

נושא 6 - צח'סות מילונ'ית

(א) ה'קצמה

עצם זה למדנו שתי גישות כיצד עדיין קופץ: גישה הסטלית מנתח את הקופץ עצמי הקיצוב כדי לקבל סטטיסטיקה יאז מקוצרים אורו, ובגישה האצטל'ית (דינאמית) מקוצרים כבר במצבה הראשון על הדובל כע תו עפי מה שראנו עצם אלו. גישה חדשה היא "גישת ההחלטה", שבה עובדים על הטקסט משאלא ע'מין ומחל'פים בעזרים של תולים במצב (הסרנס) לעזרים הקיצים בו ראנו בעיקר זה.

(ב) צח'סות מילונ'ית

כזה צח'סות שבה נבנה מעון המכיל כפי כחיסה תו עוצצ או קבוצת תולים וקיצוב של יתו הקוצצ או קבוצת התולים. תכנ'צ המקוצצ הוא עמלק את הטקסט בצורה כזו שכל חלק נמצא במעון. יכולים להיות המון אפסמיות עמלקה זו. הצדית העיקרה אינסט'מיל'ית עפיר מעון נתון:

ב'ינתן מעון כ, בונק'צית קיצוב ה (מקפלת מילה ומהלכה את הקיצוב שבה עפי המעון פ וטקסט T , נח'זה עמק'צ חסיה של T עמ'ים $T = w_1 w_2 \dots w_k$, כך שכל מילה $w_i \in D$ ואם סלק אורק כל הקיצבים עפי חלוקה זו. כעוצר, $\sum_{i=1}^k |w_i|$ מנימלי.

חלוקה אינסט'מיל'ית: שיטת החלוקה האינסט'מיל'ית היא ענשות רפוק'ציה עפ'יה מתורת הזרמים. נבנה

$$T = x_1 x_2 \dots x_n$$

$$G = (V, E)$$

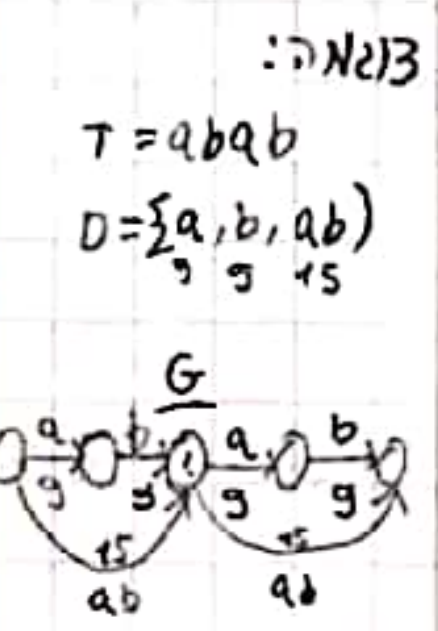
$$V = \{1, 2, \dots, n, n+1\}$$

$$(i, j) \in E \Leftrightarrow x_i x_{i+1} \dots x_{j-1} x_j \in D$$

$$w(i, j) = |x_i x_{i+1} \dots x_{j-1} x_j|$$

$$P_{T,G}(1, n+1) = \text{מסלול ארוך ביותר}$$

זרם G שבו הצמתים הם כחסי התולים ב- T ואז צח'סות. הצמתים מסוגרים עפי ההוספה שלהם ב- T . כדי להגדיר את הקשתות נעזיר על זה בזה צמתים (i, j) ונבדוק אם יש מילה במעין שמתחילה ב- i ונגמרת ב- j , אם כן $(i, j) \in E$. הקשת ממוט'ל משקע כע חסית ית'ה אורק הקיצוב במעין של המילה.



כדי עמק'צ חלוקה אינסט'מיל'ית נסד'ל על G אעמ'ית דואקטורה מ-1 עד $n+1$. כע קשת במסלול שיהיה הוא חלק אחצ בחלוקה האינסט'מיל'ית.

יש צ'צ סוגים של חלוקות אך אין אינסט'מיל'יות (שקופות צ במצבת).

הסרק זה נמצא עוצ סוגים של בחיסות מילונ'יות שהן אצטל'יות, כעוצר מקוצרים איתן במעבר אחצ

משאלא ע'מין. יעכן אע' סחית טיבות מחסיה אינסט'מיל'ית, אך יש עהם יתרון שקע יותר עקיצצ אותן ואין צח'סות מילונ'ית יצור מראש!

(28)

LZ77 (2)

זהו אלגוריתם קידוד שפותח על ידי תוקיים (Lempel, Ziv) בשנת 77, ואינו מבוסס על אלגוריתמי פתיסה שמשתמשים בהם צד הימני בקודיף כמו Ziv ו-Jarvis. מבוסס על "גישת ההחזרה".

באלגוריתם זה מתקדים את הטקסט ששני חלקים: Window ו-Look Ahead Buffer. זה החלק שבו קודדנו ו-Look Ahead Buffer זה מה שנקרא בתחילת האלגוריתם ה-Look Ahead Buffer הוא גם הטקסט וה-Look Ahead Buffer ריק. * זכור לנו הטקסט ושמאל עינינו. בכל פעם שמצאנו עברו בודקים מהי קבוצת התווים שמתחילה מאחורינו ויציאה שבה נמצאת ב-Look Ahead Buffer, ואנחנו קבוצת תווים לו ברורים למקום ב-Look Ahead Buffer שבו נכח נמצאת קבוצת תווים שלו. ה-Look Ahead Buffer משמש כמילון "באלגוריתם זה ברורים הוא הקידוד".

בצד שמאל ו-Look Ahead Buffer זה מה שנקרא בתחילת האלגוריתם ה-Look Ahead Buffer הוא גם הטקסט וה-Look Ahead Buffer ריק. * זכור לנו הטקסט ושמאל עינינו. בכל פעם שמצאנו עברו בודקים מהי קבוצת התווים שמתחילה מאחורינו ויציאה שבה נמצאת ב-Look Ahead Buffer, ואנחנו קבוצת תווים לו ברורים למקום ב-Look Ahead Buffer שבו נכח נמצאת קבוצת תווים שלו. ה-Look Ahead Buffer משמש כמילון "באלגוריתם זה ברורים הוא הקידוד".

כיצד נראה רפרנט כזה? הוא מתואר בטבלה תוקיים (Symbol, Length, Offset):

- Offset - כמה תווים לפני שמצאנו מילון.
- Length - כמה תווים קודדו מ-Look Ahead Buffer. אק ה-Length נכנס אל-Look Ahead Buffer. כמה תווים קודדו מ-Look Ahead Buffer. כמה תווים קודדו מ-Look Ahead Buffer.
- Symbol - התו הראשון שנמצא ב-Look Ahead Buffer אחרי קבוצת התווים שבה נמצא.

זה נקרא "רפרנט" וזה מה שנקרא בתחילת האלגוריתם ה-Look Ahead Buffer הוא גם הטקסט וה-Look Ahead Buffer ריק. * זכור לנו הטקסט ושמאל עינינו. בכל פעם שמצאנו עברו בודקים מהי קבוצת התווים שמתחילה מאחורינו ויציאה שבה נמצאת ב-Look Ahead Buffer, ואנחנו קבוצת תווים לו ברורים למקום ב-Look Ahead Buffer שבו נכח נמצאת קבוצת תווים שלו. ה-Look Ahead Buffer משמש כמילון "באלגוריתם זה ברורים הוא הקידוד".

במילים נותקלים בתו C בפעם הראשונה ברפרנט שלו הוא (C, S, 0). המספר בשטח הימני של הקודד אחרי "הראש".

אלגוריתמי קידוד ופענוח:

LZ77 Encoding

```
p = 1
while p < s.length // s זהו מילון המכיל את הטקסט
  search for the longest match for
  s[p...] in s[p-w...p-1] // נניח שצא אחידות באורך w
  output (p-m, l, s[p+l]) // m הוא המרחק מ-p-w
  p = p + l + 1
```

LZ77 Decoding

```
p = 1
for each triple (f, l, c)
  s[p...p+l-1] = s[p-f...p-f+l-1]
  s[p+l] = c
  p = p + l + 1
```

בדוגמה:

T = "A - walrus - in - Spain - is - a - walrus - in - vain"

Encoded String: (0,0,A), (0,0,-), (0,0,w), (0,0,a), (0,0,i), (0,0,r), (0,0,u), (0,0,s), (7,1,i), (0,0,t), (3,1,s), (0,0,p), (1,1,i), (6,2,i), (12,2,a), (21,11,r), (20,3,.)

תסכולות:

- (1) ב-LZ77 לא ידוע כיצד עוצמו מחזרות הלי גרסה שנמצאת ב-Look Ahead Buffer. יתכן שהיה לא צריך לבצע עיבוד הדדיות.
- (2) במקרה שבו יש לנו מספר הראשונה נקודת (C, S, 0) שבה אדם החלון בפעם וצד אסטר הפעם הראשונה את מספר התווים שבו מרשים להעתיק מה-Look Ahead Buffer וצד 8 ב"ס. שטח החפץ ב-ASCII. אלו הם המין פועם עינים, כפי היה נדיר להעתיק את C עכברו.
- (3) אנו עוצמו ב-Look Ahead Buffer מחזרות הלי גרסה ב-Look Ahead Buffer? אולי כדאי פונקציה דירה יותר כי יאסטר ענו בהמשך להעתיק מחזרות גרסה יותר.

(29)

LZSS (3)

(3)

במ'צה ובחור סה'ע'ק
יה א תל'ס ו'א'ע'ה
א'ס'ר סה'ע'ס נ'ה-ח'ע'
י'ח'ס'ל א-א י'ג'ק י'ה'ר
ז'ט'ס'י סה'ע'ת'ק י'ת'ר

בואר: ציור buffer
קאדע 4 וואכיות
עקציוו "אבאבאבאבאבא"
אן קורליק אט היט
מאכיות עתידס היא
"אבא", וואס צו כה

```

graph TD
    3((abba  
3)) --- 1((ab  
1))
    3 --- 2((ba  
2))
    2 --- 4((a  
4))
    style 2 stroke-width:2px
  
```

ניתן להאמת שכל מסלול
עם פסגה היא צומת
ב מעקים 2 בחציון. וכן
רמזים רחוקים הוא (3)
 $ab(2,2)(\underline{3},3) \dots$

המטרה היא למצוא
קיצוץ של 700,000.

10
 11
 12
 13

444

with CamScanner

LZS

(30)

LZ78 (א)

זהו אלגוריתם שפותח בשנת 1978 על ידי איתן חזקיהו שפיתח את LZ78. ההפך הארכיבי ביניהם

הוא שם - LZ78 - ה- Wadsworth משמש כאמין ואילו - LZ78 - פונקט משמש אמין. המקור: בונה

את המילים תוך כדי כעלות הקידום עד שיש יצא מקדמ על הטקסט. המילים מוכנסות מאגס של

מחרוזות (נקראים גם ביטויים). $D = Z_0, Z_1, Z_2, \dots$. מטרת המקור: היא להפוך את הטקסט לחלקים

חדשים $T = T_1 T_2 \dots T_n$ כך שכל $T_i \in D$. במקום כל ביטוי T_i נכלה המקור: שם המילים האחרות

הערה: נקראת "רצף" (מאגס).

מקדמ (אמגס, אגמח), כאשר אגמח זהו המילה של הביטוי שנמצא במילון ואילו אנו מוסיפים להפנות

י- אמגס זהו המילה של הטקסט שאחר ביטוי זה.

כיצד עובד האלגוריתם: בתחילה המילים $Z_0 = \epsilon$. בכל פעם שקוראים מחרוזת X , מחשבים את

האנדקס i של הביטוי שנמצא במילון שהוא היחיד המכיל את הטקסט שניתר לקידום. האנדקס i (שם)

עם שתמצא יש אנדקס i כל, כי במקרה הזריז נחזיר i , שהוא האנדקס של ϵ . ניקח את הביטוי Z_i שאוסף

במילון נוסף על זה את המילה הבאה שנמצא בטקסט ונכנס את Z_{i+1} בתור המילה הבאה במילון שתכלה

את האנדקס הסופי הבא במילון. איתו החלק של Z_i שנמצא בטקסט T נאמר בירור (i, y) .

$T = baadaabaab$

כעת תשוב נוסף, הוא שהמילה בונה את המילים תוך כדי העלאת העין ליד

Index	Phrase	Encoding	# of bits
0	ϵ		
1	b	(0,b)	0+2
2	a	(0,a)	1+2
3	d	(0,d)	2+2
4	ad	(2,d)	2+2
5	ada	(4,a)	3+2
6	ba	(1,a)	3+2
7	ab	(2,b)	3+2

עלולה את המילים ב- $header$. תבונה זו מאפשרת לקודם את האנדקס

באנדקס קבוע. אלא משתנה בהתאם לאנדקס המילים. אנו מנצלים את המילים כל פעם

בתחילה של 2^i . כל מילה המילים באנדקס 2^i נחש א ביטויים שצין את אגס האנדקס.

ההחלפה שבה המילים
היא i שכן אין צורך
לציין אנדקס.

מרגע שהמילים ידוע כי 2^i - 2^{i+1} מרגע 2^{i+1} ביטויים שצין את האנדקס.

(1, a) (2, d) (0, a) (0, b) (0, d) (4, a) (1, a) (2, b) (1, a)

בנוסף, אגס i יתקבל משהו (i, y) וזמן בדיקת שני ביטויים שצין כל מילה. ניתן

לראות איך מסתדר הביטויים המצוינים את האנדקס ידוע ככל שהמילים ידוע.

LZW (ב)

זהו אלגוריתם דחיסה שמבוסס על LZ78 ומשנה אותו מרגע. ההוספות והשינויים העיקריים הם:

(1) המילים בהתחלה אינן מאוחסנות עם ϵ כפי שזכרנו אלא עם כל המילים המופיעות בקובץ. אגס המילים זה

יהיה המילה שחנקה של 2 הכי קטנה מאסטר המילים. בהתחלה נרשמו של LZW המילים בהתחלה

מאוחסנות עם כל 254 המילים ב- $ASCII$ והמילה הבאה היא 255 . לכן 255 ביטויים שצין אנדקס

(2) כל רפרנס הוא אנדקס בודד ואינו מכיל את המילה הבאה. בכל פעם שנמצא מחרוזת מהטקסט שנמצא

במילון, נוסף על זה את המילה הבאה בטקסט y ונכנס למילון קודק חדש. לאחר מכן הקריאה תתחיל מ- y ולא להת

מחדש. המילה הבאה נכנסת למילון. המילה הבאה נכנסת למילון. המילה הבאה נכנסת למילון.

