

## סיכום אלגוריתמים – דחיסת נתונים:

### הגדרות וסימונים:

אלף בית בקובץ מקור -  $S = [s_1, s_2, \dots, s_n]$   
הסתברויות של האלפבית -  $P = [p_1, p_2, \dots, p_n]$ , כאשר  $p_1 \geq p_2 \geq \dots \geq p_n$  ובנוסף  $\sum_{i=1}^n p_i = 1$ .  
קידודי המילים -  $C = [c_1, c_2, \dots, c_n]$   
אורכי הקידודים -  $|C| = [|c_1|, |c_2|, \dots, |c_n|]$   
אורך מילת קוד ממוצעת (Expected codeword length) -  $E(C, P) = \sum_{i=1}^n p_i * |c_i|$   
קוד חסר רישות (prefix free) – אף מילת קוד היא לא תחילית של מילת קוד אחרת.  
קוד UD – כל מחרוזת של קידודי מילים ניתנת לפיענוח בצורה יחידה.  
קוד שלם – כל מחרוזת קידודים אינסופית למחיצה ניתנת לפיענוח בצורה ייחודית.  
קוד מידי – מחרוזת קידודים שהמפענח יודע לפענח גל תו ברגע שהוא סיים לקרוא את מילת הקוד שלו.  
קוד חסר יתרות (Minimum Redundancy Code) – עבור הסתברויות נתונות, אין קידוד אחר שמביא להצטמצמת expected codeword length אם קטן יותר, פורמלית: יהי  $E(C, P)$  אורך מילת קוד ממוצעת עבור קידוד  $C$ ,  $C$  ייקרא קוד חסר יתרות עבור סדרת התפלגויות  $P$  אם  $E(C, P) \leq E(C', P)$  לכל קוד חסר רישות  $C'$ .  
אינפורמציה – כמות האינפורמציה המכילה בתו  $s_i$  עם הסתברות  $p_i$  היא  $I(s_i) = -\log(p_i)$ , קוד יכול להיות מתוכנן כך שמילת הקוד של  $s_i$  מורכבת מ- $I(s_i)$  ביטים.  
אנטרופיה -  $H(P) = -\sum_{i=1}^n p_i * \log(p_i)$ , ובנוסף לכל קידוד  $C$  מתקיים  $H(P) \leq E(C, P)$ .  
אי שוויון קראפט – עבור קוד  $C$  בעל  $n$  תווים, אם  $C$  הוא UD אז מתקיים:  $K(C) = \sum_{i=1}^n 2^{-|c_i|} \leq 1$ .  
קידוד Fixed length – קידוד תווים שבו כל קידוד הוא עם אורך זהה, למשל קידוד ASCII, מייצג כל תו ב-8 ביטים.

### משפטים:

קוד UD אינו בהכרח מידי, יכול להיות קוד UD שמפוענח בצורה יחידה רק אחרי שסיימנו לקרוא את כל המחרוזת.  
קוד מידי אמ"מ קוד חסר רישות.  
אם  $K(C) \leq 1$  עבור קידוד  $C$ , אז קיים קוד  $C'$  כך ש- $E(C, P) = E(C', P)$  ובנוסף  $|C'| = |C|$  הוא קוד חסר רישות.  
אם  $K(C) > 1$  אזי הקוד עם האורכים הנתונים לא יכול להיות חסר רישות.  
קוד מידי עם מילות קוד באורכים  $l_1 \leq l_2 \leq \dots \leq l_n$  קיים אמ"מ  $\sum_{i=1}^n 2^{-l_i} \leq 1$ .  
תנאי הכרחי עבור קוד להיות UD הוא שמתקיים  $\sum_{i=1}^n 2^{-l_i} \leq 1$ .  
קוד הוא שלם אם  $\sum_{i=1}^n 2^{-l_i} = 1$ .

### אלגוריתם לבדיקת UD:

- בנה רשימה של כל מילות הקוד.
- אם קיימת מילת קוד שהיא תחילית של מילת קוד אחרת, תוסיף את הdangling suffix לרשימה (במידה ואין ברשימה מחרוזת זהה) עד ש:
  - נקבל dangling suffix שהוא מילת קוד מקורית – החזר שהקוד אינו UD.
  - אין יותר dangling suffix ייחודיים (שלא מופיעים כבר ברשימה) – החזר שהקוד הוא UD.

- Codewords {0,01,10}
- 0 is a prefix of 01 → dangling suffix is 1
- List - {0,01,10,1}
- 1 is a prefix of 10 → dangling suffix is 0 - which is an original codeword!
- the code is not UD

- Codewords {0,01,11}
- 0 is a prefix of 01 → dangling suffix is 1
- List - {0,01,11,1}
- 1 is a prefix of 11 → dangling suffix is 1 - already there
- the code is UD

### מערכת דחיסה סטטית (מודל מסדר 0):

ההסתברויות לתו אינן תלויות בהודעה שיש לדחוס (נניח בASCII אם בהודעה אין a, מבחינתנו ההסתברות לראות a היא עדיין  $1/256$ ). לכן,  $p_i = \frac{1}{256}$ ,  $H(P) = -\sum_{i=1}^{256} \frac{1}{256} * \log\left(\frac{1}{256}\right) = 8.0$ .

### מערכת דחיסה סמי-סטטית (מודל מסדר 0):

נבחן את ההודעה כדי להבין מי הם התווים שמופיעים בה – במערכות כאלו נדרש *prelude* להודעה שמכיל תיאור של התווים הנ"ל. למשל עבור ההודעה "Bring me my bow of burning gold!" נגדיר  $p_i = 1/25$  למרות שהתו b מופיע יותר מפעם אחת, כלומר

נברר כמה תווים שונים יש וניתן להם התפלגות אחידה. מכאן  $H(P) = -\sum_{i=1}^{25} \frac{1}{25} * \log\left(\frac{1}{25}\right) = 4.64$ . מה מכיל ה*prelude*?

- 8 ביטים של גודל האלפבית (מספר בין 1 ל-256 ולכן במקרה הגרוע נצטרך 8 ביטים).
  - $n * 8$  ביטים עבור התיאור של כל אחד מהתווים ( $n = 25$  בדוגמה שלנו).
- אורך מילת הקוד הממוצעת הוא בסה"כ  $4.64 + \frac{208}{|message|}$ .

### מערכת דחיסה סמי-סטטית עם הסתברויות עצמיות (מודל מסדר 0):

במקרה זה נגדיר  $p_i = \frac{v_i}{m}$  כאשר  $v_i$  הוא מספר המופעים של  $s_i$  בהודעה ו*m* זה גודל הטקסט. מה מכיל ה*prelude*?

- 8 ביטים של גודל האלפבית (מספר בין 1 ל-256 ולכן במקרה הגרוע נצטרך 8 ביטים).
- $n * 8$  ביטים עבור התיאור של כל אחד מהתווים ( $n = 25$  בדוגמה שלנו).
- $n * 4$  ביטים עבור תיאור ההסתברויות לכל תו.

### מודל מסדר 1:

ההסתברות לכל תו היא בהקשר של התו הקודם. למשל: ההסתברות שנראה i אחרי n היא  $5/8$ .

### קוד Unary (סוג של קוד סטטי):

נייצג את המספר x ע"י  $x - 1$  ביטים של "1" וביט של "0" בסוף. יתרון: קוד אינסופי (ללא דרישה שגודל האלפבית יהיה ידוע מראש).

Decimal	Unary	Binary
1	0	01
2	10	10
3	110	11
4	1110	100

### קוד בינארי (סוג של קוד סטטי):

קידוד בינארי פשוט - כל תו מקבל מילת קוד באורך  $\lceil \log_2 n \rceil$  ביטים (כאשר n זה מספר המילים).  
קידוד בינארי מינימלי – עבור אלפבית בגודל n, קוד מינימלי מכיל  $n - 2^{\lceil \log_2 n \rceil}$  מילות קוד המורכבות מ- $\lceil \log_2 n \rceil$  ביטים, ושאר  $2^{\lceil \log_2 n \rceil} - 2n$  מילות הקוד יהיו מורכבות מ- $\lceil \log_2 n \rceil$  ביטים.  
 לדוגמה: עבור  $n=5, S=[1,2,3,4,5]$  ולכן נקבל 3 מילות קוד באורך 2 ו2 מילות קוד באורך 3:  $C = [00,01,10,110,111]$ .

### קודי Elias – $C_\gamma$ :

מילת הקוד של x היא באורך  $O(\log(x))$  ביטים.

הקידוד מורכב מ2 חלקים:

- חלק ראשון – קידוד Unary עבור מספר הביטים של x בקידוד בינארי:  $Unary(\text{num of bits in binary representation of } x)$ .
- חלק שני – הקידוד הבינארי של x ללא ה1 המוביל.

Decimal	Binary	$C_\gamma$		
		First part	Second part	Total
1	1	0	-	0
2	10	10	0	100
3	11	10	1	101
4	100	110	00	11000
5	101	110	01	11001

קודי  $C_\delta$  – Elias

הקידוד מורכב מ-2 חלקים:

- חלק ראשון – קידוד  $C_\gamma$  עבור מספר הביטים של  $x$  בקידוד בינארי:  $Unary(\text{num of bits in binary representation of } x)$ .
- חלק שני – הקידוד הבינארי של  $x$  ללא ה-1 המוביל.

Decimal	Binary	$C_\delta$		
		First part	Second part	Total
1	1	0	-	0
2	10	100	0	1000
3	11	100	1	1001
4	100	101	00	10100
5	101	101	01	10101

הערה:  $C_\gamma$  ארוך יותר מקידוד  $Unary$  רק עבור  $x=2,4$ .

$C_\delta$  ארוך יותר מ- $C_\gamma$  רק עבור  $x=2,3,8,\dots,15$ .

עבור ערכים גבוהים יותר קודי  $Alias$  קצרים יותר אקספוננציאלית מקידוד  $Unary$ .

**Golomb Code**

קוד עם "דלי" בגודל קבוע  $b$ .

דוגמה לקידוד: נרצה לקודד את  $x = 8$ , כאשר  $b = 5$ .

כלומר, בכל ריישא של הקוד  $Unary$  יש 5 מילות קוד.

ו- $x = 8$ , כלומר נרצה לקודד את מילת הקוד השמינית.

$$q = (8 - 1) \div 5 = 1$$

$$r = 8 - 1 * 5 = 3$$

$$Unary(1 + 1) = 10$$

$$MBE(3,5) = 10 \text{ – כלומר מילת הקוד השלישית באלפבית בגודל 5}$$

$$(0,01,10,110,111)$$

$$C_8 = 1010$$

דוגמה לפיענוח:

$$q = Unary_{decode}(10) - 1 = 2 - 1 = 1$$

המילה (עד 0 מפריד) ויפענח.

$$r = MBD(10, 5) = 3 \text{ : שאר המילה היא } 10, \text{ וזוהי המילה השלישית באלפבית מגודל 5.}$$

$$return r + qb = 3 + 1 * 5 = 8$$

```
Golomb_encode(x,b){
    q←(x-1) div b;
    r←x-q·b;
    Unary_encode(q+1);
    Minimal_binary_encode(r,b);
}

Golomb_decode(b){
    q←Unary_decode()-1;
    r←Minimal_binary_decode(b);
    return r+q·b;
}
```

### :Rice Code

מקרה פרטי של *Golomb*, כל דלי הוא בגודל חזקה של 2:  $b = 2^k$ . במקרה זה, החלק השני של הקידוד (*MBE*) מתלכד עם קידוד בינארי פשוט. בחלק הראשון נעשה *shift right* ל- $(x - 1)$  בא ביטים ונוסיף 1 (פעולת חיבור) ואז מקודדים ב-*Unary*. דוגמה: לפי *Golomb* כאשר  $x = 8, k = 4 = 2^2$ .

$$\begin{aligned} q &= (8 - 1) \div 4 = 1 \\ r &= 8 - 1 * 4 = 4 \\ \text{Unary}(1 + 1) &= 10 \\ \text{MBE}(4, 4) &= 11 \\ X_{\text{Golomb}-4} &= 1011 \end{aligned}$$

לפי *Rice* כאשר  $x = 8, k = 4 = 2^2$ .

עבור החלק הראשון:  $x = 8 \xrightarrow{x-1} 7 \xrightarrow{(7)_2} 111 \xrightarrow{\text{shift right } k=2 \text{ bits}} 1 \xrightarrow{\text{add } 1} 2 \xrightarrow{\text{Unary}(2)} 10$   
עבור החלק השני: במקום לעשות *MBE* ניקח את  $k = 2$  הביטים הנמוכים ביותר של  $(x - 1)$ , אותם ביטים שזרקנו ב-*shift right* ← 11.  
סה"כ קיבלנו: 1011.

### :Fibonacci Code

בדיוק כמו בקידוד בינארי ניתן לקודד מספרים בפיבונאצ'י כאשר הבסיסים הם  $1, 2, 3, 5, 8, 13, \dots$ :  
למשל:  $19 = 101001$  כיוון ש  $19 = 13 * 1 + 8 * 0 + 5 * 1 + 3 * 0 + 2 * 0 + 1 * 1$ .  
התכונה המיוחדת שהקידוד מקיים הוא שאף פעם לא יהיו לנו 2 אחדות צמודות.  
לכן נוסיף 1 בהתחלה ונהפוך את הקידוד וכך נקבל קוד מייד'י.  
 $(19)_{\text{fib-code}} \rightarrow 101001 \xrightarrow{\text{add front } 1} 1101001 \xrightarrow{\text{reverse}} 1001011$   
תכונות: קוד חסר ריישות, קוד מייד'י, יש  $F_k$  מילות קוד באורך  $k - 1$  ביטים.

### :Fibonacci2 Code

ב-*Fib1* אין מילת קוד באורך 1 (הקידוד של המילה הראשונה (1) הוא 11) ולכן הציעו שינוי קטן. האלגוריתם:

- עבור כל מילת קוד ב-*Fib1* נמחק את ה-1 הימני ביותר.
- מחק את כל מילות הקוד ב-*Fib1* שמתחילות ב-0.

אלגוריתם זה:

- עבור כל מילת קוד ב-*Fib1* נמחק את ה-1 הימני ביותר.
- הוסף תחילית 10 לכל מילת קוד.
- קבע את 1 להיות מילת הקוד הראשונה.

decimal	Fib1	Delete left most 1	Add 10 prefix	Fix 1 as first cw + dropdown all cw	Fib2
1	11	1	101	1	1
2	011	01	1001	101	101
3	0011	001	10001	1001	1001
4	1011	101	10101	10001	10001
5	00011	0001	100001	10101	10101
6	10011	1001	101001	100001	100001
7	01011	0101	100101	101001	101001

תכונות: יש מילת קוד אחת באורך 1, יש  $F_{k-2}$  מילות קוד באורך  $k$ .

### :Shanon Code

מספר הביטים של כל מילת קוד  $c_i$  הוא  $\left\lceil \log_2 \frac{1}{p_i} \right\rceil$ .  
זה מבטיח לנו אורך מילת קוד ממוצע שהוא בין האנטרופיה לבין האנטרופיה+1.

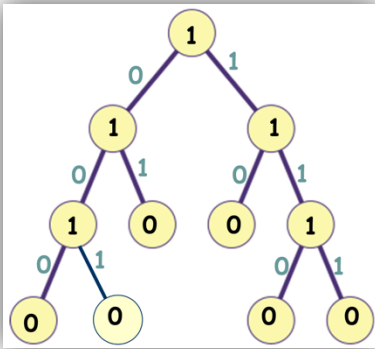
## קידוד של עץ בינארי:

כדי להעביר את הקוד למפענח לא תמיד נצטרך להעביר את כל ההסתברויות, אפשר פשוט להעביר את העץ עצמו.

עבור עץ עם  $n$  עלים, נייצג ב-1 כל צומת פנימית וב-0 עלה.

נשים לב שיש לנו  $n-1$  אפסים ו- $n$  אחדות ולכן מספר הביטים שנצטרך הוא  $2n-1$ .

לדוגמה: מייצג את העץ הבא:



Prob.	0.67	0.11	0.07	0.06	0.05	0.04
	0			1		
		0		1		
		0	1	0	1	
					0	1
Code	0	100	101	110	1110	1111

## HUFFMAN( $\Sigma$ )

```

1  $n \leftarrow |\Sigma|$ 
2  $Q \leftarrow \Sigma$ 
3 for  $i \leftarrow 1$  to  $n - 1$ 
4   do  $ALLOCATE\_NODE(z)$ 
5      $left[z] \leftarrow x \leftarrow EXTRACT\_MIN(Q)$ 
6      $right[z] \leftarrow y \leftarrow EXTRACT\_MIN(Q)$ 
7      $w[z] \leftarrow w[x] + w[y]$ 
8      $INSERT(Q, z)$ 
9 return  $EXTRACT\_MIN(Q)$ 

```

## Shanon-Fano Code

1. נסדר את התווים בסדר יורד לפי ההסתברויות.

2. כל עוד יש לנו יותר מתו אחד:

a. נחלק את הקבוצה ל-2 כך שסכום ההתפלגויות בכל קבוצה כמעט שווה.

b. נקצה 1 עבור הקבוצה הראשונה ו-0 עבור הקבוצה השנייה.

הערה: הקידוד לא תמיד יעיל.

## Huffman Encoding (זמן ריצה $O(n \log_2(n))$ ):

1. נסדר את התווים בסדר יורד לפי ההסתברויות.

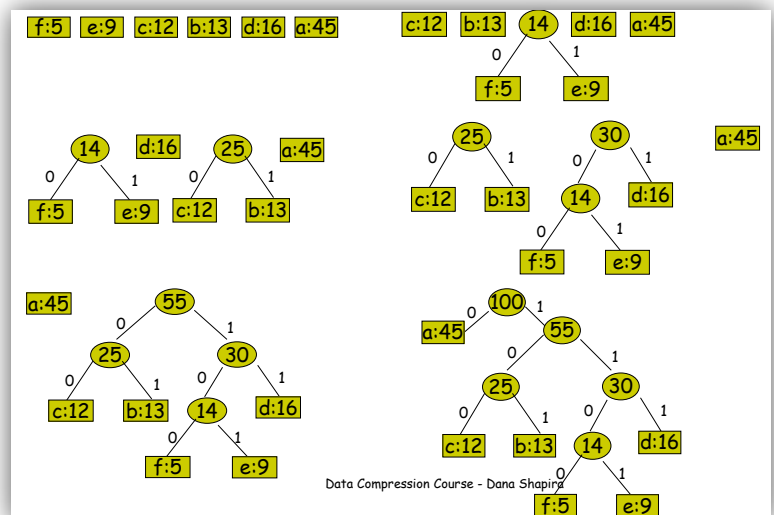
2. כל עוד יש יותר מתו אחד:

נבצע איחוד ל-2 התווים עם ההסתברויות הנמוכות ביותר תחת שורש אחד ונקצה

0 עבור הבן השמאלי ו-1 עבור הבן הימני ונוסיף את סכום ההסתברויות (שורש

העץ) לקבוצה.

דוגמה:



## תכונות:

בהינתן משקלים  $w_1, w_2, \dots, w_n$  אלגוריתם Huffman מקצה קודים באורכים  $l_1, l_2, \dots, l_n$  כך ש  $\sum_{i=1}^n w_i l_i$  הוא מינימלי.

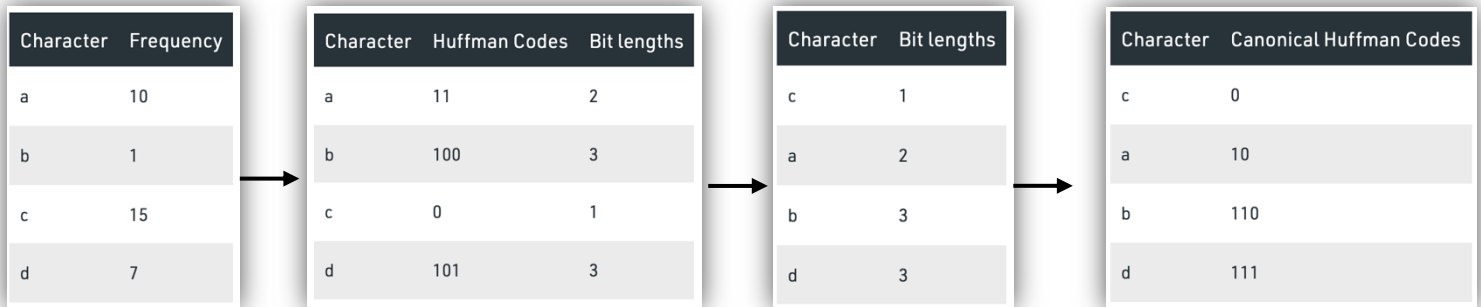
Huffman מייצר קוד שלם (אופטימלי) – קוד אופטימלי מיוצג תמיד ע"י עץ מלא (לכל צומת פנימית יש בדיוק 2 בנים).

אי שוויון קראפט מקיים עבורו שוויון  $K(C) = 1$ .

ניתן להחליף בין האפסים והאחדות בכל אחד מהבנים של הקודוקודים הפנימיים וכך אפשר לייצר  $2^{n-1}$  קידודים שונים.

### Canonical Huffman:

בקידוד האפמן קנוני, נעשה שימוש באורכי הסיביות של קודי האפמן הסטנדרטיים שנוצרו עבור כל סמל. הסמלים ממיינים תחילה לפי אורכי הסיביות שלהם בסדר עולה ולאחר מכן עבור כל אורך סיביות, הם ממיינים באופן לקסיקוגרפי. הסמל הראשון מקבל קוד המכיל את כל האפסים ובאורך זהה לזה של אורך הביט המקורי. עבור הסמלים הבאים, אם לסמל יש אורך סיביות השווה לזה של הסמל הקודם, הקוד של הסמל הקודם מוגדל באחד ומוקצה לסמל הנוכחי. אחרת, אם לסמל יש אורך סיביות גדול מזה של הסמל הקודם, לאחר הגדלה של הקוד של הסמל הקודם יתווספו אפסים עד שהאורך ישתווה לאורך הסיביות של הסמל הנוכחי והקוד יוקצה לסמל הנוכחי. תהליך זה נמשך עבור שאר הסמלים.



### הסבר לדוגמה:

- בשלב הראשון קיבלנו את התווים עם ההסתברויות.
- בשלב השני הפעלנו אלגוריתם Huffman כדי לקבל את אורכי הקידודים לכל תו.
- בשלב השלישי נסדר את התווים בסדר עולה לפי אורכי הקידודים שקיבלנו מHuffman.
- בשלב הרביעי תחיל להקצות קידודים:
- c הוא התו הראשון, צריך להקצות לו קידוד באורך 1 – לפי האלגוריתם ניתן לו אפסים בכל הביטים, הקידוד הוא בגודל 1 ולכן נקצה 0.
  - נעבור ל a, a לא באותו אורך כמו c ולכן נוסיף 1 ( $0 + 1 = 1$ ) ונרפד באפסים מימין עד שנגיע לאורך הרצוי, הקידוד הוא בגודל 2 ולכן נקצה 10.
  - נעבור ל b, b הוא לא באותו אורך כמו a ולכן נוסיף 1 ( $10 + 1 = 11$ ) ונרפד באפסים מימין עד שנגיע לאורך הרצוי, הקידוד הוא בגודל 3 ולכן נקצה 110.
  - נעבור ל d, d הוא באותו אורך כמו b, כל מה שנצטרך לעשות הוא להוסיף 1 ( $110 + 1 = 111$ ) ולכן נקצה 111.

### D-ary Huffman Code:

במקום Huffman בינארי נבצע Huffman עם D בנים.

בדיוק אותו אלגוריתם רק שבכל פעם נאחד את ה D תווים עם ההתפלגויות הכי נמוכות.

הבעיה: אם מספר התווים הוא לא כפולה של D, לשורש יהיו פחות מ D בנים ולכן בקוד שנקבל יהיו ביטים מיותרים.

פתרון: נוסיף x קודקודים עם הסתברות 0 לרשימה ההתחלתית כך שיתקיים  $x + n = 1 + (D - 1)k$ .

## Dynamic (adaptive) Huffman

נבנה את העץ בריצה אחת, כלומר בלי לדעת מראש את ההסתברויות.

דוגמה: המחזורית היא "aardvark".

- נתחיל מקודקוד NYT.
- ראינו  $a$  – נוסיף אותו בתור אח של NYT ובבנה שורש מעליהם.
- ראינו  $a$  נוסף – נעדכן את ההסתברות של  $a$  מ 1 ל 2 ולכן נעדכן את השורש להיות  $0+2$ .
- ראינו  $r$  – נוסיף אותו בתור אח של NYT ובבנה שורש מעליהם ( $1+0=1$ ), השורש החדש הוא האח החדש של  $a$ . נעדכן את השורש להיות  $1+2$ .
- ראינו  $d$  – נוסיף אותו בתור אח של NYT ובבנה שורש מעליהם ( $1+0=1$ ), השורש החדש הוא האח החדש של  $r$ . נעדכן את האבא של  $r$  ושל השורש החדש להיות  $1+1$ , הוא האח של  $a$  ונעדכן את השורש הראשי.
- ראינו  $v$  – נוסיף אותו בתור אח של NYT ובבנה שורש מעליהם ( $1+0=1$ ), השורש החדש הוא האח החדש של  $d$ . נעדכן את כל העץ כמו קודם.
- נשים לב שעד כה תמיד הבנים הימניים היו גדולים או שווים מהבנים השמאליים. ברגע שעקרון זה מופר נצטרך לתקן את העץ.
- נבצע החלפה בין כל האחים שהעיקרון הזה מופר אצלם.

## תכונות אחים:

1.  $w(v) = w(right(v)) + w(left(v))$ .
2. הקודקודים יכולים להיות ממוספרים בסדר עולה לפי המשקלים כך ש- $2j-1 \leq 2j$  הם אחים. עץ שמקיים את תכונת האחים אמ"מ הוא יכול להיות עץ Huffman.

## Sekelton Tree

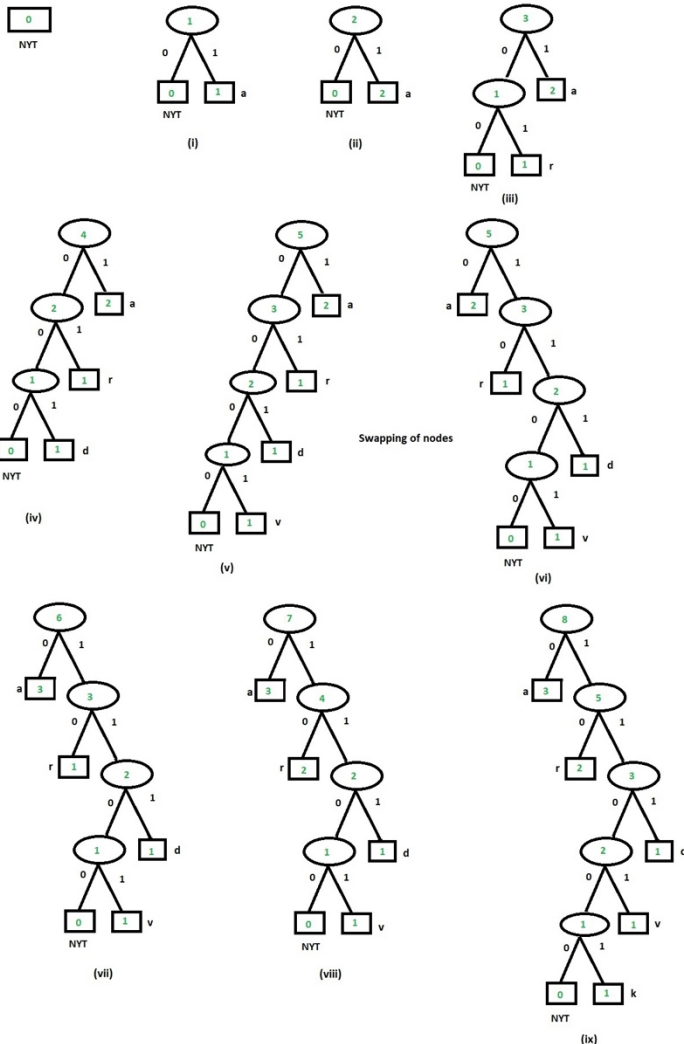
מוטבציה – לא צריך לשמור את כל העץ (חוסך מקום).

**צריך להשלים**

## Arithmetic Code

השוואה בין קוד אריתמטי לקוד האפמן: קוד האפמן:

- מחליף תו למילת קוד.
- צריך לקבל כקלט את ההסתברויות למופעים של התווים.
- קשה לו להסתגל לסטטיסטיקות משתנות.
- צריך לאחסן את טבלת מילות הקוד.
- מילת הקוד המינימלית היא באורך ביט 1.
- קוד אריתמטי: לוקח לו יותר זמן לדחוס ולפרק ולכן על אף היתרונות שלו הזמן שמושקע בדחיסה לא שווה את זה.
- מחליף את כל המחזורות במספר עשרוני יחיד (floating-point number).
- לא צריך לקבל את ההסתברויות.
- מסתגל בקלות.
- לא צריך לאחסן ולשלוח את טבלת הקידודים.
- אורכי מילות קוד שבריים.
- מגיע לאנטרופיה.



## קידוד אריתמטי:

```
low ← 0.0
high ← 1.0
while input symbols remain{
    range ← high - low
    Get symbol
    high ← low + high_bound(symbol)*range
    low ← low + low_bound(symbol)*range
}
Output any value in [low, high)
```

1. נגדיר  $low=0, high=1$ .
2. כל עוד נותרו תווים לקרוא:
  - a.  $range=high-low$
  - b. קריאת התו הבא.
  - c. עדכון high לפי החסם עליון של התו.
  - d. עדכון low לפי חסם תחתון של התו.
3. פלוט מספר כלשהו בטווח  $[low, high)$ .

M[i]	Low	High	Range
-	0.0000	1.0000	1.0000
A	0.0000	0.67	0.67
B	0.4489	0.5226	0.0737
A	0.4489	0.498279	0.049379
A	0.4489	0.48198393	0.03308393
A	0.4489	0.47106623	0.02216623
E	0.46907127	0.47017958	0.00110831
A	0.46907127	0.46981384	0.00074257
A	0.46907127	0.46956879	0.00049752
B	0.46940461	0.46945934	0.00005473
A	0.46940461	0.46944128	0.00003667

S <sub>i</sub>	Low bound	High bound
A	0.0	0.67
B	0.67	0.78
C	0.78	0.85
D	0.85	0.91
E	0.91	0.96
F	0.96	1.0

דוגמה: קלט – "ABAAAEAAABA".

לבסוף נרצה לפלוט מספר בטווח  $[0.46940461, 0.46944128)$ .  
 נעביר את הlow ואת הhigh לייצוג בינארי:  
 $low = 0.0111100000101010111010$   
 $high = 0.0111100000101101010011$   
 נחזיר את המספר הבינארי הקטן ביותר שנמצא ביניהם:  
 $return 0.0111100000101100$

```
encoded ← Get (encoded number)
do{
    Find symbol whose range contains encoded
    Output the symbol
    range ← high(symbol) - low(symbol)
    encoded ← (encoded - low(symbol)) / range
}until (EOF)
```

## פיענוח קידוד אריתמטי:

- קלט: מספר עשרוני.
1. נגדיר encoded ששווה למספר שקיבלנו מהמקודד.
  2. כל עוד לא הגענו לEOF:
    - a. מצא את התו s שנמצא בטווח של המספר.
    - b. פלוט את s.
    - c. עדכן את הrange להיות  $range = high(s) - low(s)$ .
    - d. עדכן את הencoded להיות  $encoded = \frac{encoded - low(s)}{range}$ .

## דוגמה:

בהינתן טבלת הlow והhigh של התווים נפענח את המספר העשרוני 0.109375.

- B נמצא בטווח ולכן נפלוט B.  
 $range = 0.25 - 0 = 0.25$   
 $encoded = \frac{0.109375 - 0}{0.25} = 0.4375$
- A נמצא בטווח ולכן נפלוט A.  
 מעדכן את הencoded וכן הלאה...  
 $encoded = 0$  נעצור כאשר

E = .109375  
 between 0.0, 0.25; output 'B'  
 $E = (.109375 - 0.0) / 0.25 = .4375$   
 .4375 between 0.25, 0.5; output 'I'  
 $E = (.4375 - 0.25) / 0.25 = 0.75$   
 0.75 between 0.5, 1.0; output 'L'  
 $E = (0.75 - 0.5) / 0.5 = 0.5$   
 0.5 between 0.5, 1.0; output 'L'  
 $E = (0.5 - 0.5) / 0.5 = 0.0 \rightarrow STOP$

Symbol	Low	High
B	0	0.25
I	0.25	0.50
L	0.50	1.00



דוגמה:

## האלגוריתם משתמש בנוסחה:

$$p(x) = \frac{N(x) + 1}{\sum_{i=1}^n N(s_i) + |\Sigma|}$$

- $$\begin{aligned} p(a) &= \frac{N(a)+1}{N(a)+N(b)+N(c)+3} \\ p(b) &= \frac{N(b)+1}{N(a)+N(b)+N(c)+3} \\ p(c) &= \frac{N(c)+1}{N(a)+N(b)+N(c)+3} \end{aligned}$$

- מתחילים מאינטרוול שלם  $[0,1)$ .
- כאשר  $N(a) = N(b) = N(c) = 0$  ולכן  $p(a) = p(b) = p(c) = 1/3$  נקודת את התו  $b$ :
- $b$  מוגדר בקטע  $[1/3, 2/3]$  ולכן נעדכן את האינטרוול להיות  $[0.33, 0.66]$ .
- עתה,  $N(a) = 0, N(b) = 1, N(c) = 0$  ולכן  $p(a) = \frac{1}{4}, p(b) = \frac{1}{2}, p(c) = \frac{1}{4}$  נקודת את התו  $c$ :
- $c$  מוגדר ברבע האחרון של האינטרוול.
- אם נחלק את הקטע  $\left[\frac{4}{12}, \frac{8}{12}\right)$  לרבעים נראה שם מוגדר ברבע הראשון שזה  $\left[\frac{4}{12}, \frac{5}{12}\right)$ .
- $[0.33, 0.416)$

$b$  מוגדר באופן זהה בחצי האמצעי שזה  $[0.416, 0.583]$ .

$c$  מוגדר באופן זהה ברבע האחרון של האינטרוול זה  $[0.583, 0.667]$

לכן נעדכן את האינטרוול להיות  $[0.583, 0.667]$ .

עתה,  $N(a) = 0, N(b) = 1, N(c) = 1$  ולכן  $p(a) = \frac{1}{5}, p(b) = \frac{2}{5}, p(c) = \frac{2}{5}$

- נמשך את שאר השלבים באופן זהה עד שנסיים לקודד את המחרוזת ולבסוף נפלוט מספר שנמצא בטווח של האינטרוול הסופי.

1. דיוק – כיוון שעובדים עם אינטרוול יכול להיות שהמספר שנקודד יהיה מאוד ארוך וכידוע מחשבים לא מצליחים לדייק מספרים גדולים.

2. איטי יותר מ-Huffman.
3. לא תומך בגישה אקראית.

במידה ואיבדנו ביטים נרצה לדעת האם המפענח יצליח לחזור בשלב מסויים לתווים הנכונים (כלומר יצליח לפענח החל משלב מסויים למרות איבוד הביטים).

כאשר יש קידוד למילה מסויימת שהוא סייפא של מילת קוד אחרת שנמצאת בחלק שאיבד את הביטים, נצליח להסתנכרן ברגע שתופיע מילת הקוד שהיא סייפא.

בדוגמה  $E$  היא סייפא של  $C$  ולכן הסנכרון קורה החל מקריאת  $E$ .  
אם נשתמש ב  $\text{Fixed length codes}$  או ב  $\text{Affix codes}$  לא יהיה  
סנכרון.

A 00  
 B 010  
 C 011  
 D 10  
 E 11

0 1 0 0 0 0 1 1 1 1 0 0 1 0  
 A B C D E

## מילון סטטי vs מילון אדפטיבי:

גישה סטטית (Static approach):

- המילון נבנה לפני הקידוד.
- נדרש ידע מוקדם על המקור.
- גישה אדפטיבית (Adaptive approach): LZSS
- המילון נבנה דינמית תוך כדי הקידוד.
- גישה סטטיסטית (Statistical approach):
- ננסה לחזות את התו הבא.
- גישה ההחלפה (Substitutional approach):

- החלפת בלוקים של טקסט בהפניות למופעים קודמים של טקסט זהה.

## Dictionary

$c_i$	symbol
10	a
1111110	b
111110	c
0	aa
11111110	aaaa
1111111	ab
110	baa
1110	bccb
11110	bccba <sup>3</sup>

בהינתן מילון סטטי ומחרוזת, איך נשבור אותה לבלוקים שנקודד?

נניח שקיבלנו מחרוזת "aaabccbaaaa".

- גישה חמדנית (Greedy): מסתכלת על המחרוזת מההתחלה ותיקח בכל פעם את המחרוזת הארוכה ביותר שניתן לתפוס.  
בדוגמה: aa-ab-c-c-b-aa-aa (30 bits).
- גישה longest fragment: מסתכלת על כל המחרוזת ומוצאת את המחרוזת הארוכה ביותר שניתן לתפוס, מסתכלת שוב על המחרוזת שנותרה ומחפשת את תת המחרוזת הארוכה ביותר שניתן לתפוס וכן הלאה...  
בדוגמה: aa-a-bccb-a-a-a (11 bits).
- גישה Min-words: נרצה למצוא את מספר המילים המינימליות שניתן לפרק את המחרוזת.  
בדוגמה: aa-a-bccb-aaaa (15 bits).
- אופטימלי: aa-a-bccb-aa-aa (9 bits).

## LZ77 Compression

נחליף ביטוי בטקסט למצביע למילון כדי להשיג את הדחיסה.

Window	Look Ahead Buffer
--------	-------------------

חלון הטקסט (window): הקוד שקודדנו עד כה.

look ahead buffer: התווים שותר לקודד.

הערה: המילון כאן הוא מרומז, זה למעשה החלון.

פלט האלגוריתם יהיה שלשות כאשר כל שלשה מורכבת מ:

- offset – כמה אחורה צריך ללכת.
- length – אורך הביטוי שיש להעתיק.
- הביטוי עצמו.

## אלגוריתם קידוד:

```

1. p ← 1 // The next character to be coded
2. while there is text remaining to be coded{
  1. search for the longest match for S[p...]
    in S[p-W...p-1] suppose that the match occurs
    at position m with length l
  2. Output the triple (p-m, l, S[p+l])
  3. p ← p+l+1
}
```

1. נגדיר מצביע p שיצביע למקום הראשון.

2. כל עוד יש עוד תווים בטקסט שצריך לקודד:

- חפש את המחרוזת הארוכה ביותר בlook ahead שקיימת בwindow.
- נוציא את השלשה המתאימה.
- נקדם את הפוינטר שיצביע על המקום הבא.

בדוגמה נרצה לקודד את המחרוזת הנ"ל.

כל שלשה בקידוד אומרת כמה ללכת אחורה? כמה להעתיק? ואיזה תו לשים בסוף?

למשל השלשה (21,11,v') אומרת לך אחורה 21 תווים תעתיק 11 תווים ותוסיף בסוף v.

בפועל לא נשלח למפענח את השלשות בצורה כזו, נשלח 6 ביטים עבור offset (הגדול ביותר הוא 21 וזה דורש 6 ביטים), 4 ביטים עבור length (הגדול ביותר הוא 11 וזה דורש 4 ביטים) ו-8 ביטים עבור התו (ASCII) – סה"כ 18 ביטים לכל שלשה.

String:

A\_walrus\_in\_Spain\_is\_a\_walrus\_in\_vain.

Encoded String

(0,0,'A')(0,0,'\_')(0,0,'w')(0,0,'a')(0,0,'l')(0,0,'r')(0,0,'u')  
(0,0,'s')(7,1,'i')(0,0,'n')(3,1,'S')(0,0,'p')  
(11,1,'i')(6,2,'i')(12,2,'a')(21,11,'v')(20,3,'.')

## אלגוריתם פנענח:

```

1. p ← 1 // The next character to be decoded
2. For each triple (f, l, c) in the input{
  1. S[p..p+l-1] ← S[p-f..p-f+l-1]
  2. S[p+l] ← c
  3. p ← p+l+1
}

```

### Encoded String

```

(0,0,'A')(0,0,'_')(0,0,'w')(0,0,'a')(0,0,'l')(0,0,'r')(0,0,'u')
(0,0,'s')(7,1,'i')(0,0,'n')(3,1,'S')(0,0,'p')
(11,1,'i')(6,2,'i')(12,2,'a')(21,11,'v')(20,3,':')

```

### String:

A\_walrus\_in\_Spain\_is\_a\_walrus\_in\_vain.

1. נגדיר מצביע  $p$  שיצביע למקום הראשון.
2. לכל שלשה בקלט:
  - a. לך אחורה  $offset$  צעדים.
  - b. תעתיק  $length$  תווים ותוסיף את התו.
  - c. נקדם את הפוינטר שיצביע על המקום הבא.

בדוגמה המפענח מקבל את השלשות ומתחיל לפענח כל שלשה. בדוגמה כל השלשות הראשונות הן עם  $offset$  0 וגם  $length$  0 ולכן המפענח פשוט ירשום את התווים  $A\_walrus$ . לאחר מכן יפענח את השלשה  $(7,1,'i')$  יילך 7 תווים אחורה, יעתיק תו 1 ('\_') וישים 'i' בסוף.

סה"כ נקבל  $A\_walrus\_i$ . נמשיך ככה עד שייגמרו השלשות.

## חסרונות של LZ77:

- צוואר בקבוק - בעת הקידוד, עליו לבצע השוואות של מחרוזות אל מול  $look-ahead$  buffer עבור כל מיקום בחלון הטקסט.
- שימוש בשלשה כאשר אין התאמה – כשלא נמצא התאמה ב  $window$  נכתוב  $(0,0,'c')$  וזה תוספת מיותרת של ביטים. נניח שהחלון בגודל 4096 תווים וה  $look-ahead$  buffer בגודל 15 תווים. כדי לקודד את השלשה הנ"ל נצטרך 12 סיביות ל  $offset$  ( $2^{12} = 4096$ ), 4 סיביות ל  $length$  ו 8 תווים לתו עצמו. למעשה אנחנו מעבירים 24 סיביות במקום להעביר רק 8.

## LZSS:

גרסה משופרת ל LZ77.

- מאפשר לערבב מצביעים (שלשות) ותווים בקידוד.
- חיפוש המחרוזות הזהות נעשה באמצעות בניית עץ מאוזן.
- כדי שהמפענח ידע לפענח את הקידוד נוסיף ביט זיהוי לפני קידוד כל תו. במידה והקידוד הוא של תו נוסיף בהתחלה 0 ובמקרה וזה מצביע (שלשה) נוסיף בהתחלה 1.
- לכן, במקרה שהחלון שלנו בגודל 4096 תווים וה  $look-ahead$  buffer בגודל 16 תווים, נצטרך 9 ביטים לקודד תו  $0+character$  ו 17 ביטים לשלשה  $1+offset(12bit)+length(4bit)$ .

## דוגמה:

נשים לב שלא משתלם לנו לכתוב שלשה עבור תו בודד אך לעיתים גם לא משתלם לנו לכתוב שלשה בשביל להעתיק תו או 2 כי הקידוד ינו להיות יותר ביטים מאשר אם נקודד את התווים ב ASCII ולא נשלח הצבעה.

לדוגמה לא משתלם לנו להעתיק את ה' השני ע"י מצביע לראשון כי זה ייקח יותר ביטים מאשר לכתוב את התו עצמו. במידה וגודל החלון וה  $look-ahead$  הם כמו למעלה, העתקה של 2 ביטים תיקח לנו 18 ביטים אם נכתוב אותם כתווים ב ASCII ו 17 ביטים עם מצביעים – לכן נבנה שלשה עבור העתקה של 2 תווים ומעלה.

## LZS:

חידוש ל LZ77.

- נגדיר מספר אופציות ל  $offset$ , לפעמים נרצה ללכת הרבה תווים אחורה ולפעמים נרצה ללכת מעט תווים אחורה.
- אין סיבה ששתי האופציות ייקחו אותם כמות תווים ולכן נוסיף  $flag$  שידע כמה הוא צריך לקרוא מה  $offset$ .
- אם אני רוצה ללכת מעט תווים אחורה נגביל את החלון ל 7 ביטים ונוסיף  $flag$  של 10 כדי שהמפענח ידע כמה תווים לקרוא.
- אם נרצה ללכת יותר תווים אחורה ניתן חלון בגודל 11 ביטים ונוסיף  $flag$  של 11.

ניתן להוסיף יותר גדלים (ולכן  $flag$ ים יותר ארוכים) ולשחק עם גודלי ה  $offset$ .

### String:

A\_walrus\_in\_Spain\_is\_a\_walrus\_in\_vain

### Encoded String

A\_walrus\_in\_Spa(6,3)is\_a(21,11)v(20,3)

### LZS stacker

- Char 0 8 bit
- Offset 11 7 bit 128
- 10 11 bit 2048
- Length

בLZ77 מילון הביטויים הוגדר על ידי חלון קבוע של טקסט שנראה בעבר – window. בLZ78 מילון הוא רשימה בלתי מוגבלת של ביטויים שנראו בעבר. כל אסימון ב LZ78 מורכב מקוד שבוחר ביטוי נתון ותו בודד שאחריו הביטוי. שלא כמו LZ77, אורך הביטוי אינו עובר מאחר שהמפענח יודע זאת.

#### אלגוריתם קידוד:

1. מתחילים ממילון ריק – בכניסה 0 יש את המילה הריקה  $\epsilon$ .
2. מצא את ההתאמה הארוכה ביותר בין look-ahead לבין הכניסה הארוכה ביותר שקיימת במילון שפותחת בתו הראשון בlook-ahead.
3. פלוט את הזוג  $(id, c)$  כאשר  $c = next\ char$ .
4. אם  $c = EOF$  עצור.
5. אחרת, הוסף למילון  $(\max id + 1, T(id) * c)$ .
6. חזור ל2.

#### • T = badadadabaab

index	Phrase	Encoding	# of bits
0	$\epsilon$		
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

דוגמה:  $\Sigma = \{a, b, c, d\} \xrightarrow{\text{binary}} \{00, 01, 10, 11\}$

- מתחילים ממילון ריק עם כניסה 0 שמכילה את המילה הריקה.
- ההתאמה הארוכה ביותר במילון היא המילה הריקה ולכן נפלוט את הזוג  $(0, b)$ . ונוסיף למילון בכניסה ה1 את b.
- נקפוץ לאיטרציה הרביעית: עד כה קידדנו את bad. מסתכלים בlook-ahead והמחרוזת הארוכה ביותר במילון שפותחת את המחרוזת ad היא a (כניסה 2). לכן, נפלוט את הזוג  $(2, d)$  ונוסיף למילון  $(4, ad)$ .

#### כמה ביטים זה לוקח?

תחילה נציין שלא מעבירים את המילון כיוון שהמפענח ידע לבנות אותו לבד.

המפענח יודע שהכניסה ה0 היא המילה הריקה ולכן לא נצטרך לקודד אותו (ייקח לנו 0 ביטים).

בכניסה ה1 המפענח יודע שהכניסה שהאיבר הראשון בזוג הוא הכניסה ה0 ולכן לא נצטרך לקודד אותה, נצטרך לקודד רק את התו b, יש 4 תווים באלפבית ולכן הקידוד של b יעלה לנו 2 ביטים.

כרגע במילון של המקודד (ושל המפענח) יש רק 2 כניסות ולכן נוכל לציין את מספר הכניסה ע"י ביט בודד + 2 ביטים עבור התו עצמו בבינארי – סה"כ 3 ביטים.

עכשיו יש במילון 3 כניסות ולכן קידוד של כניסה יעלה לנו 2 ביטים ומסיבה זו הקידוד של הזוג הבא יעלה 4 ביטים, 2 עבור הכניסה 21 עבור התו עצמו וכן הלאה.

$$\left\{ \begin{array}{l} \text{encoding :} \\ \underbrace{(0, b)}_b \underbrace{(0, a)}_a \underbrace{(0, d)}_d \underbrace{(2, d)}_{ad} \underbrace{(4, a)}_{ada} \underbrace{(1, a)}_{ba} \underbrace{(2, b)}_{ab} \\ \text{in bits :} \\ \underbrace{(\epsilon^*, 01)}_b \underbrace{(0, 00)}_a \underbrace{(00, 11)}_d \underbrace{(10, 11)}_{ad} \underbrace{(100, 00)}_{ada} \underbrace{(001, 00)}_{ba} \underbrace{(010, 01)}_{ab} \end{array} \right\}$$

בסופו של דבר נקבל את הקידוד שמופיע משמאל. ניתן לראות את הקידוד בביטים.

## LZW

### אלגוריתם קידוד:

```

1. Dictionary ← single Characters
2. w ← first char of input
3. repeat{
  1. k ← next char
  2. if(EOF)
    1. output code(w)
  3. else if (w · k) ∈ Dictionary
    1. w ← w · k
  4. else
    1. output code(w)
    2. Dictionary ← w · k
    3. w ← k
}

```

Data Compression Course - Dana Shapira

1. נאתחל את המילון עם כל התווים הבודדים שבאלפבית.
2. נאתחל את w להיות התו הראשון בטקסט.
3. כל עוד לא הגענו לסוף הטקסט:
  - a. נאתחל את k להיות התו הבא.
  - b. אם  $wk$  נמצא במילון:  $w = wk$ .
  - c. אחרת:
    - i. נפלוט את הקידוד של w.
    - ii. נכניס למילון תחת כניסה חדשה את  $wk$ .
    - iii.  $w$  יקבל את התו הבא בטקסט -  $w = k$ .

### דוגמה:

- נרצה לדחות את המחרוזת T.
- תחילה המילון ריק (כניסה 0 עם ביטוי  $\epsilon$ ).
- נאתחל את המילון באלפבית שלנו (כניסות 1-4 יקבלו את התווים a, b, o, w בהתאמה).
- נאתחל את w להיות האות הראשונה בטקסט  $w = w$ .
- ואת k להיות האות הבאה  $k = a$ .
- הביטוי  $wa$  לא מופיע במילון ולכן נפלוט 4 (הקידוד של w), ונכניס לכניסה 5 את הביטוי  $wa$ , ונגדיר את w להיות שווה לא (כלומר  $w = a$ ).
- נמשיך באותה דרך עד שנגיע לאיטרציה השביעית.
- עתה,  $w = w, k = a$ .
- הביטוי  $wa$  כן מופיע בטבלה בכניסה 5 ולכן נעדכן את  $w = wa$ .
- k מקבל את התו הבא:  $k = b$ .
- הביטוי  $wab$  לא מופיע בטבלה ולכן נפלוט את הקידוד של  $w = wa$  שהוא 5, נכניס את  $wab$  למילון לכניסה הבאה (כניסה 11) ונגדיר את w להיות שווה לא.
- נמשיך כך עד שנסיים לקרוא את כל הקובץ.
- בסופו של דבר הקידוד שנקבל הוא 41221057...

T=wabba\_wabba\_wabba\_wabba\_woo\_woo

code	phrase	w	k	∈	update (n, wk)	out
0	—	w	a	×	(5, wa)	4
1	a	a	b	×	(6, ab)	1
2	b	b	b	×	(7, bb)	2
3	o	b	a	×	(8, ba)	2
4	w	a	—	×	(9, a_)	1
		—	w	×	(10, _w)	0
		w	a	✓	—	
		wa	b	×	(11, wab)	5
		b	b	✓	—	
		bb	a	×	(12, bba)	7
		⋮				

- הביטוי  $wab$  לא מופיע בטבלה ולכן נפלוט את הקידוד של  $w = wa$  שהוא 5, נכניס את  $wab$  למילון לכניסה הבאה (כניסה 11) ונגדיר את w להיות שווה לא.
- נמשיך כך עד שנסיים לקרוא את כל הקובץ.
- בסופו של דבר הקידוד שנקבל הוא 41221057...

### אלגוריתם מפענח:

```

1. Initialize table with single character strings
2. OLD = first input code
3. output translation of OLD
4. while not end of input stream{
  1. NEW = next input code
  2. if NEW is not in the string table
    1. S = translation of OLD
    2. S = S · C
  3. else
    1. S = translation of NEW
  4. output S
  5. C = first character of S
  6. Translation(OLD) · C to the string table
  7. OLD = NEW
}

```

Data Compression Course - Dana Shapira

1. נאתחל את המילון בכל התווים הבודדים שבאלפבית.
2. נגדיר את OLD להיות התו הראשון בקלט.
3. נפלוט את התרגום של OLD.
4. כל עוד לא הגענו לסוף הטקסט:
  - a. נגדיר את NEW להיות הקידוד הבא.
  - b. אם NEW לא מופיע במילון: נגדיר את S להיות שווה לפירוש של OLD משורשר עם C.
  - c. אחרת, נגדיר את S להיות שווה לפירוש של NEW.
  - d. פלוט את S.
  - e. נגדיר את C להיות התו הראשון של S.
  - f. נכניס בכניסה הבאה במילון את הפירוש של OLD משורשר עם C.
  - g. נגדיר  $OLD = NEW$ .

## דוגמה:

נניח שהמפענח קיבל את הקידוד הבא: 0,1,3,2,4,7,0,9,10,0.

code	phrase		old	new	$\in$	s	out	c	update (n, T(old) · c)		code	phrase
0	a		0	—	✓	—	a	—	—		0	a
1	b		0	1	✓	b	b	b	$(4, T(0) \cdot b) = (3, ab)$		1	b
2	c		1	3	✓	ab	ab	a	$(4, T(1) \cdot a) = (4, ba)$		2	c
			3	2	✓	c	c	c	$(5, T(3) c) = (5, abc)$		3	ab
			2	4	✓	ba	ba	b	$(6, T(2) b) = (6, cb)$		4	ba
		⇒	4	7	×	ba · b	bab	b	$(7, T(4) b) = (7, bab)$	⇒	5	abc
			7	0	✓	a	a	a	$(8, T(7) b) = (8, baba)$		6	cb
			0	9	×	a · a	aa	a	$(9, T(0) a) = (9, aa)$		7	bab
			9	10	×	aa · a	aaa	a	$(10, aaa)$		8	baba
			10	0	✓	a	a	a	$(11, aaaa)$		9	aa
											10	aaa
											11	aaaa

deocded = a b ab c ba bab a aa aaa a

## Run-Length Codes

Runs – בלוקים של תווים שחוזרים על עצמם ברצף.

נייצג את המחזורות ע"י האורכים של תווים רצופים.

דוגמה: בהינתן המחזורות: acccbbaaabb נקודד אותה ל:  $(a, 1)(c, 3)(b, 2)(a, 3)(b, 2)$ .

## Binary Run-Length Codes

במידה והמחזורות היא בינארית, יש רק שני תווים יהיה לנו קל יותר לקודד.

דוגמה: בהינתן המחזורות 011110011100000101011 נקודד אותה ל: 1423511112.

הערה: נניח שתמיד נפתח מחזורות ב0, במידה והמחזורות מתחילה ב1 אז יש רצף של אפס אפסים בהתחלה ולכן הקידוד של התו

הראשון יהיה 0. דוגמה: בהינתן המחזורות 111001 נקודד אותה ל: 0321.

בדר"כ יהיה מוגדר מספר  $k > 0$ , שיהיה מספר הביטים המקסימלי שאפשר לקודד בו run בודד.

איך נטפל בruns גדול שהקידוד שלו דורש יותר מא ביטים, כלומר החrun גדול יותר מ  $2^k$ ?

עבור המקרים הנדירים האלה נשתמש בescape code שיהיה רצף של אחדות על כל הא ביטים וכך המפענח ידע לא להחליף צד,

כלומר, אנחנו נשארים בקידוד של אותו מספר.

עכשיו יש לנו בעיה חדשה, נניח שהתו 0 מופיע  $2^k - 1$  פעמים, איך נקודד אותו? הרי אם נקודד אותו ע"י k אחדות המפענח יחשוב

שזה escape code, כדי לפתור את זה נוסיף 0 אחרי רצף האחדות.

דוגמה: נניח ש  $k=3$  ואנו רוצים לקודד את 011110011100000101011.

נקודד את זה ל: 1423511112.

נכתוב את זה ע"י 3 סיביות לכל קידוד ונקבל את הקידוד הסופי: 001 100 010 011 101 001 001 001 010.

דוגמה נוספת: נניח ש  $k=3$  ואנו רוצים לקודד את 0000000 (7 אפסים).

הקידוד יהיה: 111 000.

דוגמה נוספת: נניח ש  $k=3$  ואנו רוצים לקודד 14 אפסים.

הקידוד יהיה: 111 111 000.

דוגמה נוספת: נניח ש  $k=3$  ואנו רוצים לקודד 16 אפסים.

הקידוד יהיה: 111 111 010.

דוגמה נוספת: נניח ש  $k=3$  ואנו רוצים לקודד  $m(2^k - 1)$  אפסים.

הקידוד יהיה: m פעמים 111 ובסוף 000.

כדי לשפר אפשר לבצע בחירה אדפטיבית של k:

• המקודד יכול לבצע תחילה מעבר על הקלט כדי לבחור את הערך של k הממזער את האורך הכולל של קידוד אורך הריצה.

• אפשר להשתמש בקודים בעלי אורך משתנה לספירות: להשתמש בקודי Huffman וכו' לספירה.

## אלגוריתם מפענח:

```

MaxCount  $\leftarrow 2^k - 1$ 
parity  $\leftarrow 0$ 
while input remains
    read k bits from the input stream to get the integer x.
    if parity = 0
        output x "0" bits
    else
        output x "1" bits
    if x < MaxCount
        parity := 1 - parity
    
```

## אלגוריתם מקודד:

```

MaxCount  $\leftarrow 2^k - 1$ ; count  $\leftarrow 0$ ; PrevBit  $\leftarrow "0"$ ;
while input remains
    CurBit  $\leftarrow$  next input bit;
    if PrevBit = CurBit and count < MaxCount
        count++
    else
        output count using k bits
        if count = MaxCount and PrevBit  $\neq$  CurBit
            output 0 using k bits
        count  $\leftarrow 1$ ;
    PrevBit  $\leftarrow$  CurBit
output count using k bits.
    
```

1. Generate a matrix of all the cyclic shifts of T
2. Sort the matrix rows in lexicographic order
3. The output of BWT:
  - The final column of the matrix
  - The number of the row corresponding to the original input

mississippi  
ississippim  
ssissippimi  
sissippimis  
issippimiss  
ssippimissi  
sippimissis  
ippimississ  
ppimississi  
pimississip  
imississipp

F L  
imississipp  
ippimississ  
issippimiss  
ississippim  
mississippi  
pimississip  
ppimississi  
sippimissis  
sissippimis  
ssippimissi  
ssissippimi

```

BWT(L, Index) {
    F  $\leftarrow$  sort(L)
    I  $\leftarrow$  Index
    for (i=0; i<n; i++){
        T[i]  $\leftarrow$  F[I];
        I  $\leftarrow$  S[I];
    }
}
    
```

## BWT:

אלגוריתם שמבצע טרנספורמציה לטקסט ומחזיר טקסט שניתן לדחוס אותו בצורה יותר טובה.

## אלגוריתם המקודד:

1. נבנה מטריצה של כל הטרנספורמציות הציקליות של T.
2. נמיין את שורות המריצה לפי סדר לקסיקוגרפי.
3. נפלוט את העמודה האחרונה במטריצה ואת מספר השורה שבה נמצא הטקסט המקורי במטריצה הממויינת.

## דוגמה:

T=mississippi

בשלב הראשון נבנה את המטריצה עם כל הטרנספורמציות הציקליות של T ונקבל את המטריצה השמאלית.  
בשלב השני נמיין את המטריצה לקסיקוגרפית ונקבל את המטריצה הימנית.  
לבסוף נפלוט את העמודה האחרונה ואת האינדקס שבו נמצא הטקסט המקורי:  
 $BWT(T) = (pssmipissii, 4)$

## תכונות המטריצה הממויינת:

1. מיון של L (העמודה האחרונה) יניב לנו את F (העמודה הראשונה).
2. תווים זהים בL מסודרים באותו סדר בF.

## אלגוריתם המפענח:

1. נגדיר את F להיות הקלט L בצורה ממויינת.
2. נגדיר את I להיות האינדקס שקיבלנו.
3. עבור  $i = 0, \dots, n$ :
  - a.  $T[i] = F[I]$
  - b.  $I = S[I]$

דוגמה לפיענוח: קלט: (pssmipissii, 4)

נגדיר את L להיות הטקסט שקיבלנו ואת F להיות L ממוריד.

נגדיר  $I = 4$  ונגדיר את T להיות וקטור בגודל n, בסוף נחזיר את T. לאחר מכן נבנה את S להיות מיופיו של כל תו ב F לתו המתאים ב L.

נתחיל מאינדקס 4, ולכן נוסף ל T את  $F[4] : T=m$ .

S ממפה את 4 ב L ל 3 ולכן נוסף ל T את  $F[3] : T=mi$ .

S ממפה את 3 ב L ל 10 ולכן נוסף ל T את  $F[10] : T=miss$ .

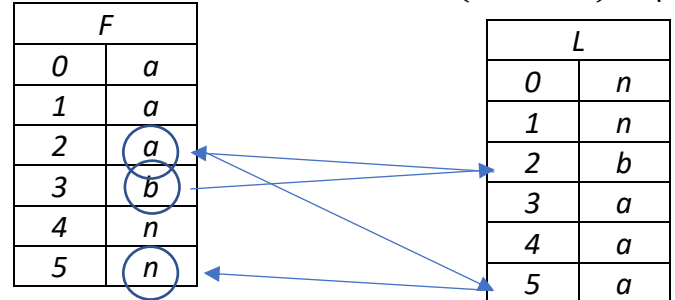
S ממפה את 10 ב L ל 8 ולכן נוסף ל T את  $F[8] : T=miss$ .

S ממפה את 8 ב L ל 2 ולכן נוסף ל T את  $F[2] : T=missi$ .

וכן הלאה.

דוגמה נוספת ללא שימוש ב S:

קלט: (nnbaaa, 3).



התחלנו מאינדקס 3 ב F, נפלוט b.

נעבור ל המתאים ב L, נמצא באינדקס 2, ומשם נחזור ל L לאינדקס הזה ונפלוט את a (נמצא באינדקס 2 ב F).

נעבור ל המתאים ב L, זהו ה a השלישי ב F ולכן נעבור ל a בשלישי ב L, נמצא באינדקס 5, ומשם נחזור ל L לאינדקס הזה ונפלוט את n (נמצא באינדקס 5 ב F) וכן הלאה...

### MoveToFront Compression:

כל תו בקלט מקודד לפי מספר התווים השונים שהופיעו מאז הופעתו האחרונה של תו זה.

מיושם באמצעות רשימה של תווים ממיינים לפי עדכניות השימוש.

הפלט מכיל מספרים קטנים רבים אם הטקסט הומוגני מקומי – BWT פולט טקסט שהוא הומוגני מקומי.

דוגמה:

קלט:  $\Sigma = \{d, e, h, l, o, r, w\}$ ,  $T = \text{helloworld}$ .

לאורך כל הקידוד נתחזק את רשימת התווים ונעדכן אותה בכל איטרציה.

קידוד התו הראשון:

$MTF - List = \{d, e, h, l, o, r, w\}$ ,  $char = 'h'$

התו 'h' מופיע במקום השני ברשימה ולכן נפלוט 2 ונבצע הזזה קדימה לתו 'h' ברשימה.

קידוד התו השני:

$MTF - List = \{h, d, e, l, o, r, w\}$ ,  $char = 'e'$

התו 'e' מופיע במקום השני ברשימה ולכן נפלוט 2 ונבצע הזזה קדימה לתו 'e' ברשימה.

קידוד התו השלישי:

$MTF - List = \{e, h, d, l, o, r, w\}$ ,  $char = 'l'$

התו 'l' מופיע במקום השלישי ברשימה ולכן נפלוט 3 ונבצע הזזה קדימה לתו 'l' ברשימה.

קידוד התו הרביעי:

$MTF - List = \{l, e, h, d, o, r, w\}$ ,  $char = 'l'$

התו 'l' מופיע במקום האפס ברשימה ולכן נפלוט 0 ונבצע הזזה קדימה לתו 'l' ברשימה – בפועל הוא כבר נמצא שם.

נמשיך כך ובסוף נקבל את הקידוד: 2230461...

### האלגוריתם הכללי לדחיסה:

קלט – מחרוזת T.

1.  $BWT(T)$  – נבצע טרנספורמציה לטקסט.

2.  $MTF$  – נבצע קידוד לטקסט שקיבלנו מה BWT.

3. נבצע דחיסה סטטיסטית לקידוד ע"י Huffman או Arithmetic Coding.