$ אשר איבריו הם $1 \leq c \leq 1000$ ויצוגם כ-Single הוא \tilde{c} :

נזכור כי לפי שאלה קודמת, Single מחזיק ב-6 ספרות מדויקות לפחות. מכיוון שהמספר חסום על ידי 1000. כלומר הדיוק לכל מספר הוא לכל הפחות 10^{-2} . נסתכל על המקרה שבו x מקסימלי קטן ככל האפשר עם שגיאה כדי להגדיל את השגיאה היחסית. מקרה קרוב הוא:

$$x = 1 + 10^{-8} \quad \tilde{x} = 1$$

באופן טריוויאלי, השגיאה היחסית תהיה קטנה מ- 10^{-7} , לכן נסתכל על הקצה השני (במקרה הגרוע).

כאן החזקה הבסיסית היא $2^9 = 512$, לכן הדיוק הוא עד 10^{-4} למעשה. אם כן

$$x = 999 + 9 \cdot 10^{-5} \quad \tilde{x} = 1000$$

אם כן השגיאה היחסית:

$$\delta\tilde{x} = \frac{|x - \tilde{x}|}{|x|} \leq \frac{|2^n \cdot 10^{-4}|}{|2^n \cdot 999.00009|} = 10^{-4} / (9.9900009 \cdot 10^3) \leq 10^{-7}$$

ב. חסם מקורב על השגיאה בסכום של 10^6 ערכים בין 1 ל-2:

לערכים בין 1 ל-2 במשתני single יש שגיאה של לכל היותר $\frac{1}{8688608} = \frac{1}{2^{23}}$. אם כן למיליון ערכים כאלו תהיה

$$\frac{1}{8688608} \cdot \frac{1000000}{8688608} = \frac{10^6}{2^{23}} \sim 0.12$$

והשגיאה היחסית תהיה לכל היותר $\frac{10^6}{2^{23}}$.

ג.

```
function sum = anss(array)
    sum = double(0);
    for ii = 1:numel(array)
        sum = sum + array(ii);
    end
    array = single(array);
    array = double(array);
    for ii = 1:numel(array)
        sum = sum - array(ii);
    end
```

end

5. $f(x) = x^2$, $\tilde{x} = x + \Delta x$, $x, \tilde{x} > 0$.

אם כן:

$$f(\tilde{x}) = f(x + \Delta x) = f(x \pm 1) = (x \pm 1)^2 = x^2 \pm (2x + 1).$$

ידוע כי בכפל השגיאה: $\delta(\tilde{x} \cdot \tilde{y}) = \delta\tilde{x} + \delta\tilde{y}$ מכאן השגיאה בכפל היא:

$$\delta(\tilde{x} \cdot \tilde{x}) = \delta\tilde{x} + \delta\tilde{x} = 2\frac{\Delta x}{x}$$

נשים לב של $\Delta x \geq 1$ החסם לא מתקיים, כאשר בעצם הנחנו ש $\tilde{x} \cdot \tilde{x}$ הוא מספר קטן ביותר (אשר מתכנס 0), אך

לא ניתן להניח זאת כאשר בכפל הוא לא קטן - כלומר כאשר $\Delta x \geq 1$.

והשגיאה לפי גרדיאנט:

$$\delta(\tilde{x}^2) = \frac{(x^2)' \Delta x}{x^2} = 2\frac{\Delta x}{x}$$