

## בעיית התאמת מחרוזות

**בעיה:** נתון טקסט של תווים ומטרתנו למצוא בטקסט הנתון את כל המופעים של "מילה" מסוימת.

**קלט:** מערך  $T$  באורך  $n$  המכיל טקסט (מחרוזת).  
 מערך  $P$  באורך  $m$  (כאשר  $m < n$ ) המכיל תבנית שהיא מילה (מחרוזת) מסוימת.

**פלט:** מספר המופעים של התבנית ( $P$ ) במחרוזת הנתונה ( $T$ ).

טרם נדון בבעיית התאמת מחרוזות נדון בהגדרות, סימונים מקובלים ופעולות על מחרוזות של תווים.

כידוע מידע אינו בהכרח מספרים בלבד, אלא גם צירופי אותיות ומילים. מידע לא מספרי מיוצג לרוב בצורת מחרוזות של תווים.

בהמשך נשתמש בשני מונחים: מילה או מחרוזת, אשר שקולים זה לזה.

## מילים וקבוצות מילים

הגדרה: אלפבית (א"ב) - זוהי קבוצה סופית לא ריקה של סימנים.

סימון:  $\Sigma$  - מסמן את האלפבית (א"ב).

דוגמאות: הקבוצה הבאה:  $\Sigma = \{1, 0\}$  היא א"ב של 2 סימנים אפס ואחד.  
 הקבוצה הבאה:  $\Sigma = \{א, ב, ג, \dots, ת\}$  היא א"ב בן 22 סימנים (אותיות).  
 הקבוצה הבאה:  $\Sigma = \{0, 1, 2, \dots, 9\}$  היא א"ב בן 10 סימנים (ספרות).

הגדרה: מילה (מחרוזת) - זוהי סדרה סופית של אותיות מתוך הא"ב ( $\Sigma$ ).

הערה חשובה! בהמשך נשתמש הרבה פעמים במונח "מעל" אשר מציין "מתוך האותיות של".

דוגמא: נתון:  $\Sigma = \{א, ב, ג, \dots, ת\}$  (א"ב בן 22 אותיות).  
 "אמא" היא מילה מעל (מתוך האותיות של)  $\Sigma$ .  
 "אבא" היא מילה מעל (מתוך האותיות של)  $\Sigma$ .  
 "אאמ" היא מילה מעל (מתוך האותיות של)  $\Sigma$ .  
 "באא" היא מילה מעל (מתוך האותיות של)  $\Sigma$ .

נשים לב לכך ש-2 מילים "אמא" ו-"אאמ" שונות זו מזו ולכן מילה היא סידרה של אותיות ולא קבוצה של אותיות.

הגדרה: מילה ריקה – זוהי מילה שאינה מכילה שום אות או זוהי סידרה סופית של אפס אלמנטים (כלומר זוהי סידרה שאינה מכילה שום איבר).

סימון:  $\varepsilon$  או  $\Lambda$  מסמן את המילה הריקה. (כלומר ישנם שני סימנים מקובלים לציון מילה ריקה).

הגדרה: אורך של מילה – הינו מספר האותיות שבמילה.

סימון:  $|w|$  מסמן את אורך המילה  $w$ .

דוגמאות: נתון:  $\Sigma = \{0,1\}$ , אז

$|0|=1$ , כי במילה 0 ישנו רק סימן אחד והוא אפס.

$|010|=3$ , כי במילה 010 ישנן 3 אותיות בסידרה (במילה).

$|1101|=4$  כי במילה 1101 ישנן 4 אותיות בסידרה (במילה).

$| \varepsilon | = 0$ , כי במילה ריקה ישנם 0 תווים.

הגדרה: שרשור של שתי מילים – היא המילה הנוצרת מ- "הדבקת" המילה השנייה מימין למילה הראשונה.

סימון:  $w_1 \bullet w_2$  – מסמן את שרשורן של שתי מילים  $w_1$  ו-  $w_2$ .  
נקודה " " - מציינת את אופרטור השרשור.

דוגמאות:

1. אם  $w_1$  היא המילה aba כלומר  $w_1 = aba$  ו-  $w_2 = ima$

אזי  $w_1 \bullet w_2 = aba \bullet ima = aba \bullet ima$  ו-  $w_2 \bullet w_1 = ima \bullet aba$

**מסקנה:  $w_1 \bullet w_2$  לא בהכרח זהה ל-  $w_2 \bullet w_1$ , כלומר סדר השרשור חשוב.**

2. אם  $w_1$  היא המילה aba, כלומר  $w_1 = aba$  ו-  $w_2$  היא המילה הריקה, כלומר

$w_2 = \varepsilon$

אזי  $w_1 \bullet w_2 = aba \bullet \varepsilon = aba$

וגם  $w_2 \bullet w_1 = \varepsilon \bullet aba = aba$

**מסקנה: לכל מילה  $w$ :  $w \bullet \varepsilon = \varepsilon \bullet w = w$**

3. אם  $w_1 = ab$  ו-  $w_2 = ab$

אזי  $w_1 \bullet w_2 = ab \bullet ab = abab$

וגם  $w_2 \bullet w_1 = ab \bullet ab = abab$

הערה: בהמשך אנו נוותר על הנקודה " " . המציינת את אופרטור השרשור ובמקום  $w_1 \bullet w_2$  נרשום  $w_1 w_2$  (בלי ציון הנקודה).

דוגמא: אם  $w_1 = ab$  ו-  $w_2 = ba$

אזי  $w_1 w_2 = abba$  ו-  $w_2 w_1 = baab$

שים לב לעובדות הבאות:

1. כל סימן (אות שבאלפבית) הינו מילה באורך 1.

דוגמא: אם  $\Sigma = \{0,1\}$

אז  $|1| = 1$  ו  $|0| = 0$  היא גם מילה באורך 1, וגם

$|11| = 1$  ו  $|10| = 1$  היא גם מילה באורך 1.

2. לכל שתי מילים מעל (מתוך האותיות של)  $\Sigma$  ו-  $w_1$  ו-  $w_2$  מתקיים:

$$|w_1 w_2| = |w_1| + |w_2| \quad \text{א.}$$

ב. לכל סימן (אות)  $a$ , כלומר  $a \in \Sigma$ ,  $|w_1 a| = |w_1| + 1$ , כיוון שהאורך של כל סימן (אות) שבא"ב הוא 1.

דוגמא: אם  $\Sigma = \{0,1\}$  וניקח 2 מילים:  $w_1 = 010$  ו  $w_2 = 10$

$$|w_1 w_2| = |01010| = 5 \quad \text{נקבל:}$$

$$|w_1| = |010| = 3$$

$$|w_2| = |10| = 2$$

ואכן:  $3 + 2 = 5$  כנדרש.

$$|w_1 1| = |0101| = 4 \quad \text{וכן: } |1| = 1 \quad \text{וגם } |w_1| = 3$$

$$\text{ולכן } 3 + 1 = 4.$$

הגדרה: חזקה של מילה

נתונה מילה  $w$  מעל  $\Sigma$ . החזקה של  $w$  מוגדרת ומסומנת כדלהלן:  $w^0 = \epsilon$   
לכל  $i > 0$   $w^i = w^{i-1} \cdot w$

השרשור של מילה  $w$  לעצמה מסומן ב-  $w^2$ . באופן דומה נסמן ב-  $w^3$  את שרשור המילה 3 פעמים.

כלומר  $w^3 = w \cdot w \cdot w$  ו-  $w^3$  מורכבת משלושה מופעים של  $w$ .

דוגמא: אם  $\Sigma = \{a,b\}$  אזי  $(abb)^3 = abb \cdot abb \cdot abb = abbabbabb$

באופן כללי  $w^n = \underbrace{w \cdot w \cdot w \cdot \dots \cdot w}_n$

$n$  מופעים של  $w$

שים לב ל-  $w^0$  אשר מציינת 0 מופעים של  $w$ . זוהי בעצם סדרה של 0 אלמנטים שהיא מילה ריקה.

הגדרה: היפוך של מילה – הוא המילה הנוצרת מהיפוך סדר האותיות שבמילה.

סימון:  $w^r$  – מסמן את האופרטור היפוך של המילה  $w$ .

דוגמאות: אם  $\Sigma = \{0,1\}$  ו-  $w = 1101$  אזי  $w^r = 1011$

הערה: אם  $w = \epsilon$  אזי גם  $w^r = \epsilon$

הגדרה: רישא (prefix) של מילה – אם  $x, w, y$  מילים כך ש-  $x=wy$  אזי המחרוזת  $w$  נקראת רישא של המחרוזת  $x$ .  
סימון:  $w=\text{prefix}(x)$  – מציין שהמחרוזת  $w$  הינה רישא של המחרוזת  $x$ .

דוגמא: נתון  $\Sigma = \{0,1\}$  והמילה  $x=010$  מעל  $\Sigma$ .

$$010 = \epsilon \cdot 010$$

$$010 = 0 \cdot 10$$

$$010 = 01 \cdot 0$$

$$010 = 010 \cdot \epsilon$$

ולכן  $\epsilon, 0, 1, 010$  הן הרישואות של  $x$ .

הערה: כאשר  $w=\text{prefix}(x)$  ברור כי  $|w| \leq |x|$

הגדרה: סיפא (suffix) של מילה – אם  $x, w, y$  מילים כך ש-  $x=wy$  אזי המחרוזת  $y$  נקראת הסיפא של המחרוזת  $x$ .

סימון:  $w=\text{suffix}(x)$  – מציין שהמחרוזת  $w$  סיפא של המחרוזת  $x$ .

דוגמא: נתון  $\Sigma = \{0,1\}$  והמילה  $x=010$  מעל  $\Sigma$  כידוע:

$$010 = 010 \cdot \epsilon$$

$$010 = 01 \cdot 0$$

$$010 = 0 \cdot 10$$

$$010 = \epsilon \cdot 010$$

ולכן  $\epsilon, 0, 10, 010$  הן סיפואות של  $x$ .

הערה: כאשר  $w=\text{suffix}(x)$  ברור כי  $|w| \leq |x|$

**מסקנה: המילה הריקה ( $\epsilon$ ) היא גם רישא וגם סיפא של כל מחרוזת.**

הגדרה: תת-מילה של מילה (substring) – אם  $x, w, y, z$  מילים כך ש-  $x=wyxz$ , אז  $y$  נקראת תת-מילה של  $x$ .

דוגמא: נתון  $\Sigma = \{0,1\}$  והמילה  $x=010$  מעל  $\Sigma$ . כידוע:

$$010 = 0 \cdot 10 \cdot \epsilon$$

$$010 = 0 \cdot 1 \cdot 0$$

$$010 = 0 \cdot \epsilon \cdot 10$$

$$010 = \epsilon \cdot \epsilon \cdot 010$$

$$010 = \epsilon \cdot 010 \cdot \epsilon$$

$$010 = \epsilon \cdot 01 \cdot 0$$

$$010 = \epsilon \cdot 0 \cdot 10$$

$$010 = \epsilon \cdot \epsilon \cdot 010$$

$$010 = \epsilon \cdot 010 \cdot \epsilon$$

$$010 = 0 \cdot 010 \cdot \epsilon$$

$$010 = 01 \cdot 0 \cdot \epsilon$$

$$010 = 010 \cdot \epsilon \cdot \epsilon$$

וכך הלאה.

קל לראות כי כל רישא או כל סיפא היא תת מילה אך לא כל תת-מילה היא רישא או סיפא.

- שים לב! עבור שתי מחרוזות כלשהן  $x$  ו- $y$  ותו כלשהו  $a$  מתקיימים:
  - $x$  היא סופא של  $y$   $[x=\text{suffix}(y)]$  אם ורק אם  $xa$  היא סופא של  $ya$   $[xa=\text{suffix}(ya)]$ .
  - $x$  היא רישא של  $y$   $[x=\text{prefix}(y)]$  אם ורק אם  $ax$  היא רישא של  $ay$   $[ax=\text{prefix}(ay)]$ .

#### הגדרה: שרשור של קבוצות מילים

שרשור של שתי קבוצות מילים הוא הקבוצה שמכילה את כל אפשרויות השרשור של מילה מהקבוצה הראשונה עם מילה מהקבוצה השנייה.

תהיינה  $L_1$  ו- $L_2$  קבוצות מילים מעל  $\Sigma$ . כאמור פעולת השרשור מוגדרת ומסומנת כך:  $L_1 \cdot L_2 = \{xy \mid x \in L_1, y \in L_2\}$

דוגמא: אם  $L_1 = \{0,1,\epsilon\}$  ו- $L_2 = \{1,01\}$

$$0 \cdot 1 = 01$$

$$0 \cdot 01 = 001$$

$$1 \cdot 1 = 11$$

$$1 \cdot 01 = 101$$

$$\epsilon \cdot 1 = 1$$

$$\epsilon \cdot 01 = 01$$

$$L_1 \cdot L_2 = \{01, 001, 11, 101, 1\}$$

המילה 01 נוצרה פעמיים באופנים הבאים:  $0 \cdot 1 = 01$  ו- $\epsilon \cdot 01 = 01$  אך היא תרשם רק פעם אחת ב- $L_1 \cdot L_2$  כיוון שהתוצאה של פעולת השרשור  $L_1 \cdot L_2$  הינה קבוצה וכל איבר בקבוצה מופיע רק פעם אחת (לפי הגדרת הקבוצה).

שים לב! יש חשיבות לסדר השרשור (בדומה לשרשור המילים).

הגדרה: קבוצת כל המילים מעל  $\Sigma$  – זוהי קבוצת כל המחרוזות שאפשר ליצור מן האלפבית  $\Sigma$ .

סימון:  $\Sigma^*$  (כוכב קליני) – מסמן את קבוצת כל המילים מעל  $\Sigma$ .

$$\Sigma^* = \bigcup_{i \geq 0} \Sigma^i$$

דוגמא: נניח כי  $\Sigma = \{0,1\}$  נקבל:

$$\Sigma^0 = \{\epsilon\}$$

$$\Sigma^1 = \Sigma = \{0,1\}$$

$$\Sigma^2 = \Sigma^1 \cdot \Sigma^1 = \Sigma \cdot \Sigma = \{0,1\} \cdot \{0,1\} = \{00, 01, 10, 11\}$$

$\Sigma^2$  מציין אוסף כל המחרוזות, שאורכן בדיוק 2, וש אפשר ליצור אותן מן האלפבית  $\Sigma$ .

$$\Sigma^3 = \Sigma^2 \cdot \Sigma = \{00, 01, 10, 11\} \cdot \{0, 1\} = \{000, 001, 010, 011, 100, 101, 110, 111\}$$

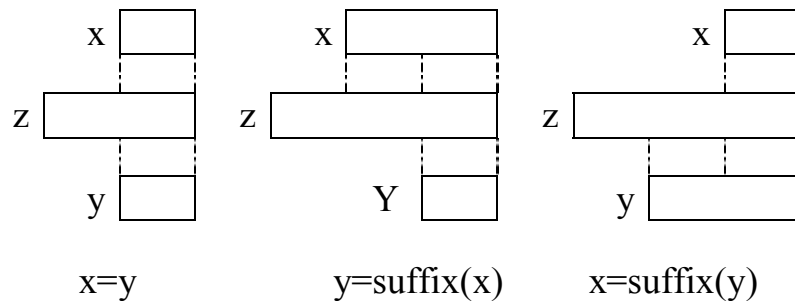
$\Sigma^3$  מציין אוסף כל המחרוזות, שאורכן בדיוק 3, וש אפשר ליצור אותן מן האלפבית  $\Sigma$ .

לכן נסיק כי:  $\Sigma^i$  - תצוין אוסף כל המחרוזות, שאורכן בדיוק  $i$ , וש אפשר ליצור אותן מן האלפבית  $\Sigma$ .

לכן  $\Sigma^* = \bigcup_{i \geq 0} \Sigma^i$  - זוהי קבוצת כל המחרוזות בכל האורכים האפשריים שאפשר ליצור אותן מן האלפבית  $\Sigma$ .

למה: תהינה  $x, y$  ו- $z$  מחרוזות המקיימות:  
 $y = \text{suffix}(z)$  וגם  $x = \text{suffix}(z)$   
 אם  $x = \text{suffix}(y)$  אזי  $|x| \leq |y|$   
 ואם  $y = \text{suffix}(x)$  אזי  $|x| \geq |y|$   
 ואם  $x = y$  אזי  $|x| = |y|$

להלן התיאור הגרפי של הלמה:



טענה: תהי מילה  $x$  כאשר  $x \in \{a,b\}^*$  אם המילה  $x$  מקיימת את התכונה הבאה:  
 $xa=ax$ , אזי  $x=a^n$  עבור  $n$  כלשהו כאשר  $n \in \mathbb{N}$ .

הוכחה:

נוכיח את הטענה בדרך השלילה. נניח ש  $xa=ax$  והמילה  $x$  מכילה תו  $b$ . לכן נוכל לכתוב כי  $x=a^n b u$  ואז מתקיים:  $a^n b u \cdot a = a \cdot a^n b u$   
 לכן מתקיים:  $bua = abu$   
 אך לא יתכן ששתי המילים מקיימות את השוויון, כלומר  $abu = bua$  כיוון ש  $a \neq b$ , כי אם שתי מחרוזות שוות אזי גם התו הראשון שלהן צריכים להיות זהים.  
 לכן הגענו לסתירה. הסתירה נובעת מההנחה שהמילה  $x$  מכילה תו  $b$  ולכן המסקנה היא שהמילה  $x$  לא מכילה תו  $b$  כלומר  $x = a^n$  עבור  $n$  כלשהו. מש"ל

טענה: לא קיימת מילה  $x$  כאשר  $x \in \{a,b\}^*$  כך שיתקיים  $ax=xb$ .

הוכחה: אם המילה  $x \in \{a,b\}^*$  כך ש  $|x|=n$  (האורך של המילה  $x$ ) אזי עלינו להוכיח כי לכל  $n \in \mathbb{N}$   $ax \neq xb$ .

נוכיח את הטענה באינדוקציה מתמטית.

בסיס: כאשר  $n=0$  כלומר  $|x|=0$  אזי המילה  $x$  היא מילה ריקה ולכן  $ax \neq xb$  כיוון ש  $x = \_$  ו  $a \neq b$ .

הנחת האינדוקציה: נניח את נכונות הטענה לכל  $m$  כאשר  $m < k$  עבור  $k$  כלשהו נתון, כלומר אם  $|x|=m$  ו  $m < k$  אזי  $ax \neq xb$ .

עתה נוכיח את נכונות הטענה עבור המילה  $x$ , כאשר  $|x|=k$ . נראה זאת בדרך השלילה, כלומר נניח ש  $|x|=k$  ו  $ax=xb$  מאחר ו  $a \in \Sigma$ ,  $b \in \Sigma$  ו  $x \in \Sigma^*$  כאשר  $\Sigma = \{a,b\}$  ו  $ax=xb$  כלומר:

a	x
שווה ל-	
x	b

אזי ברור כי: המילה  $x$  מתחילה באות  $a$  ומסתיימת ב  $b$ .  
 אי לכך נוכל לכתוב כי  $x=awb$  כאשר  $w \in \{a,b\}^*$  ו  $|w| < |x|$ . אז נקבל:  
 $ax=xb$   
 $a \cdot awb = awb \cdot b$   
 כלומר  $aw = wb$

אך כאמור  $|w| < |x|$  לכן לפי הנחת האינדוקציה היינו צריכים לקבל כי:  $aw \neq wb$   
 קיבלנו סתירה עקב הנחתנו ש  $ax=xb$  ו נסיק כי  $ax \neq xb$  לכל מילה  $x$ . מש"ל.

## תרגילים:

1. נתון כי  $A_1 = \{a, b\}$  ו  $A_2 = \{1, 2, 3\}$ 
  - א. מצא את:  $A = A_1 \cdot A_2$ .
  - ב. בהמשך לסעיף א' מצא את  $A^r$  כאשר  $A^r = \{w^r | w \in A\}$ .
  - ג. מצא את כל הרישיות (prefixes) ואת כל הסיפות (suffixes) של כל מילה שבקבוצה  $A$ .
2. נתון כי:
  - א.  $\Sigma = \{0, 1\}$
  - ב.  $A = \{0^n 1^n | n \geq 1\}$  כלומר  $A$  היא אוסף המילים שהן שרשרת של מספר כלשהו של  $a$ 'ים ולאחריה שרשרת של אותו מספר  $b$ 'ים.
  - ג.  $B = \{w \in \{0, 1\}^* | \#_0(w) = \#_1(w)\}$  כלומר  $B$  הוא אוסף של המילים מעל  $\{0, 1\}$  שבהן מספר האפסים שווה למספר האחדים.  
הסימון  $\#(w)$  פירושו מספר המופעים של האות  $a$  במילה  $w$ .  
כך למשל  $\#_a(ba) = 2$  ו  $\#_a(bb) = 0$ .
1. עבור כל המילים הבאות קבע אם היא שייכת לקבוצה  $A$  או לאו.  
 $0011, 0101, 1100, 00111, 00011$ . נמק תשובתך.
2. אותה שאלה כמו בסעיף א' כאשר נבדקת השייכות לקבוצה  $B$ .
3. האם  $A \subset B$  או  $B \subset A$  או  $A = B$ ? נמק תשובתך.
4. האם  $A \subset \Sigma^*$ ? נמק
3. נניח ש-  $x, y, v, w \in \Sigma^*$  ונניח ש:  $xy = vw$  אזי קיימת  $z \in \Sigma^*$  כך ש:
  1.  $x = vz$  ו  $w = zy$
  - או
  2.  $v = xz$  ו  $y = zw$
4. נניח ש-  $x, y \in \Sigma^*$  הראה שאם  $xy = yx$  אזי קיימת  $z \in \Sigma^*$  כך ש-  $x = z^m$  ו  $y = z^n$  עבור  $m, n \geq 0$ .
5. נניח  $x, y, z \in \Sigma^*$  הראה שאם  $xz = zy$  אזי  $x^k z = zy^k$  לכל  $k \geq 0$ .
6. שתי מילים  $x, y \in \Sigma^*$  נקראות צמודות אם קיים  $z \in \Sigma^*$  כך ש-  $xz = zy$ . הראה ש-  $x$  ו-  $y$  מילים צמודות אם ורק אם קיימות  $u, v \in \Sigma^*$  כך ש:  $x = uv$  ו  $y = vu$ .
7. נניח ש-  $x, y, z \in \Sigma^*$  הראה שאם  $x^2 y^2 = z^2$  אזי  $xy = yx$ .



8. הראה כי לא קיימות שתי מחרוזות  $x, y \in \{a, b\}^*$  כך ש:  $xay = ybx$ .

9. לפניכם טענה: אם  $u, v \in \{a, b\}^*$  ו-  $u \neq v$ , אזי לכל שתי מילים  $x, y \in \{a, b\}^*$  מתקיים  $xuy \neq yvx$ .

א. הראו בעזרת דוגמה נגדית שהטענה איננה נכונה.

ב. נתונה ה"הוכחה" הבאה לטענה:

נעיין בפרדיקט הבא: אם  $x, y \in \{a, b\}^*$  ו-  $|xy| = n$ , אז לכל שתי מילים  $u$  ו-  $v$  מעל  $\{a, b\}$ , מתקיים  $xuy \neq yvx$ .

נוכיח את נכונותו של הפרדיקט לכל  $n \in \mathbb{N}$ :

עבור  $n=0$  הוא נכון.

נניח את נכונותו לכל  $m < k$ , עבור  $k$  נתון כלשהו, ונוכיח את נכונותו ל-  $k$ :

נניח בשלילה שקיימות שתי מילים  $x, y$  כך ש-  $|xy| = k$ , ושתי מילים שונות  $u$  ו-  $v$  כך ש-  $xuy = yvx$ .

מן השוויון הזה נובע שהאות הראשונה של  $x$  זהה לאות הראשונה של  $y$ , והאות

האחרונה של  $x$  זהה לאות האחרונה של  $y$ .

נסמן:  $x = \alpha x_1 \beta$ ,  $y = \alpha y_1 \beta$ ,  $(\alpha, \beta \in \{a, b\}^*)$ .

נקבל:  $\alpha x_1 \beta u \alpha y_1 \beta = \alpha y_1 \beta v \alpha x_1 \beta$

ומכאן:  $x_1 \beta u \alpha y_1 = y_1 \beta v \alpha x_1$

מכיוון ש-  $u \neq v$ , גם  $\beta u \alpha \neq \beta v \alpha$

נסמן:  $u_1 = \beta u \alpha$  ו-  $v_1 = \beta v \alpha$

נקבל:  $x_1 u_1 y_1 = y_1 v_1 x_1$ .

אבל  $|x_1 y_1| < k$  ולכן על-פי הנחת האינדוקציה,  $x_1 u_1 y_1 \neq y_1 v_1 x_1$ . סתירה!

מהסתירה ניתן להסיק שהפרדיקט נכון גם ל-  $k$ , ועל כן הוא נכון לכל  $n \in \mathbb{N}$ .

מה כאן לא בסדר?

עתה נחזור לבעיה המקורית – בעיית התאמת מחרוזות (string – matching problem). בבעיה זו נתון טקסט (מחרוזת של תוים) באורך  $n$  שנמצא במערך  $T[0 \dots (n-1)]$  ונתונה תבנית (pattern) (מחרוזת של תוים) שאורכה  $m$ , כאשר  $m < n$ , והיא מוחזקת במערך  $P[0 \dots (m-1)]$ . המחרוזות  $T$  ו- $P$  מעל אותו אלפבית  $\Sigma$ .

המטרה העיקרית בבעיה זו למצוא את כל המופעים של התבנית  $P$  בטקסט הנתון  $T$ , כלומר עלינו למצוא את כל האינדקסים  $k$  במערך  $T$  כך ש:  $T[k..k+m-1] = P[0..m-1]$  כאשר  $0 \leq k \leq n-m$

דוגמא : נתון :

T טקסט	b	a	b	c	a	b	a	b	b	a	b	a	b
P תבנית	0	1	2	3	4	5	6	7	8	9	10	11	12
	aba												

התבנית  $P$  מופיע בטקסט  $T$  פעמיים. החל מאינדקס 4, כלומר  $T[aba] = P[aba]$   
 456      012

והחל מאינדקס 9, כלומר  $T[a b a] = P[aba]$   
 91011      012

כלומר עבור כל  $j$  מ-0 עד 2 מתקיים:  $T[k+j] = P[j]$  כאשר  $k=4$  או  $k=9$ .

### אלגוריתם א' – אלגוריתם הנאיבי

עתה נתאר פונקציה (קטע קוד) אשר מקבלת שתי מחרוזות  $T$  ו- $P$ . הפונקציה תדפיס את כל השלמים ( $k$ ) המציינים את מיקום המופע של התבנית  $P$  בטקסט  $T$ . עבור כל אינדקס  $k$  שמקבל ערכים מאפס עד  $n-m$  נבדוק האם  $T[k..k+m-1] = P[0..m-1]$ .

האיור הבא מתאר את הרעיון של האלגוריתם הנאיבי להתאמת מחרוזות. לפעמים תמונה אחת שווה אלף מילים וגם במקרה זה האיור שמטה מתאר את ה"אלף מילים" ללא מילים.

T[0]	T[1]	T[2]	T[3]	T[4]	T[5]	T[6]	T[7]	T[8]	T[9]	T[10]
A	B	C	E	F	G	A	B	C	D	E
P[0]	P[1]	P[2]	P[3]							
A	B	C	D							

A	B	C	D
---	---	---	---

A	B	C	D
---	---	---	---

...

A	B	C	D
---	---	---	---

```

n=length(T);
m=length(P);
for(k=0;k<=n-m;k++)
{
    /* החל מהנקודה הזו מתחילה לולאה המבצעת בדיקה האם P[0..m-1]
    /*1]=T[k...k+m-1]

    count=0;
    for(j=0;j<=m-1;j++)
    if(P[j]==T[k+j]) count ++;

    /*עתה נבדוק האם היתה התאמה מלאה*/
    if(count==m)printf("%d\n",k);
}

```

הלולאה החיצונית רצה בזמן  $0(n-m+1)$  ו- הלולאה הפנימית רצה בזמן  $0(m)$ . כיוון שהלולאות מקוננות זה בזה, אזי לפי עיקרון הכפל סיבוכיות זמן הריצה של האלגוריתם הינה  $0((n-m+1)m)$ .

#### שיפור לאלגוריתם א'

נתבונן על המקרה הבא כאשר  $k=6$ .

T	a	b	c	a	b	c	a	b	a	a	c	d	...
	0	1	2	3	4	5	6	7	8	9	10	11	....

P	a	b	a	a	b	e
	0	1	2	3	4	5

באלגוריתם א' אנו מבצעים את ההשוואה הבאה:  $T[6..11]=P[0..5]$ . ניתן לראות כי:  $P[3]=T[9], P[2]=T[8], P[1]=T[7], P[0]=T[6]$ . אך  $P[4]=b$  שונה מ-  $T[10]=c$ .

לכן מיותר להמשיך את ההשוואות הבאות ואין טעם לעבור על ה-  $j$  - ים הגדולים מ-  $j_1$  כאשר  $P[j_1] \neq T[k+j_1]$ .

לאור זאת נציע אלגוריתם משופר יותר אשר משתמש במשתנה בוליאני הבא:

כאשר נמצא אינדקס כלשהו  $j_1, 0 \leq j_1 \leq m$  כך ש-  $P[j_1] \neq T[k+j_1]$   $\left\{ \begin{array}{l} \text{true} \\ \text{false} \end{array} \right.$   $\text{finis}$  אחרת

ברור כי בתחילת האלגוריתם לפני שמתחילים את הבדיקה האם  $P[0..m-1]=T[k..k+m-1]$  המשתנה  $\text{finis}$  יקבל את הערך  $\text{false}$  כיוון שלפני הסריקה עדין לא מצאנו אינדקס  $j_1$  כך ש-  $P[j_1] \neq T[k+j_1]$ .

להלן האלגוריתם המשופר

```
#define true 1
#define false 0
n=length(T); /* n ← T המחרוזת */

m=length(P); /* m ← P המחרוזת */
P[0..m-1]=T[k..k+m-1] החל מנקודה זו מתחילה הלולאה המבצעת בדיקה האם
for (k=0;k<=n-m;k++) /* 1 */
{
    finis=false; j=0;
    while ((j<=m-1)&&(!finis))
        if (T[k+j]==P[j]) j++;
        else finis=false;
    if (j>m-1)
        /* האם היתה התאמה מלאה */
        printf("%d\n",k);
}
```

נתבונן על קלט המייצג את המקרה הגרוע ביותר.

נניח  $\Sigma=\{0,1\}$

$$T = \underbrace{11111\dots 1}_n = 1^n$$

פעמים n

$$P = \underbrace{11\dots 1}_m = 1^m$$

פעמים m

וברור  $m \leq n$

בהינתן קלט כזה, עבור כל ערך של  $k$ , לולאת while תתבצע בכל מקרה  $m$  פעמים. לכן באלגוריתם המשופר אנו לא מקטינים את מספר הצעדים במקרה הגרוע ביותר ולכן זמן הריצה של האלגוריתם המשופר נשאר  $O((n-m+1)m)$ .

הערה: כאשר  $m=O(n)$  אז סיבוכיות זמן הריצה של אלגוריתם אי' (וגם של הגירסה המשופרת) היא  $O(n^2)$ .

נראה בהמשך אלגוריתמים אחרים הפותרים את הבעיה הנתונה בזמן שהוא יותר קטן מ-  $O((n-m+1)m)$  ולכן אלגוריתם אי' אינו אופטימלי עבור בעית התאמת מחרוזות.

הערות:

א. ניתן לממש את האלגוריתם המשופר, ביחס לאלגוריתם א', גם ללא שימוש במשתנה עזר – משתנה בוליאני (finis) כדלקמן:

```
n=length(T);
m=length(P);
for (k=0;k<=n-m;k++)
{
    i=0;
    while((i<m)&&(P[i]==T[k+j])) i++;
    if(i==length(P)) printf("%d\n",i);
}
```

ב. אם כל התאים של התבנית  $P$  שונים זה מזה, אז ניתן לפתור את בעיית התאמת מחרוזות בזמן לינארי –  $O(n)$ . (מדוע? חשוב היטב!)