

CS-234218

מעבד נתונים

נאדר בשותי

חדר 737

שעות קבלה: ג 10:30 - 11:30

[/http://webcourse.cs.technion.ac.il/234218/](http://webcourse.cs.technion.ac.il/234218/)

ספר הלימוד העיקרי (קיים גם תרגום):

.Cormen, Leiserson, Rivest, Introduction to Algorithms

מטרות הקורס

1. הַפְּרוֹת עִם מְבִנֵי נִתּוּנִים וּמִיִּמּוּשֵׁיהֶם הִיעִילִים
2. פִּיתּוּחַ כְּלִים לְנִיתּוּחַ יַעִילוֹת
3. בְּחִירַת מְבִנֵי נִתּוּנִים לְפִתְרוֹן בַּעִיּוֹת

דוגמאות למבני נתונים: מחסנית, תור, מילון, תור עדיפויות, טבלת ערבול...

דוגמאות לשימושים: מיון, מימוש שפות תכנות, ארגון קבצים, מְנוּעֵי חִיפּוּשׁ
ועוד, ועוד

תוכנית הקורס

1. מבני נתונים בסיסיים וסימונים אסימפטוטיים
2. מערכים ורשימות מקושרות
3. עצים ועצי חיפוש
4. עצי AVL
5. עצי 2-3
6. רשימות דילוגים
7. טבלאות ערבול
8. אחזקת קבוצות זרות
9. מיון
10. טיפול במחרוזות
11. גרפים
12. איסוף אשפה

הגדרות בסיסיות

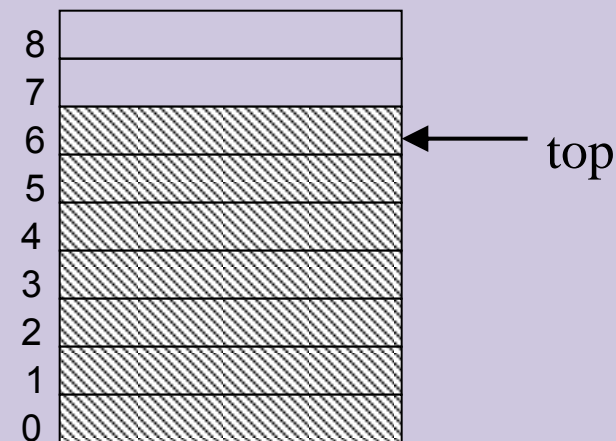
טיפוס נתונים אבסטרקטי (Abstract Data Type = ADT): הוא אוסף של

פעולות ו**כללים** על **קבוצת נתונים**

מימוש של מבנה נתונים הוא אוסף פרוצדורות, אחת לכל פעולה, הממשות את הפעולות של מבנה הנתונים.

דוגמא: מהי מחסנית?

האם היא מערך עם מציין ל- top ?



<code>top = -1</code>	אתחול:
<code>top = top + 1; A[top] = x;</code>	הוספת איבר:
<code>if top > -1 return A[top]</code>	ראש המחסנית:
<code>top = top - 1</code>	הוצאת איבר:

לא ! זהו מימוש של מחסנית באמצעות מערך ומצביע.

מחסנית כטיפוס נתונים אבסטרקטי

אחרון נכנס - ראשון יוצא Last In -- First Out : LIFO

מחסנית (Stack) מוגדרת ע"י הפעולות הבאות:

$create(S)$ - מחזיר מחסנית S ריקה חדשה.

$push(S, x)$ - מכניס איבר בעל ערך x לראש המחסנית S .

$top(S)$ - מחזיר את האיבר שבראש המחסנית S (המחסנית אינה משתנה).

$pop(S)$ - מוציא את האיבר שבראש המחסנית S .

$is-empty(S)$ - מחזיר $true$ אם המחסנית S ריקה ו- $false$ אחרת.

פונקציה: $is-empty:\{Stacks\} \rightarrow \{true, false\}$

מחסנית כטיפוס נתונים אבסטרקטי (המשך)

פעולות המחסנית מקיימות את הכללים הבאים:

1. אפשר לבצע `pop`, `top` רק על מחסנית לא ריקה.
2. מיד לאחר `create(S)`, `is-empty(S)` מחזיר ערך `true`.
3. לאחר ביצוע `push`, ואחריו `pop`, המחסנית לא משתנה.

כלל 2: $is_empty(create(S)) = true$

כלל 3: $pop(push(S, x)) = S$

יש טענות הנובעות רק מהכללים. **לדוגמא:** `create(S);`

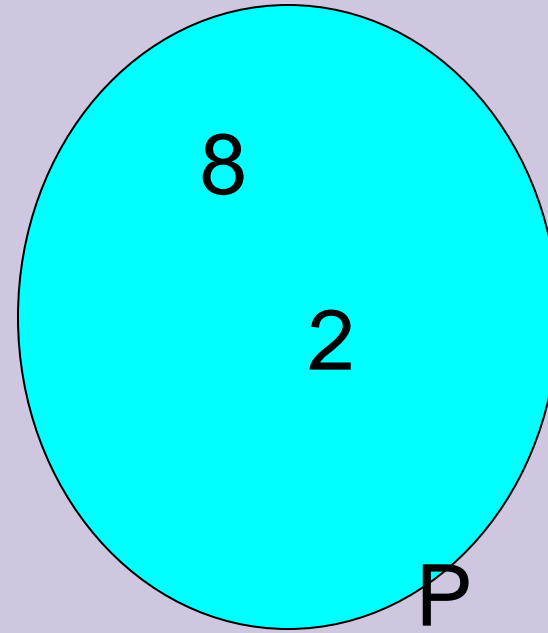
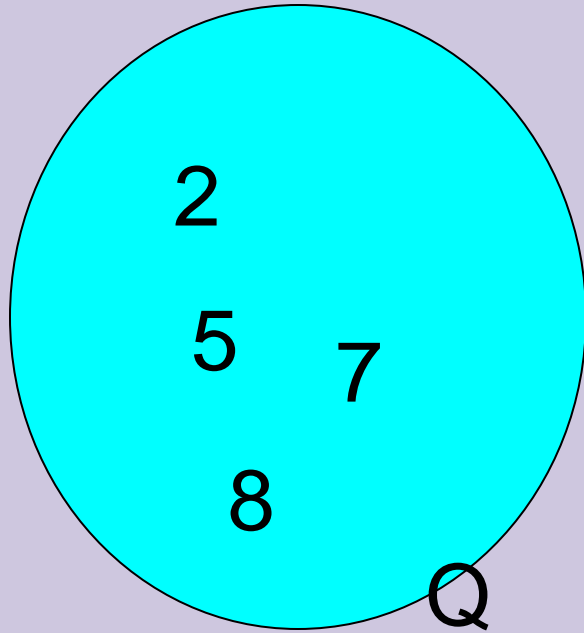
`push(S, 17);`

`pop(S);`

`print is_empty(S);` מה יודפס - הוכח?

`is_empty(pop(push(create(S)))=true;`

כל מימוש חייב לאפשר את הפעולות ולקיים את הכללים.



Create(Q) Push(Q,2) Push(Q,5) Push(Q,7)

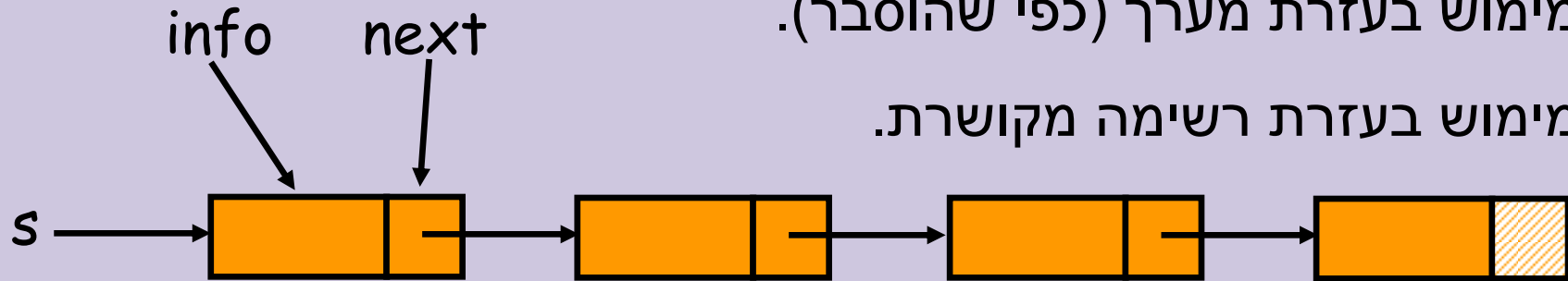
Pop(Q) Create(P) Push(Q,8) Is-Empty(P)

Push(P,8) Is-Empty(Q) Push(P,2) Pop(Q)

מימושי מחסנית

1. מימוש בעזרת מערך (כפי שהוסבר).

2. מימוש בעזרת רשימה מקושרת.



```
#define NULL 0
typedef struct node {
    DATA_TYPE info;
    struct node *next;
} NODE;
typedef NODE *STACK;
```

הגדרת צומת:

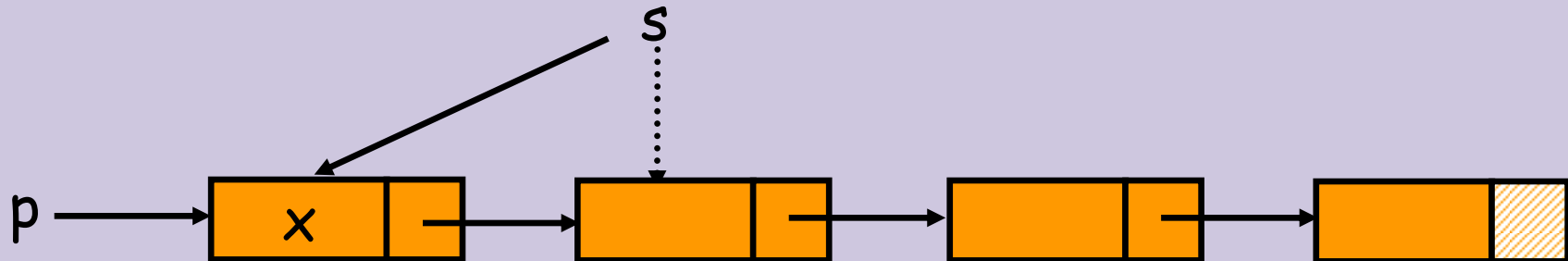
```
void create (STACK *s) {
    (*s) = NULL;
}
```

יצירת מחסנית ריקה:

הכנסת איבר

```
void push ( STACK *s, DATA_TYPE x){  
    NODE *P;  
    P = malloc (sizeof (NODE));  
    P -> info = x;  
    P -> next = (*s) ;  
    (*s) = p ;  
}
```

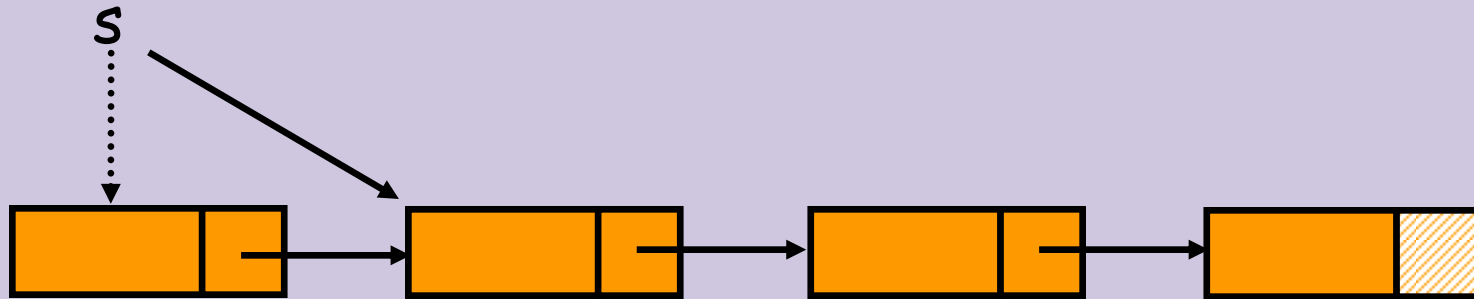
פעולת $\text{push}(s,x)$:



אחזור והוצאת איבר

```
void pop ( STACK *s){
    STACK t ;
    t = (*s) → next;
    free (*s) ;
    (*s) = t ;
}
```

פעולת pop(s):



```
DATA_TYPE top ( STACK *s){
    return (*s) → info ;
}
```

פעולת top(s):

תור כמבנה נתונים

ראשון נכנס - ראשון יוצא First In -- First Out : FIFO

תור (Queue) מוגדר ע"י הפעולות הבאות:

create(Q) - מחזיר תור ריק.

head(Q) - מחזיר את ערך האיבר שבראש התור Q (התור אינו משתנה).

enqueue(Q, x) - מכניס איבר עם ערך x לסוף התור Q.

dequeue(Q) - מוציא את האיבר שבראש התור Q.

is-empty(Q) - מחזיר true אם התור Q ריק ו-false אחרת.

מימוש של תור בעזרת מערך

```

head(Q):      Q[f];
enqueue(Q,x): Q[r] = x;
              r = (r+1) % n;
dequeue(Q)    f = (f+1) % n;
is_empty(Q):  f == r;
create(Q):    f = r = 0;
  
```

מערך Q בן n איברים עם שני מציינים

$r =$ מציין את מקום האיבר שאחרי סוף התור
(rear)

$f =$ מציין את מקום האיבר שבראש התור
(front)

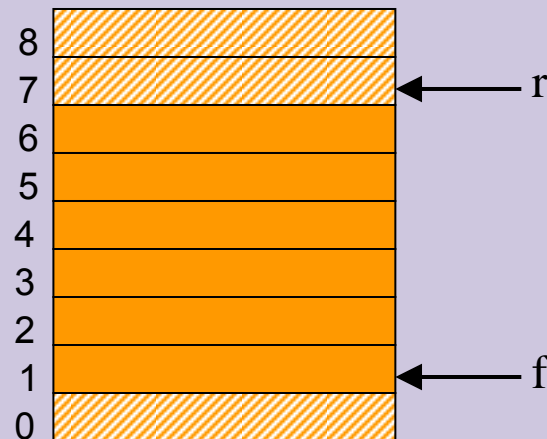
נקודות חשובות:

הפעולות האריתמטיות נעשות $\% n$ (mod n).
(n).

התור מלא אם $f == (r + 1) \% n$.

יש לבדוק "לא מלא" לפני enqueue

ו"לא ריק" לפני dequeue.



n=9

Creat(Q)

7	5	
6	3	← r
5	9	
4	1	
3	2	
2	9	← f
1	5	← r
0	4	← f ← r

enqueue(Q,4)

enqueue(Q,5)

enqueue(Q,9)

enqueue(Q,2)

enqueue(Q,1)

enqueue(Q,9) dequeue(Q)

dequeue(Q)

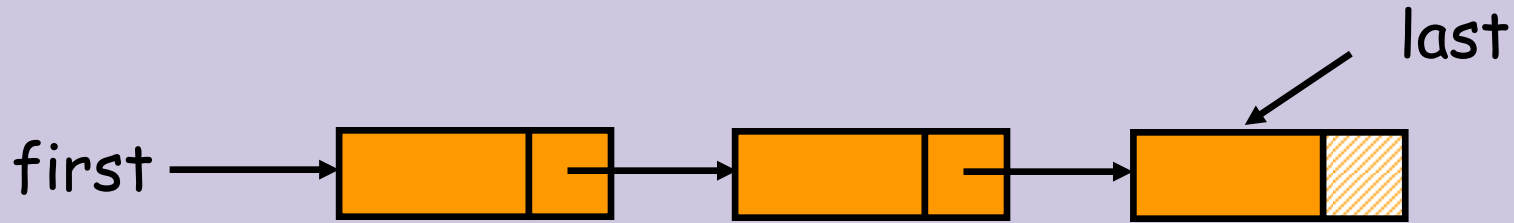
head(Q)?

enqueue(Q,3)

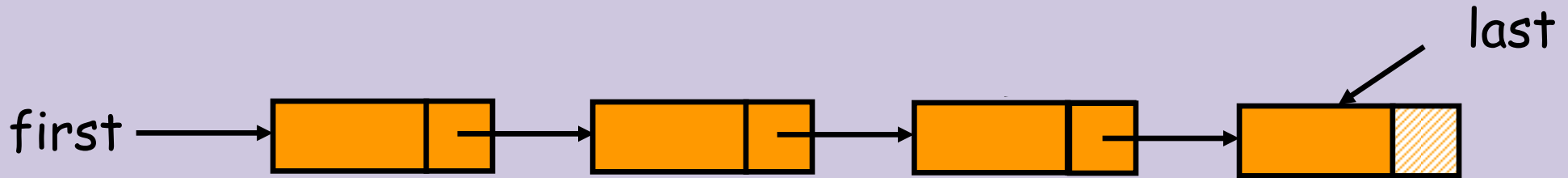
enqueue(Q,5)

enqueue(Q,1)

מימוש תור ע"י רשימה מקושרת



הכנסה בסוף התור (last) :



הוצאה מראש התור (first).

שימו לב שלא ניתן להוציא איבר באמצעות המצביע last.

הערות לגבי ADT ומימושים

Abstract Data Type = ADT

המשתמש ב-ADT:

מכיר את הפעולות והשפעתן על הנתונים.
אינו נדרש להכיר את המימוש.

איכות המימוש מאופינת ע"י:

• יעילות:

• זמן - מספר צעדי החישוב הנדרשים, לכל פעולה.

• מקום - כמות הזיכרון הנדרשת.

• פשטות התוכנית: (מאפשר מימוש נכון ואחזקה יעילה).

זמן ריצה של אלגוריתם

זמן ריצה של אלגוריתם A עבור קלט x יסומן ב- $T_A(x)$.

זמן הריצה נמדד ע"י מספר פקודות מכונה שהאלגוריתם מבצע על קלט נתון.

מדד זה מתעלם מהבדלי המהירות בין הפקודות. (למשל, חבור לעומת כפל).

הגודל של קלט x יסומן ב- $|x|$.

לדוגמא: בתוכנית המסכמת איברי מערך x, גודל הקלט הוא מספר איברי המערך.

זמן הריצה הגרוע ביותר (worst case) של אלגוריתם A

עבור קלט שגודלו n מוגדר ע"י

$$T_A(n) = \max_{|x|=n} T_A(x)$$

דוגמא 1:

נתון מערך $a[]$

```
A: sum = 0
   for (i = 0; i < n; i++)
     sum = sum + a[i];
```

זמן הריצה של אלגוריתם זה עבור כל קלט

$$T_A(n) = c_1 \cdot n + c_2 \text{ הוא } n$$

כאשר c_1, c_2 הם קבועים התלויים במימוש

הפקודות בשפת מכונה.

זמן ריצה של אלגוריתם

דוגמא 2:

```
sum = 0
for (i = 0; i < n; i++)
    if sum < a[i+1]
        sum = sum + a[i];
    else
        terminate;
```

מהו זמן הריצה של האלגוריתם?

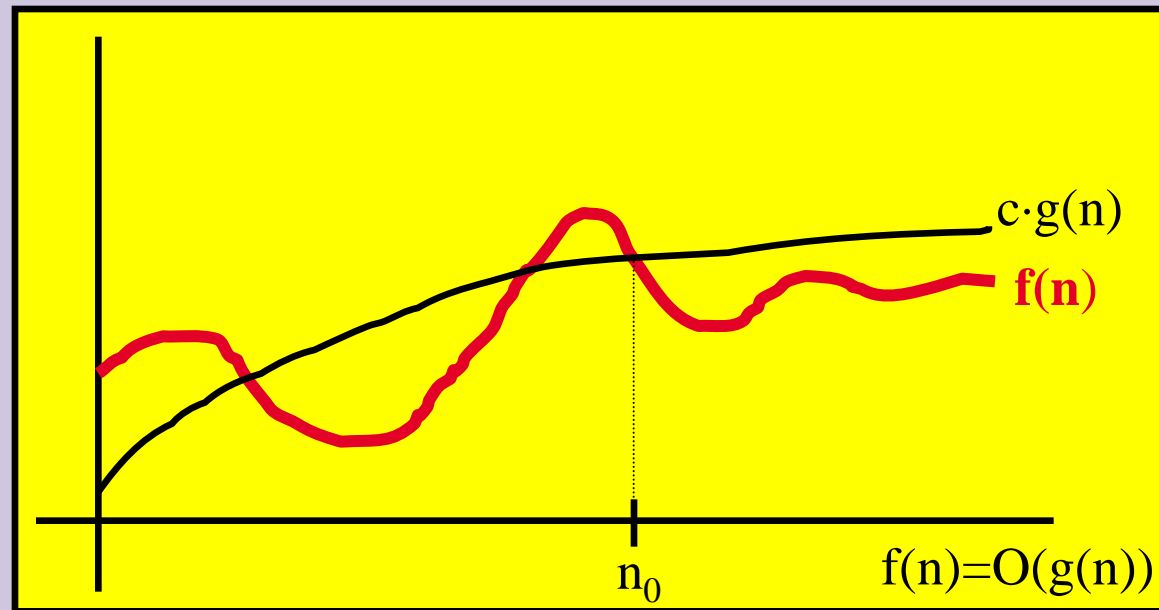
זמן הריצה **הגרוע ביותר** של אלגוריתם זה מקיים: $T_A(n) = c_1 \cdot n + c_2$

סיבוכיות והסימון O

הגדרה: יהיו $f(n), g(n)$ פונקציות חיוביות. נאמר שהפונקציה $f(n)$ נמצאת בקבוצת הפונקציות $O(g(n))$ אם קיימים קבועים c, n_0 כך שלכל $n \geq n_0$ מתקיים:

$$f(n) \leq c \cdot g(n)$$

כמו כן נאמר ש- $g(n)$ מהווה חסם עליון אסימפטוטי לפונקציה $f(n)$ ונסמן זאת ע"י $f(n) = O(g(n))$ במקום הסימון הרגיל $f(n) \in O(g(n))$.



נשתמש בסימון זה כאשר הפונקציה $f(n)$ היא פונקציה שקשה לתאר במדויק, למשל זמן הריצה של אלגוריתם, בעוד $g(n)$ פשוטה יותר לתיאור.

דוגמאות פולינומיאליות

טענה: עבור k קבוע מתקיים $f(n) = a_0 + a_1n + a_2n^2 + \dots + a_kn^k = O(n^k)$

הוכחה: נמצא קבועים c, n_0 , כך שלכל $n \geq n_0$ יתקיים $f(n) \leq c \cdot n^k$.

$$\begin{aligned} f(n) &\leq |a_0| + |a_1|n + |a_2|n^2 + \dots + |a_k|n^k \\ &\leq |a_0|n^k + |a_1|n^k + |a_2|n^k + \dots + |a_k|n^k \\ &= (|a_0| + |a_1| + |a_2| + \dots + |a_k|)n^k = cn^k \end{aligned}$$

$$f(n) = 10006 = O(1)$$

דוגמאות: קבוע

$$f(n) = 4n + 27 = O(n)$$

ליניארי

$$f(n) = n \log_2 n + 3n^2 = O(n^2)$$

ריבועי

הדוגמא האחרונה נובעת מאי השוויון $\log_2 n < n$.

דוגמאות נוספות

$$f(n) = \log_a n = O(\log_b n) = \boxed{O(\log n)}$$

לוגריתמי

מכיוון שמתקיים $\log_a n = \log_a b \cdot \log_b n$.

$$f(n) = \log_a n = O(n^\varepsilon)$$

לכל ε חיובי (שברים ושלמים) מתקיים:

$$f(n) = 8 + 15n + 9n \log_2 n = O(n \log n)$$

בין ליניארי וריבועי

כיוון שמתקיים $f(n) < 32n \log_2 n$.

$$f(n) = n^k + a^n = O(a^n) \quad (a > 1)$$

אקספוננציאלי

$$f(n) = a^n = O(b^n) \quad (b \geq a)$$

דוגמאות

$$f(n) = 3^n = O(2^n)$$

האם מתקיים ?

$$f(n) = 2^n = O(n^k) \quad (k \text{ קבוע})$$

$$f(n) = \log n = O(\log \log n)$$

התשובה שלילית בשלושת המקרים. הוכיחו!

סיבוכיות זמן

דוגמא ראשונה: סכום איברי מערך. ראינו שמתקיים $T(n) < c \cdot n + 1$

ולפיכך $T(n) = O(n)$.

דוגמא שניה: כפל מטריצות ריבועיות בגודל $m \cdot m$.

גודל הקלט $n = 2m^2$.

$$\begin{array}{c} m \\ \square \\ m \end{array} A \times \begin{array}{c} m \\ \square \\ m \end{array} B = \begin{array}{c} m \\ \square \\ m \end{array} C$$

מחשבים m^2 איברים במטריצת התוצאה C . כאשר כל איבר מחושב ע"י:


$$C[i, j] = \sum_{k=1}^m A[i, k] \cdot B[k, j]$$

סיבוכיות הזמן כפונקציה של m היא $O(m^3)$.

סיבוכיות הזמן כפונקציה של גודל הקלט n היא $O(n^{3/2})$ כיוון שמתקיים:

$$T(n) = O(m^3) = O\left(\frac{n^{3/2}}{2^{3/2}}\right) = O(n^{3/2})$$

$m = \sqrt{n/2}$



קבוע

סיבוכיות זמן (המשך)

דוגמא שלישית: חיפוש בינרי של x במערך ממוין בן n איברים.

בכל צעד משווים את x לאיבר האמצעי במערך העֵכְשׁוּי.

• אם האיבר האמצעי שווה ל- x החיפוש נגמר.

• אם האיבר האמצעי גדול מ- x ממשיכים עם חלק המערך המכיל את המספרים הקטנים.

• אם האיבר האמצעי קטן מ- x ממשיכים עם חלק המערך המכיל את המספרים הגדולים.

נסמן ב- T את סיבוכיות הזמן. מתקיימת משוואת הנסיגה הבאה:

$$T(1) = c_1 \quad T(n) = c + T(\lceil n/2 \rceil)$$

בהנחה ש- n חזקה של 2 נקבל:

$$\begin{aligned} T(n) &= c + c + T(n/4) = \overbrace{c + c + \dots + c}^i + T(n/2^i) \\ &= c + c + \dots + c + T(n/2^{\log_2 n}) = c \log_2 n + T(1) = O(\log n) \end{aligned}$$

הטענה נכונה גם אם n לא חזקה של 2. פרטים על שיטות פתרון למשוואות נסיגה ניתן למצוא

בספר הלימוד (פרק 4): Cormen, Leiserson, Rivest, Introduction to Algorithms.

סיבוכיות זמן (המשך)

דוגמא רביעית:

```
S = 0;
for ( i = 1; i < n; i ++ )
    for ( j = 0; j < n; j + = i )
        S++ ;
```

אנליזה גסה: שתי לולאות מקוננות. $T(n) \leq n(n-1)$.

אנליזה עדינה:

מתבצעות n פעולות.

$$\lceil n/2 \rceil$$

$$\lceil n/3 \rceil$$

$$n/n = 1$$

כאשר $i = 1$,

$$i = 2$$

$$i = 3$$

...

$$i = n$$

H_n נקרא המספר ההרמוני ה- n

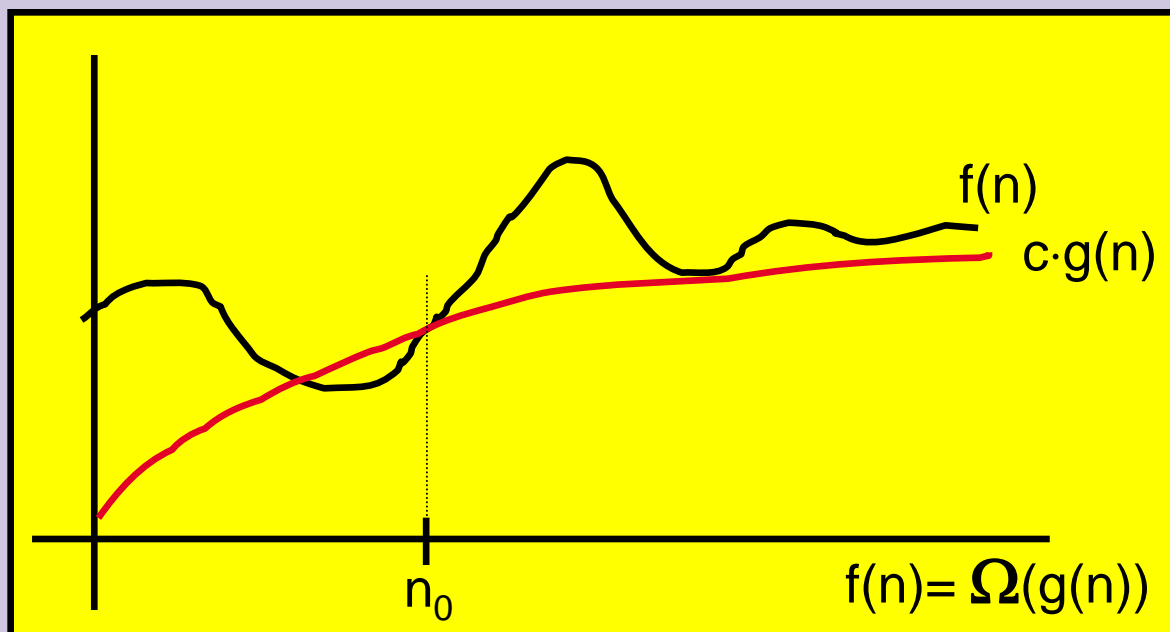
$$T(n) \leq n + \frac{n}{2} + \frac{n}{3} + \dots + \frac{n}{n} = n \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \right) = nH_n \quad \text{סה"כ}$$

$$\ln(n+1) = \int_1^{n+1} \frac{dx}{x} \leq 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \leq 1 + \int_1^n \frac{dx}{x} = 1 + \ln(n) \quad \text{כיוון שמתקיים:}$$

הערה: לשם דיוק היה עלינו לכפול את כל המשוואות בקבוע ולפיכך: $T(n) = O(n \cdot \log n)$

חסם תחתון אסימפטוטי

הגדרה: יהיו $f(n)$, $g(n)$ פונקציות חיוביות. נאמר שהפונקציה $f(n)$ נמצאת בקבוצת הפונקציות $\Omega(g(n))$ (אומגה) אם קיימים קבועים c, n_0 כך שלכל $n \geq n_0$ מתקיים:

$$f(n) \geq c \cdot g(n)$$


דוגמאות פולינומיאליות

טענה: אם $a_k > 0$ עבור k קבוע מתקיים

$$f(n) = a_0 + a_1n + a_2n^2 + \dots + a_kn^k = \Omega(n^k)$$

הוכחה: נמצא קבועים c , n_0 כך שלכל $n \leq n_0$ יתקיים $c \cdot n^k \leq f(n)$.

$$\begin{aligned} f(n) &\geq -|a_0| - |a_1|n - |a_2|n^2 - \dots - |a_{k-1}|n^{k-1} + a_kn^k \\ &\geq -|a_0|n^{k-1} - |a_1|n^{k-1} - |a_2|n^{k-1} - \dots - |a_{k-1}|n^{k-1} + a_kn^k \\ &= (-(|a_0| + |a_1| + |a_2| + \dots + |a_{k-1}|) + a_kn)n^{k-1} \end{aligned}$$

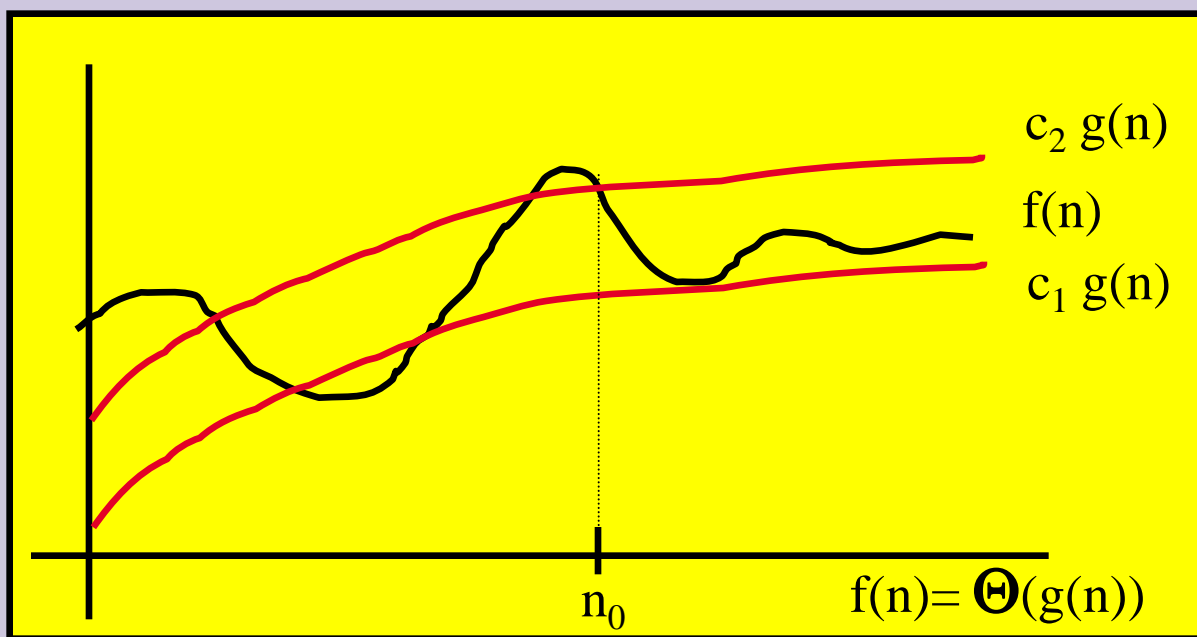
$$n \geq n_0 = \frac{2(|a_0| + |a_1| + |a_2| + \dots + |a_{k-1}|)}{a_k} \quad \text{עבור}$$

$$\geq \frac{a_k n}{2} n^{k-1} = cn^k$$

$$c = \frac{a_k}{2}$$

חסם הדוק אסימפטוטי

הגדרה: יהיו $f(n)$, $g(n)$ פונקציות חיוביות. נאמר שמתקיים $f(n) = \Theta(g(n))$ (תטה) אם $f(n) = O(g(n))$ וגם $f(n) = \Omega(g(n))$.



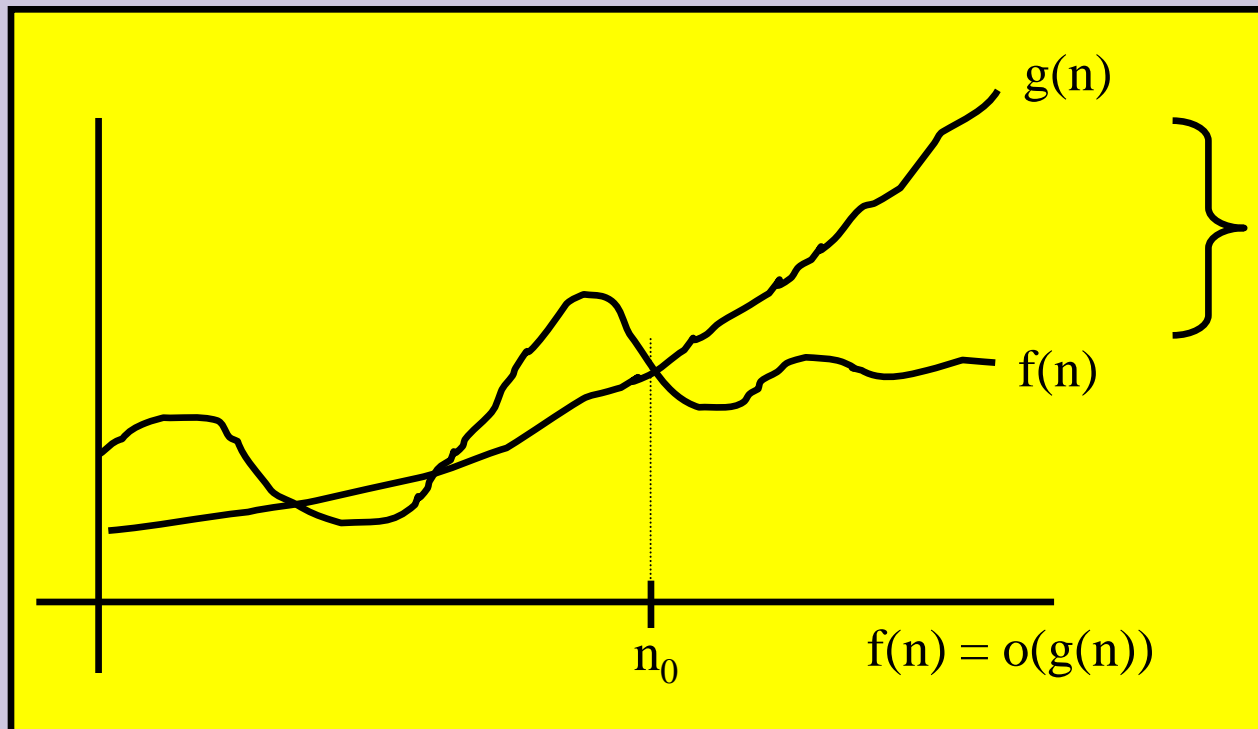
הגדרה (אקויוולנטית): נאמר שמתקיים $f(n) = \Theta(g(n))$ אם קיימים קבועים c_1, c_2, n_0 כך שלכל $n \geq n_0$ מתקיים $c_1 \cdot g(n) \leq f(n) \leq c_2 \cdot g(n)$.

$$f(n) = a_0 + a_1 n + a_2 n^2 + \dots + a_k n^k = \Theta(n^k) \quad \text{דוגמא}$$

הסימון o קטן

הגדרה: יהיו $f(n), g(n)$ פונקציות חיוביות. נאמר שהפונקציה $f(n)$ נמצאת בקבוצת הפונקציות $(g(n))$ אם לכל קבוע c קיים קבוע n_0 כך שלכל $n \geq n_0$ מתקיים:

$$f(n) \leq c \cdot g(n)$$



$f(n)$ זניחה אסימפטוטית
יחסית לפונקציה $g(n)$.

$$\lim_{n \rightarrow \infty} f(n)/g(n) = 0$$

הגדרה (אקויוולנטית): נאמר שמתקיים $f(n) = o(g(n))$ אם

$$\log n = o(n), \quad n-100 \neq o(n) \quad \text{דוגמאות:}$$

מגבלות הסימון האסימפטוטי

נוח להשתמש בסימונים אסימפטוטיים מפני

1. הסימון מתעלם מקבועים

2. מאפשר ניתוח זמנים פשוט יותר.

אנו נשתמש בסימונים אלה לאורך הקורס.

אבל לסימון יש מגבלות מסוכנות...

ברור שנעדיף תוכנית הרצה בזמן $T(n) = n^2$

על פני תוכנית הרצה בזמן קבוע של $T(n) = 10^{80}$

כיון שבתוכניות ממשיות אנו משתמשים בגודל קלט n סופי הקטן באופן משמעותי מ-
 10^{40} .

מסקנה: צריך לוודא שהקבועים c, n_0 המתחבאים בהגדרות אסימפטוטיות O, Θ, Ω
אמנם "סבירים".

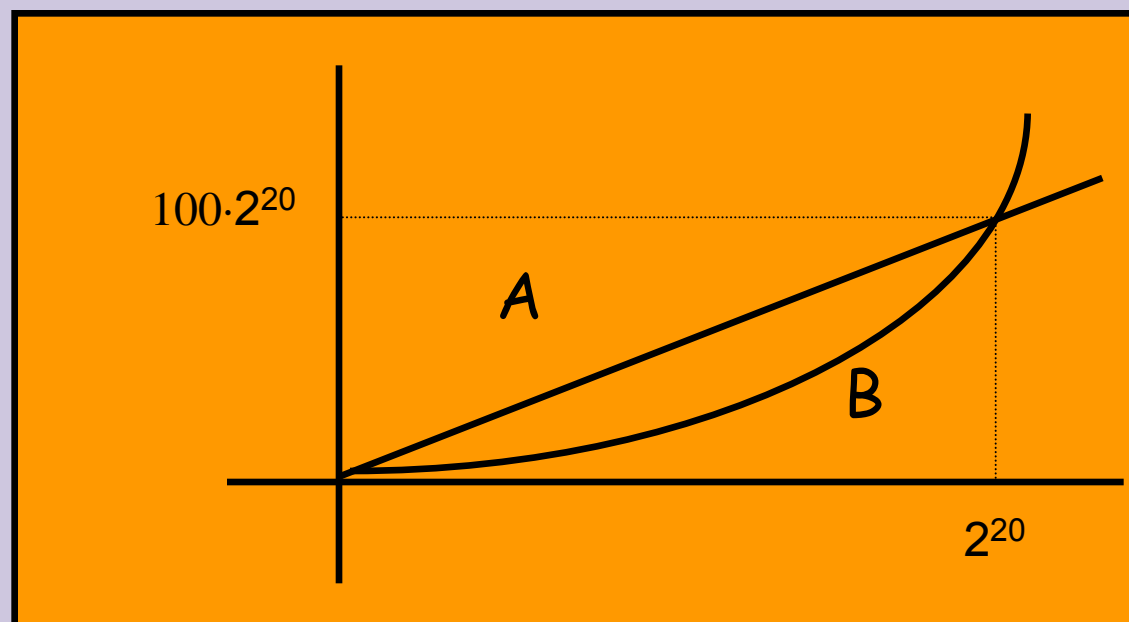
דוגמא.

נניח שאלגוריתם A רץ בזמן $T_A(n) = 100n$

ואלגוריתם B רץ בזמן $T_B(n) = 5n \log_2 n$.

בניתוח אסימפטוטי עדיף אלגוריתם A כיון שהאלגוריתם ליניארי,

אבל עבור קלטים המקיימים $n < 2^{20}$ עדיף אלגוריתם B .



שאלה מבוחן 1993

נתונה הפונקציה הבאה ב-C - נתחו את הסיבוכיות שלה.

```

Void func (int n)
{
    float x, delta;
    int i;
    for (i = 1 ; i <= n ; i++) {
        delta = 1.0 / i;
        x = i;
        while ( x > 0 ){
            x = x - delta;
        }
    }
}

main()
{
    int n;
    scanf("%d",&n);
    func(n)
}

```

$$i = 1, 2, 3, \dots, n$$

$$\text{delta} = 1/i$$

$$x = i, i - \frac{1}{i}, i - \frac{2}{i}, i - \frac{3}{i}, \dots, 0$$

$$\sum_{i=1}^n i^2 + 1$$

$$\sum_{i=1}^n i^2 + 1 \leq n \cdot n^2 + n = n^3 + n$$

$$\sum_{i=1}^n i^2 + 1 \geq \sum_{i=n/2}^n i^2 \geq \frac{n}{2} \cdot \left(\frac{n}{2}\right)^2 = \frac{n^3}{8}$$

$$\sum_{i=1}^n i^2 + 1 = \theta(n^3)$$