

דחיסת נתונים – סיכום

הגדרה (קוד חסר-ראשות). אף מילה אינה רישא של מילת קוד אחרת.
הגדרה (קוד ניתן-לפיענוח). כל סדרה של מילות-קוד ניתן לקודד לפענוח באופן אחד.
הגדרה (קוד שלם). כל סדרה אינסופית למחצה ניתנת לפיענוח בצורה יחידה. **קוד שלם** \Leftrightarrow עץ מלא.
הגדרה (קוד מיידי)-בסוף כל מילת קוד יודעים (המפענח יודע) שזה סופה. **מיידי** \Leftrightarrow חסר ראשות.

משפט (אי שוויון קראפט)

עבור קוד $C = \{c_1, \dots, c_n\}$ נגדיר $K(C) = \sum_{i=1}^n 2^{-|c_i|}$ אזי מתקיים

1. אם $K(C) \leq 1$ אזי C ניתן לפיענוח.

2. אם $K(C) \leq 1$ אזי קיים קוד C' חסר רישות כך ש- $\mathbb{E}(C, P) = \mathbb{E}(C', P)$.

מודלים

ישנם מספר סוגים של מודלים:

מודל מרקוב מסדר ראשון

ישנם מספר סוגי מודלים הנופלים תחת קטגוריה זאת:

- **קודים סטטיים** – שיטות קידוד סטטיות הם פשטניות למדי, הקידודים הללו ידועים לכולם (כגון קידוד הסקי). במודל זה יש שימוש מינימלי בשכיחויות של התווים, השכיחויות נתונות רק לצורך מיון בסדרה מונוטונית יורדת. יעילות הדחיסה אינה מרשימה במיוחד אולם היות והיא פשוטה אז היא מאוד שימושית והיתרון הוא מהירות הפיענוח.

○ **קידוד אונארי** – מילת הקוד $1^{i-1}0$ מייצגת את מילת הקוד ה- i , $i \in \mathbb{N}$, כלומר את c_i .

$$c_1 = 0, c_2 = 10, c_3 = 110, c_4 = 1110, \dots$$

אם מספר המילים שאנחנו מקודדים הוא סופי $n \in \mathbb{N}$, ואנחנו יודעים מראש את כל המילים אותם אנו מקודדים נוכל לקחת את המילה שמופיעה הכי מעט כלומר c_n ובמקום לקודד אותה ב- $1^{n-1}0$ נקודד אותה ב- 1^{n-1} , זאת אומרת שחסנו מילת תו אחד בכל קידוד של המילה. מצד שני, זה אינו באמת משמעותי שכן זוהי המילה הכי פחות שכיחה.

נשים לב: קוד אונארי הוא בעל יתירות מינימלית אם ההסתברויות הם חזקות של 2, כלומר $p_i = 1/2^{\ell_i}$, כאשר ℓ_i היא אורך מילת הקוד ה- i .

$$\mathbb{E}(C) = \sum_{i=1}^n p_i \cdot \ell_i = - \sum_{i=1}^n p_i \cdot \log_2 p_i = H(C) \Leftrightarrow \sum_{i=1}^n \ell_i = \sum_{i=1}^n -\log_2 p_i$$

ע"פ ההנחה ש- $p_i = 2^{-\ell_i}$ נובע כי $-\log_2 p_i = -\log_2 2^{-\ell_i} = \ell_i$, ולכן צד ימין מתקיים, לפיכך התוחלת מינימלית שכן היא האנטרופיה והאנטרופיה הינה חסם תחתון על אורך מילת הקוד הממוצעת. **הערה.** כל סידור של ההסתברויות הנ"ל יגדיר קוד בעל יתירות מינימלית.

- **קוד דיאדי** – קוד שבו ההסתברויות הם חזקות של 2 וכמות האינפורמציה שווה בדיוק לגודל מילת-הקוד. בעצם ראינו שקוד אונארי הוא אופטימלי כאשר הוא קוד דיאדי.

$MBE(3,5)$, כלומר מילת הקוד השלישית בקידוד בינארי מינימלי בעל 5 מילות קוד. היות ומילת הקוד

החמישית היא בדלי הראשון צריך מינוס אחד ב- q , וכיו"ב.

- **פיענוח גולומב – מקבלים מילה לפיענוח α ואת גודל הדלי b , נשים לב שהאינדקס הוא**

$$MBD(\alpha, b) + (UnaryD(\alpha) - 1) \cdot b$$

- **קוד רייס – מקרה פרטי של גולומב כאשר הדליים הם חזקות של 2, ואז העץ בינארי המינימלי הוא גם העץ בינארי הפשוט וזה יקל עלינו קצת. לדוגמה ניקח $b = 2^2$ ואת האינדקס 8, אז המילה בדלי השני, במיקום 4 כלומר הקידוד 1011. מסיבות כאלו ואחרות יותר יעיל לבצע את החישוב באופן אחר שמתואר במחברת של תום. במקרה לעיל מתקבל**

$$8 \rightarrow \langle 8 - 1 \rangle_2 \rightarrow 111 \rightarrow 1 \rightarrow \langle 1 \rangle_{10} + 1 = 2 \rightarrow \langle 2 \rangle_1 = 10 \rightarrow 1011$$

כאשר בסוף מוסיפים את האיברים שמחקנו, דוגמה נוספת:

$$i = 30, b = 2^3 \rightarrow \langle 30 - 1 \rangle_2 \rightarrow 11101 \rightarrow 11 \rightarrow \langle 11 \rangle_{10} + 1 = 3 + 1 = 4 \rightarrow \langle 4 \rangle_1 = 1110 \rightarrow 1110101$$

הגדרה (קוד בעל יתירות מינימלית)

נאמר שקוד C הוא בעל יתירות מינימלית אם לכל קוד אחר C' מתקיים $\mathbb{E}(C) \leq \mathbb{E}(C')$.

- **קוד שנון – הקוד מקיים $\ell_i = -\lceil \log_2 p_i \rceil$, וניתן להוכיח שאורך מילת הקוד הממוצעת היא לכל היותר אחד ועוד האנטרופיה. מצד אחד זה די טוב, מצד שני במקרה הגרוע כל מילת קוד ארוכה בסיבית מהחסם התחתון.**

למטה נסמן ב- $Binary^*$ את הפונקציה שמחזירה את הקידוד הבינארי ללא ה-MSB. וב- $\langle n \rangle_2$ את הקידוד הבינארי של n .

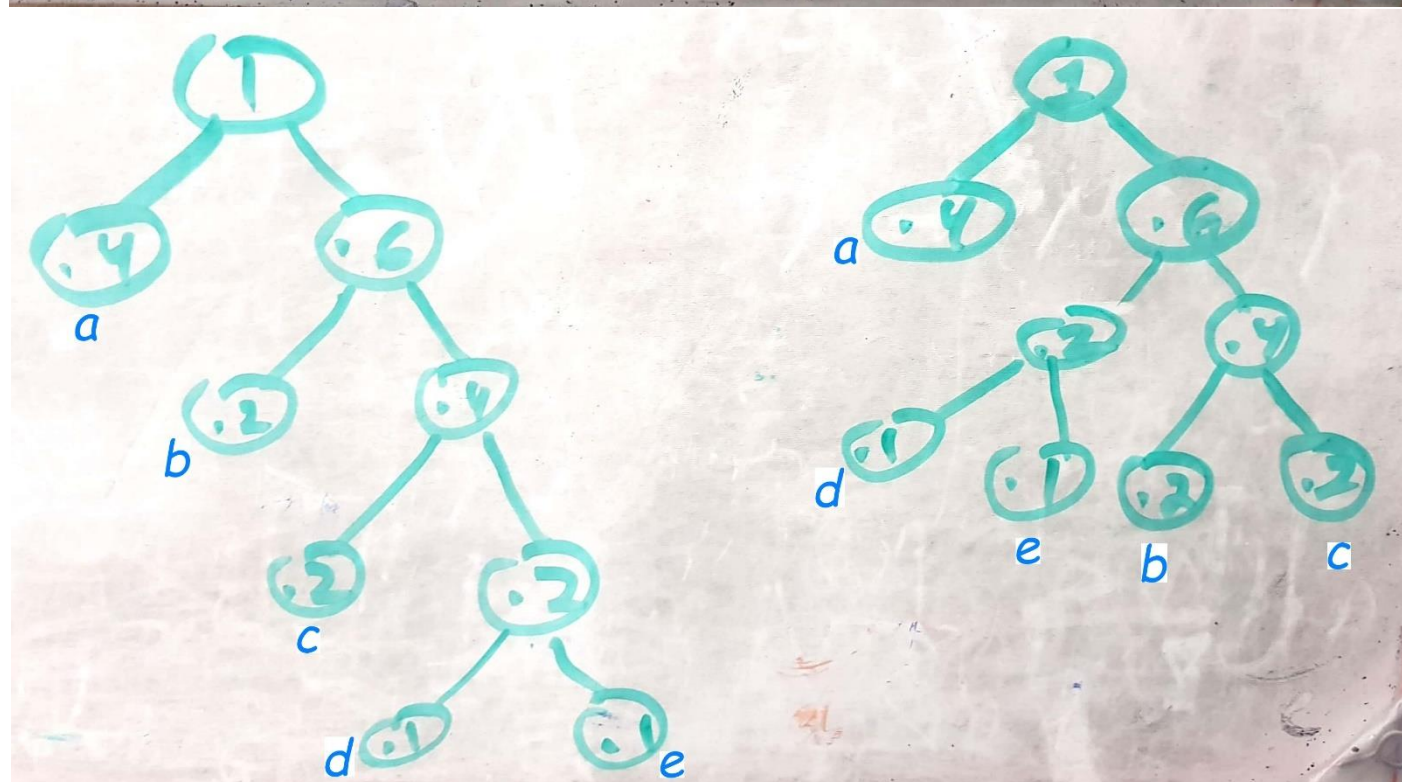
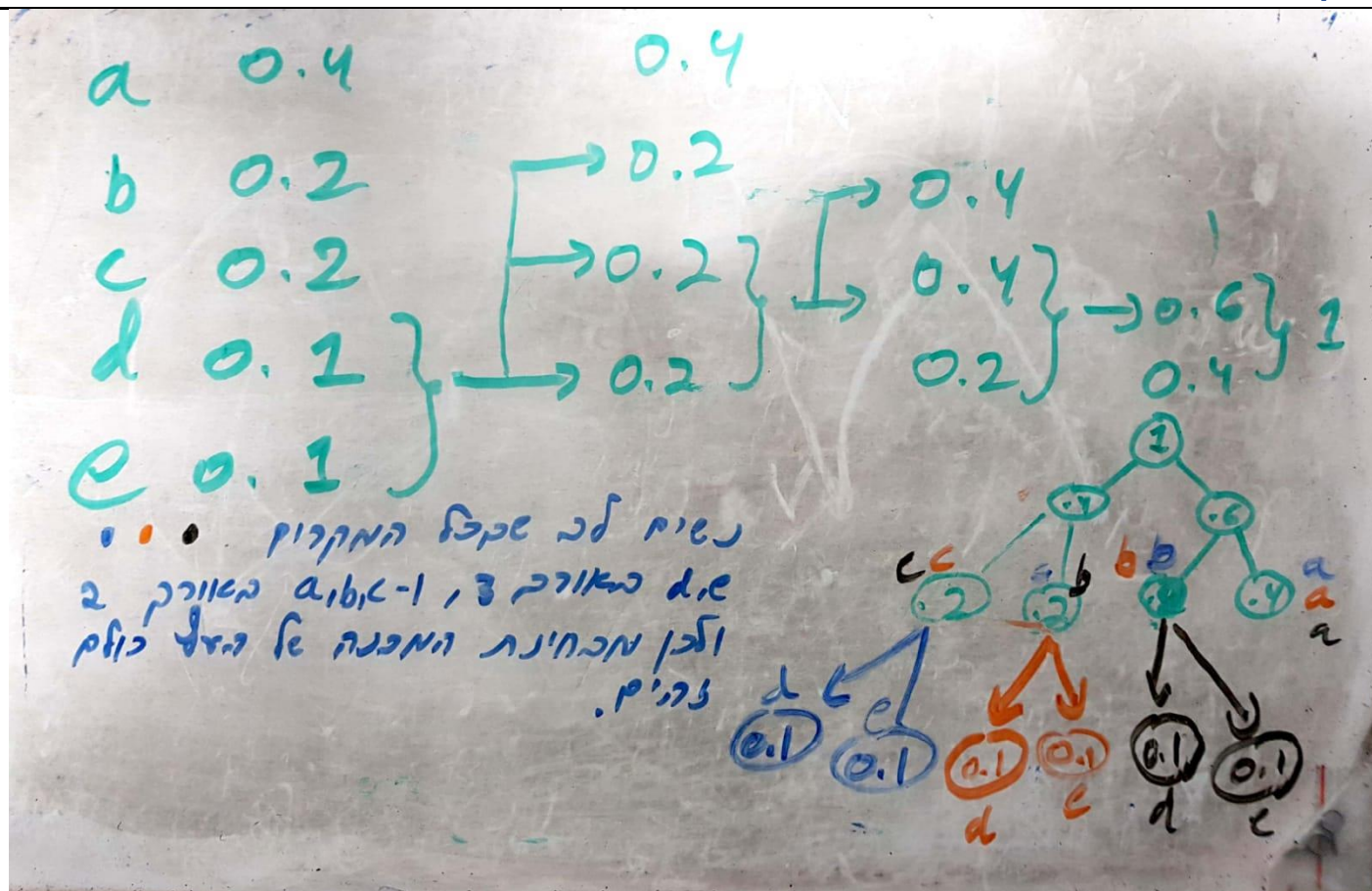
אלגוריתם לקידוד מילה (מספר) σ באמצעות C_γ : $Unary(|\langle \sigma \rangle_2|) \circ Binary^*(\sigma)$.

אלגוריתם לקידוד מילה (מספר) σ באמצעות C_δ : $C_\gamma(|\langle \sigma \rangle_2|) \circ Binary^*(\sigma)$.

פתרון מבחן לדוגמה

שאלה 1

סעיף א'



סעיף ב'

נסמן ב- C_1, C_2, C_3 את הקידודים לעיל כאשר הסדר מלמעלה למטה ומשמאל לימין. מספיק לחשב פעם אחת כי כולם שווים אומנם אנו ננצל זאת על מנת לוודא כי לא שגינו.

$$E[C_1] := \sum_{i=1}^n p_i \ell_i = 0.1 \cdot 3 \cdot 2 + 0.2 \cdot 2 + 0.2 \cdot 2 + 0.4 \cdot 2 = 2.2[b/s]$$

$$E[C_2] := \sum_{i=1}^n p_i \ell_i = 0.4 \cdot 1 + 0.2 \cdot 2 + 0.2 \cdot 3 + 2 \cdot (0.1 \cdot 4) = 2.2[b/s]$$

$$E[C_3] := \sum_{i=1}^n p_i \ell_i = 1 \cdot 0.4 + 2 \cdot (3 \cdot 0.1) + 2 \cdot (2 \cdot 0.3) = 2.2[b/s]$$

ואכן מתקבל כי $E[C_1] = E[C_2] = E[C_3]$, כפי שצפינו. חשוב לא לשכוח את יחידות המידה ביט למחרוזת.

סעיף ג'

$$H(P) = - \sum_{i=1}^n p_i \log_2 p_i = -(0.4 \log_2 0.4 + 2 \cdot (0.2 \cdot \log_2 0.2) + 2 \cdot (0.1 \cdot \log_2 0.1)) \\ \approx 2.12[b/s]$$

נראה שתוצאה זאת אכן הגיונית שכן $E[C] \geq H(P)$ ואכן $2.2 \geq 2.12$. בדיקה נוספת שרצוי לבצע היא שהפמן נותן מילת קוד ממוצעת שרחוקה לכל היותר בביט מהאנטרופיה ואכן מתקיים כי $2.2 \leq 3.12$.

סעיף ד'

נזכיר שלכל קידוד כנ"ל אנחנו נצטרך להוסיף את ה-prelude ובאלגוריתם הקנוני הוא היה מילת הקוד הראשונה בכל בלוק ואורכה בכל בלוק. ולכן לצמצם את התקורה ככל הניתן וזה יקרה כאשר נקטין את מספר הבלוקים. תשובה נוספת הא שבעץ ה-skeleton המצומצם אנחנו רוצים שההפרש בין העצים השלמים יהיה לכל היותר אחד, הקידוד היחיד שיקיים זאת הוא קידוד C_1 .

שאלה 2

סעיף א'

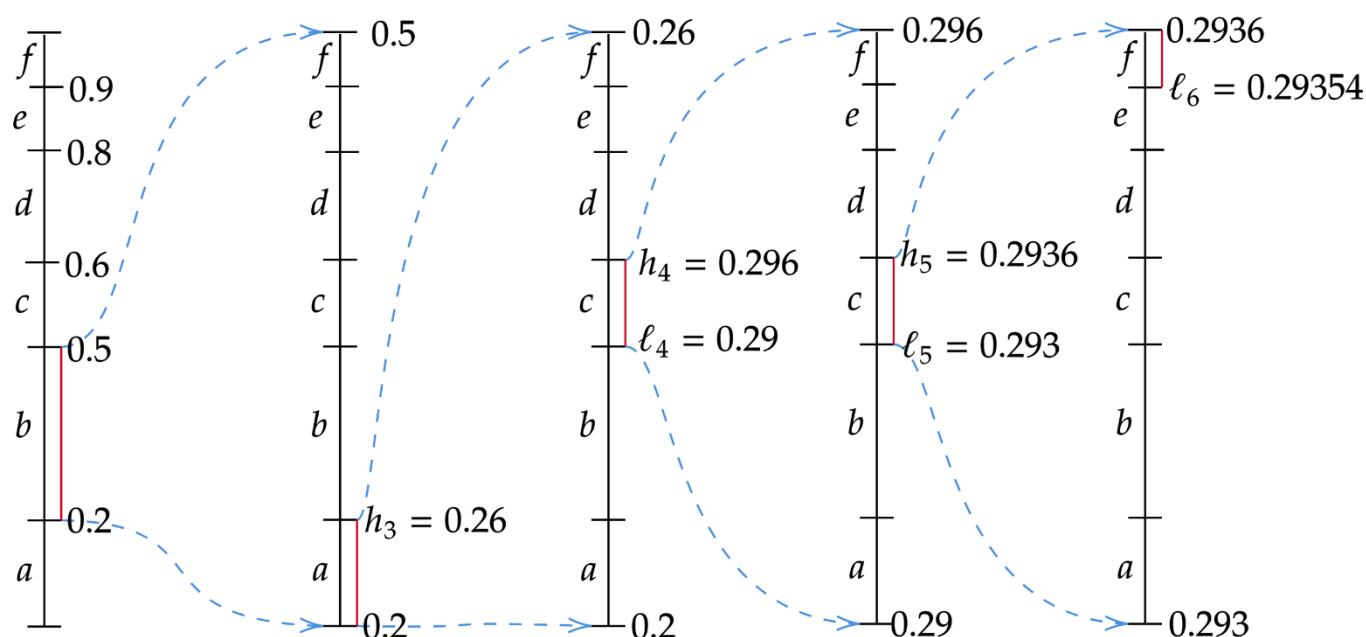
$\Pr[baccf] = \Pr[b] \cdot \Pr[a] \cdot \Pr[c] \cdot \Pr[c] \cdot \Pr[f] = 0.3 \cdot 0.2 \cdot 0.1^2 \cdot 0.1 = 0.00006 = 6 \cdot 10^{-5}$
 כאשר השוויון הראשון נובע מהאי תלות בין המאורעות.

סעיף ב'

אנחנו מתחילים את ההרצה כאשר $\ell = 0, h = 1$. קוראים את התו $\sigma = b$ מעדכנים את הטווח להיות $r = 1$ ומחשבים

σ	$hb(\sigma)$	$lb(\sigma)$
a	0.2	0
b	0.5	0.2
c	0.6	0.5
d	0.8	0.6
e	0.9	0.8
f	1	0.9

כעת ניתן לקודד באופן הבא:



$$h_3 = 0.2 + 0.2 \cdot 0.3 = 0.26$$

$$\ell_5 = 0.29 + 0.5 \cdot 0.006 = 0.293$$

$$\ell_4 = 0.2 + 0.5 \cdot 0.06 = 0.29$$

$$h_5 = 0.29 + 0.6 \cdot 0.006 = 0.2936$$

$$h_4 = 0.26 + 0.6 \cdot 0.06 = 0.296$$

$$\ell_6 = 0.293 + 0.9 \cdot 0.0006 = 0.29354$$

אנחנו צריכים להוציא מספר out כך ש- $\ell_6 \leq out < h_6$, נבחר לדוגמה $out = 0.29355$.

סעיף ג'

באופן כללי בקידוד אריתמטי הרעיון הוא לעשות scaling כך שבכל פעם הקידוד של $\sigma = \sigma_1 \dots \sigma_n$ הוא באינטרוול המתאים ל- $\prod_{i=1}^n \Pr[\sigma_i]$, ואכן $r_6 = h_6 - \ell_6 = 0.2936 - 0.29354 = 6 \cdot 10^{-5}$

שאלה 3

סעיף א'

נזכיר כי מודל סטטי למחצה אנו שולחים בתקורה (prelude) את מספר התווים את הקידוד של התווים ב-*ascii*. סה"כ נקבל:

$$\text{prelude}(C) = 8 + 8 \cdot |\Sigma| = 32[\text{bits}]$$

סעיף ב'

כעת נצטרך לקודד גם את ההסתברויות ע"פ נתוני השאלה ניתן לייצג הסתברות באמצעות 4 סיביות ולכן נקבל

$$\text{prelude}(C') = 32 + 4 \cdot 3 = 44[\text{bits}]$$

סעיף ג'

נזכיר ש- C_γ מורכב משני חלקים: חלק ראשון - $\text{Unary}(|\text{Binary}(i)|)$, חלק שני - $\text{Binary}^*(i)$. נסמן ב- $w = ababbac$ ונשים לב שמתקיים

i	σ_i	C_γ	ℓ_i	$f_i(w)$
1	a	$\langle B_1 \rangle_1 B_1^* = \langle 1 \rangle_1 \epsilon = 0$	1	3
2	b	$\langle B_2 \rangle_1 B_2^* = \langle 10 \rangle_1 0 = 100$	3	3
3	c	$\langle B_3 \rangle_1 B_3^* = \langle 11 \rangle_1 1 = 101$	3	1

לפיכך מתקיים:

$$E[C, P] = \frac{\text{prelude}(C) + \sum_{i=1}^n \ell_i \cdot f_i(w)}{7} = \frac{32 + 3 \cdot 1 + 3 \cdot 3 + 1 \cdot 3}{7} = \frac{47}{7} \approx 6.714$$

$$E[C', P] = \frac{\text{prelude}(C') + \sum_{i=1}^n \ell_i \cdot f_i(w)}{7} = \frac{44 + 3 \cdot 1 + 3 \cdot 3 + 1 \cdot 3}{7} = \frac{59}{7} \approx 8.428$$

סעיף ד'

באמצעות LZSS מנקודת את ab ואז נבצע הצבעה עצמית $(2, 2n - 2)$. בסדר גודל נקבל $O(1)$. אם אי אפשר להשתמש בהצבעות עצמיות אז נבצע $ab(2, 2)(4, 4)(8, 8), \dots$ כלומר קידוד בסדר גודל $O(\log n)$.

שאלה 5

סעיף א'

נסמן $C_1 = \{101, 1101, 1011, 1100\}$ ונחפש סופיות מתנדדות נשים לב ש-101 היא תחילית של 1011 ולכן 1 היא סיפא מתנדדת, נוסיף אותה למילון ונקבל $C_2 = \{101, 1101, 1011, 1100, 1\}$ כעת נבחין כי 1 היא תחילית של 101, 1101, 1011, 1100 ולכן 01, 101, 011, 100 הם סיפות מתנדדות. היות ו-101 היא מילת קוד, כלומר $101 \in C_1$ אזי הקוד אינו UD. נראה כיצד לבנות את הדוגמה:

$$\underbrace{101}_a \underbrace{1101}_b = \underbrace{1011}_b \underbrace{101}_a$$

בעצם לקחנו את המילה המתנגשת הוספנו לה את הסיפא המתנדדת והשלמנו למילה השניה.

סעיף ב'

תנאי הכרחי על מנת שקוד יהיה UD הוא ש- $K(C) \leq 1$. בנוסף אנו יודעים ש- $\ell_1 = \ell_2 = \ell_3 = 3$ וכל שאר $\ell_4 = \dots = \ell_{x+3} = 8$, לפיכך נקבל את האילוץ:

$$\sum_{i=1}^{x+3} 2^{-\ell_i} = 3 \cdot 2^{-3} + (x + 3 - 4 + 1) \cdot 2^{-8} \leq 1 \iff x \leq 160 \iff \boxed{x + 3 \leq 163}$$

ולכן גודל הא"ב המקסימלי הוא 163.

מבחן 2015 מועד א'

שאלה 1

סעיף א'

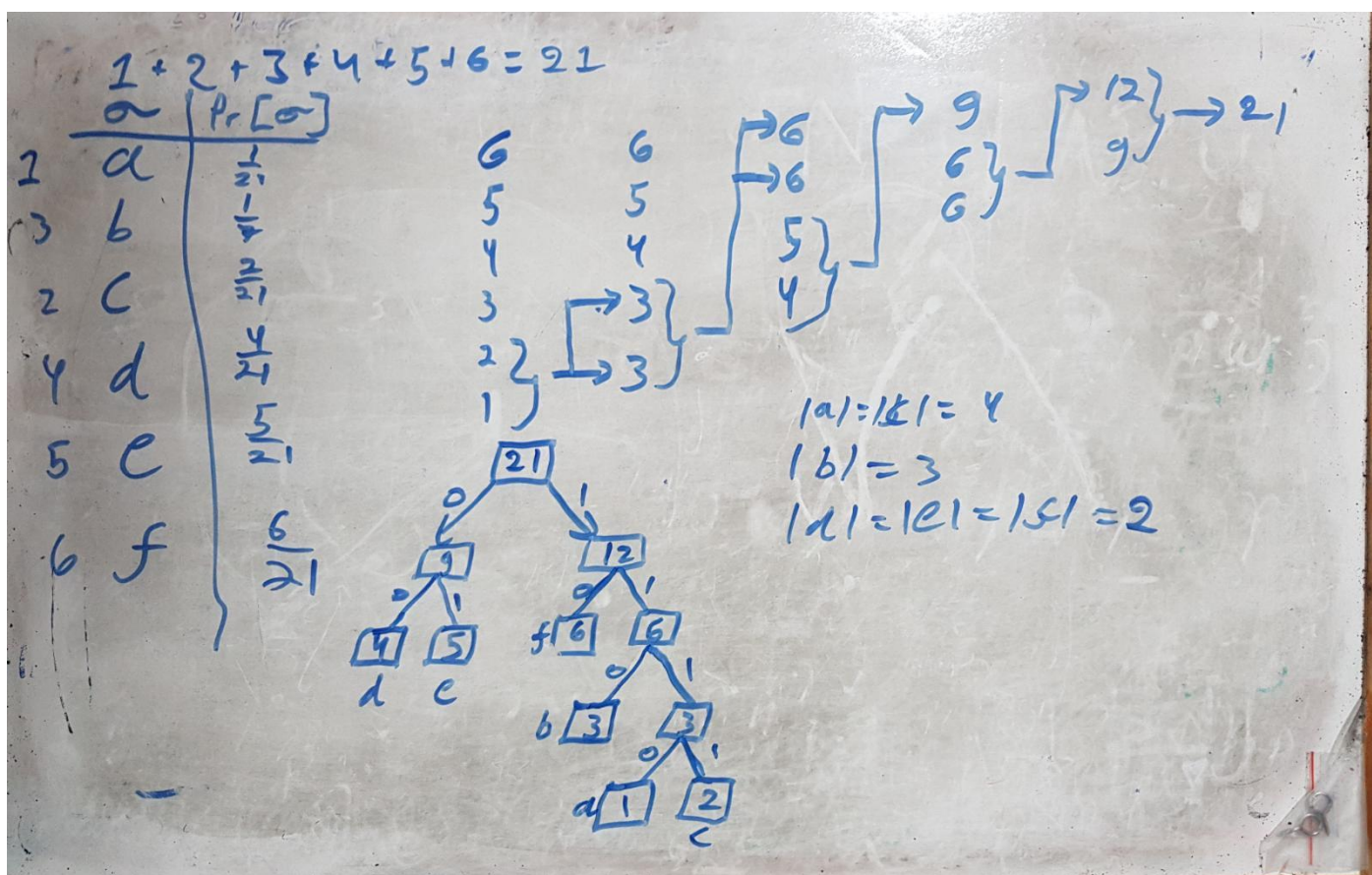
ע"פ משפט קיים קוד פרפיקסי בעל האורכים 1,2,3,3,4 אם ורק אם

$$2^{-1} + 2^{-2} + 2^{-3} + 2^{-3} + 2^{-4} \leq 1$$

היות ואגף שמאל שווה ל-1, $\frac{17}{16} > 1$, האי שוויון אינו מתקיים ומכך נובע שלא קיים קוד פרפיקסי כנ"ל.

סעיף ב'

הוכחנו בכיתה כי קוד הפמן הוא קוד אופטימלי לכן נחשב את קוד מילת הקוד הממוצעת בקידוד של הופמן ונראה אם הם שווים:



$$E[C_H, P] = 4 \cdot \left(\frac{1}{21} + \frac{2}{21} \right) + 3 \cdot \frac{3}{21} + 2 \cdot \left(\frac{4+5+6}{21} \right) = \frac{17}{7} \approx 2.428$$

כאשר C_H הוא הקוד ההתקבל מהפמן המתואר בציור.

$$E[C, P] = 3 \cdot \left(\frac{1+3+2+4}{21} \right) + 2 \cdot \left(\frac{5+6}{21} \right) = \frac{52}{21} \approx 2.476$$

כאשר C הינו הקוד המופיע בשאלה, לפיכך נסיק כי הקוד C אינו אופטימלי, שכן מצאנו קוד אופטימלי יותר.

סעיף ב'

נרחיב את השאלה וניתן דוגמה שבה $|C_\delta(n)| \leq |C_\gamma(n)|$ ולהפך. העניין הוא שכלל שהמספר גדול יותר ככה C_δ טובה יותר מ- C_γ , ננצל עובדה זאת נתחיל מהמקרה הראשון

$$n = 2^{12} - 1 \Rightarrow \langle n \rangle_2 = 1^{12} \Rightarrow C_\gamma(n) = \text{Unary}(|\langle n \rangle_2|)B^*(12) = 1110100$$

$$C_\delta(n) = C_\gamma(|\langle n \rangle_2|) \circ B^*(n) = 1110100 \circ 1^{11} \Rightarrow |C_\delta(n)| = 11 + 7 = 18$$

$$C_\gamma(n) = \text{Unary}(|\langle n \rangle_2|) \circ B^*(n) = 1^{11}0 \circ 1^{11} \Rightarrow |C_\gamma(n)| = 12 + 11 = 23$$

כעת נענה על המקרה השני:

$$n' = 2 \Rightarrow \langle 2' \rangle_2 = 10 \Rightarrow C_\gamma(2) = 10 \circ 0$$

$$\Rightarrow C_\delta(n') = C_\gamma(|\langle n' \rangle_2|) \circ B^*(n') = C_\gamma(2) \circ 0 = 100 \circ 0$$

הטענה נובעת.

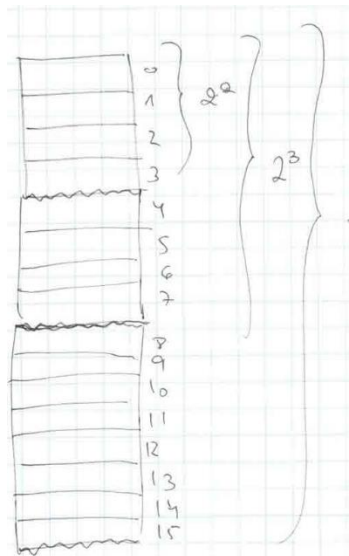
שאלה 5

סעיף א'

נתאר את הריצה באמצעות טבלת מעקב

OLD	NEW	<	<	0 4 6	Code	Symbol
0	1	B	B	A	0	A
1	3	AB	A	AB	1	B
3	5	ABA	A	ABA	2	C
5	6	ABAA	A	ABAA	3	AB
6	2	C	C	C	4	BA
2	7	ABAAC	A	ABAAC	5	ACA
7	9	ABAACA	A	ABAACA	6	ABAA
9	2	C	C	C	7	ABAAC
					8	CA
					9	ABAACA
					10	ABAACA

סעיף ב'



פתרון אחד הוא לשים לב שהמילון מוכפל רק כאשר אנחנו באים לכתוב את הכניסה ב- 2^i , והמילון של המקודד והמפענח זהים. לפיכך, כאשר נכתוב את 0,1,3 השתמשנו בשני סיביות, באופן דומה 5,6,2 נרשמים באמצעות שלושה סיביות, וכן 7,9,2 נרשמים באמצעות 4 סיביות. נסכום ונקבל:

$$3 \cdot 2 + 3 \cdot 3 + 3 \cdot 4 = 27bits$$

פתרון הנכון יהיה לקודד שוב את ההודעה

ABABABAABAACABAACABAACAC

ולספור את מספר הסיביות שנדרשים לקידוד.

codeword	σ
0	A
1	B
2	C
3	AB
4	BA
5	ABA
6	ABAA
7	ABAAC
8	CA
9	ABAACA

w	k	Output	Enter	Bits
A	B	0	AB	2
B	A	1	BA	2 (★)
A	B			
AB	A	3	ABA	3
A	B			
AB	A			
ABA	A	5	ABAA	3
A	B			
AB	A			
ABA	A			
ABAA	C	6	ABAAC	3
C	A	2	CA	3
A	B			
AB	A			
ABA	A			
ABAA	C			
ABAAC	A	7	ABAACA	4
A	B			
AB	A			
ABA	A			
ABAA	C			
ABAAC	A			
ABAACA	C	9	ABAACAC	4
C	EOF	2		4

ואכן הקידוד הוא 0,1,3,4,6,2,7,9,2, כאשר נסכום את העמודות של הביטים נקבל שנצטרך

$$2 + 2 + 3 + 3 + 3 + 3 + 4 + 4 + 4 = 28bits$$

נשים לב שבשלב (★) קודם מוציאים את הקוד ורק אחר כך מכניסים את המילה למילון, ולכן ההגדלה של המילון מתבצעת רק לאחר שלב שרשמנו את 3 בשתי סיביות כלומר 11. יש פער של סיבית בין התשובות, כרגע עוד לא מצאתי טעות וככל הנראה הפער נובע מכך שברגע שהמפענח מפענח את שלוש הוא עוד לא הכניס את BA בניגוד למקודד.

מבחן 2017 מועד ב' – שחזור

שאלה 1 – סעיף א' – האם הקוד $\{0,01,110\}$ הוא UD ?

פתרון. נריץ את המבחן

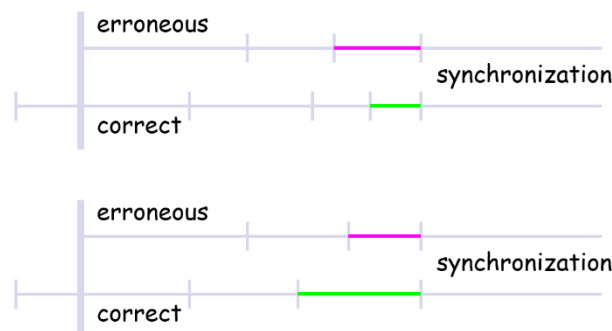
$\{0,01,110\} \rightarrow \{0,01,110,1\} \rightarrow \{0,01,110,1,10\} \rightarrow \{0,01,110,1,10,0\}$

קיבלנו את מילת הקוד 0 ולכן הקוד אינו UD , נראה דוגמה נגדית:

0110

בדיקות שכדאי לבצע במבחן

- כל התוצאות של קוד הפמן הם אופטימליות כלומר מחזירות מילת קוד ממוצעת מינימלית.
- אם מריצים הפמן מתקיים $E[C, P] \leq H(P) + 1$.
- $E[C, P] \geq H(P)$.
- לשים לב שעץ skeleton קשור להפמן אז אם שואלים על אחד מהם יעזור לשני.
- בקוד אריתמטי מספר הסיביות שנדרשות הוא $-\log_2(\prod_{i=1}^n p_i)$
- בקידוד LZW קודם כל מוצאים את הפלט ואחר כך מוסיפים את המילה החדשה, כלומר הקידוד שלפני ההכנסה של המילה החדשה נשאר בגודל באותו מספר סיביות שהיה.
- $2|fc[i] + num[i]$, כלומר הערך הדצימאלי של מילת הקוד הראשונה בכל בלוק ועוד מספר המילים בכל בלוק תמיד זוגי.
- דרך לשלול שקוד נתון הינו קוד הפמן היא להראות שהקוד לא שלם, כלומר העץ לא מלא.
- דרך נוספת היא להראות שיש צמתים פנימיים שיש להם מילות קוד.
- עוד דרך היא להסתכל על לשים לב שבגלל שבכל פעם אנחנו ממיינים את השכיחויות אז ישנם עצים שלא יכולים להתקבל, זוהי גם דרך ליצור עץ שאינו יכול להתקבל ע"י הפמן.
- בקידוד בינארי מינימלי לזכור שהמילות הקצרות הם עם השכיחויות הנמכות.
- הפמן ובפרט הפמן קנוני חייב להיות עץ מלא ועץ מלא אם ורק אם $K(C) = 1$, ולכן אם נותנים לנו את המערך num של הפמן קנוני, $num = [x_1, \dots, x_n]$ אזי חייב להתקיים $K(C) = x_1 \cdot 2^{-1} + x_2 \cdot 2^{-2} + \dots + x_n \cdot 2^{-n} = 1$
- על מנת לבנות קוד אפיקסי נרשום את כל האופציות למילים בינאריות באורך 1,2,3,4 ונבחון מה מתאים.
- בהפמן יש סינכרון כי יש מילות קוד שהם סיפא של מילות קוד אחרות וזהו הדרך היחידה שיתבצע סינכרון:



כלומר או שהסיפא של הטעות מסתנכרנת עם הנכון או להפך.

- לכן קוד אפיקסי לא מסתנכר
- קוד בעל אורך k מסתנכר רק אם k סיביות נאבדות.

מבחן 2017 מועד ב' שאלה 4: קודד את Mississippi באמצעות הפמן דינמי.