

דחיסת נתונים א' - דף נוסחאות

משפטים והגדרות

קוד UD: קוד שניתן לפענחו בצורה יחידה
קוד מיד: קוד שבו כל מילת קוד מפוענחת בסיום קריאתה

- קוד מיד אינו הכרחי עבור תכונת UD (למשל במחרוזת סופית נשים מחיצות בסיום קריאתה)

קוד שלם: קוד שהוא אינסופי למחצה כך שכשנפענח אותו משמאל לימין הוא תמיד יפוענח באופן UD

וגם לא נתקע כאשר נפענח אותו.

אם קוד הוא שלם אזי הוא UD.

עבור כל קוד שהוא שלם מתקיים $1 \leftrightarrow K(C)$. הוא לא מבוזבז ואין צמתים פנימיים בעלי בן אחד.

אורך מילת קוד ממוצעת: $E(C, P) = \sum_{i=1}^n p_i \cdot |c_i|$

- קוד מיד \leftrightarrow קוד חסר רישות ניתן לבנות לו עץ
- אם קוד חסר רישות אז הוא UD
- קוד מיד \leftrightarrow קוד חסר רישות

קוד בעל יתירות מינימלית: קוד C הוא קוד בעל יתירות מינימלית אם מתקיים $E(C, P) \leq E(C', P)$ עבור כל קוד C' חסר רישות בעל n אותיות

אינפורמציה: $I(s_i) = -\log_2(p_i)$

- צריך לתכנן את הקוד כך שמילת קוד עבור s_i מכילה $I(s_i)$ תווים
- אם מתקיים $p_i = 1$ אז מתקיים $I(s_i) = 0$.
- אם הרצף $s_i s_j$ מופיע בהסתברות $p_{s_i s_j}$ אז $I(s_i s_j) = I(s_i) + I(s_j)$.

אנטרופיה: $H(P) = -\sum_{i=1}^n p_i \cdot \log_2(p_i)$

- אנטרופיה מהווה חסם תחתון (של סיביות) לאורך מילת הקוד
- את האנטרופיה מעגלים כלפי מעלה כי חייב לקחת מספר שלם של סיביות
- אנטרופיה זה בעצם ממוצע משוקלל של כמויות האינפורמציה
- לכל קוד UD מתקיים $H(P) \leq E(C, P)$
- קוד חסר רישות הוא קוד UD (הפוך לא בהכרח)

קראפט: $K(C) = \sum_{i=1}^n 2^{-|c_i|}$

- אם $K(C) \leq 1$ אזי קיים קוד חסר רישות C' כך ש: $|C| = |C'|$ וגם $E(C, P) = E(C', P)$
- אם $K(C) > 1$ אזי לא קיים קוד חסר רישות עם אותם אורכים
- אם קוד הוא UD אז $K(C) \leq 1$
- קיים קוד מיד $1 \leftrightarrow K(C)$

UD TEST

נגדיר: Dangling Suffix = המשלים של רישה של מילה מסוימת. אלגוריתם:

- הכנס את כל מילות הקוד לקבוצה (נקרא למילים אלו "מילים מקוריות").
- חזור על התהליך כל עוד שלב 1 מתקיים:
 - בדוק אם קיימת מילה שהיא רישה של מילה אחרת
 - אם קיימת מילה כזו, הוסיף לקבוצה את Dangling Suffix
 - אם Dangling Suffix שהוספנו היא מילת קוד מקורית אזי הקוד לא UD

Static

ההנחה היא שההסתברויות בלתי תלויות ואחידות. הקוד הוא סטדרטי ולכן אין prelude.

כמו כן, השימוש הוא בתווים האסקי, ולכן ההסתברות לכל תו הינה $\frac{1}{256}$ נשים לב כי האנטרופיה היא 8.

Semi-Static

כאן ההגבלה היא רק על התווים שמופיעים בטקסט במקום להשתמש בכל התווים באסקי. לכן, לאחר מעבר על הטקסט וספירת כמות התווים (ללא ספירת ההופעות של כל תו כי מניחים שההסתברות לכל תו היא $\frac{1}{n}$)

ובהנחה שמצאנו n תווים, אזי, prelude היא יראה כך:

- 8 סיביות ראשונות: יהיו מספר התווים
- $n \cdot 8$ הסיביות הבאות: יהיו קודי האסקי של התווים
- שאר הסיביות: מילות הקוד באורך קבוע

$$E(C, P) = H(P) + \frac{8n + 8}{|text size|}$$

Semi-Static (Self-probabilities)

בשיטה זו עוברים על הטקסט וסופרים את כמות התווים + כמות ההופעות של כל תו.

הסתברות של תו תהיה $\frac{\text{number of occurrences in the text}}{|text size|}$

בהנחה שמצאנו n תווים, אזי, prelude היא יראה כך:

- 8 סיביות ראשונות: יהיו מספר התווים
 - $n \cdot 8$ הסיביות הבאות: יהיו קודי האסקי של התווים
 - שאר הסיביות: ביטים המייצגים את ההסתברויות של כל תו
 - נניח שמצאנו n תווים ונניח שצריך q סיביות בשביל לקודד הסתברות אזי
- $$E(C, P) = H(P) + \frac{8n + pn + 8}{|text size|}$$
- נעביר את התווים עם ההסתברויות השכיחות תחילה.

Unary Code

- כל קידוד של תו במיקום i יהיה $1^{i-1}0$ (אינדקס נמוך להסתברות שכיחה) אלגוריתם להפיכה לקוד שלם:
- הורד את 0 במילת הקוד האחרונה

Binary Code

קוד בינארי פשוט: מס' הסיביות המינימלי עבורו ניתן לייצג את כמות הא"ב
קוד בינארי מינימלי: נאפשר לחלק ממילות הקוד להיות באורך $\lceil \log_2 n \rceil$ ולחלק להיות באורך $\lceil \log_2 n \rceil - 1$ בצורה הבאה:

מילות	אורך
$2n - 2^{\lceil \log_2 n \rceil}$	$\lceil \log_2 n \rceil$
$n - 2^{\lceil \log_2 n \rceil}$	$\lceil \log_2 n \rceil$

כאשר מילות הקוד הקצרות בצד שמאל של העץ.

Elias

קוד C_γ :

אלגוריתם:

- סדר את האותיות לפי ההסתברויות שלהן (הכי נפוץ יהיה ראשון עם אינדקס 1).
- כל מילת קוד תראה כך: XY
- כאשר X הוא **מספר הסיביות** הנדרש על מנת לייצג את האינדקס באונארי וי הוא האינדקס עצמו בבינארי ללא ה1 המוביל (הכי שמאלי).

קוד C_δ :

זהה לחלוטין לקוד C_γ פרט לכך שא מקודד בעזרת C_γ במקום באונארי.

Golomb + Rice

קוד Rice הוא מקרה פרטי של Golomb כאשר הדלי הוא חזקה של 2. כלומר $b=2^k$.

Encode
$Golom_encode(x, b):$ $q = \lfloor (x - 1) / b \rfloor$ $r = x - q * b$ $p1 = Unary_encode(q + 1)$ $p2 = Minimal_binary_encode(r, b)$ $return p1 \cdot p2$
Decode
$Golom_decode(xy, b):$ $q = Unary_decode(x) - 1$ $r = Minimal_binary_decode(y)$ $return r + q * b$

p2 היא מילת הקוד ה-ב"ב בינארי **מינימלי** בגודל b. (כאשר המילים מסודרות בסדר לקסיקוגרפי ומהקצרה לארוכה).

כאשר xy היא מילת הקוד אותה מפענחים.
 x הוא הרישה (האונארי) של מילת הקוד
 (ניתן לזהות כי מסתיים ב0)
 וי זה שאר מילת הקוד

Fibonacci Code

אלגוריתם להפיכה לקוד מיד:

- חשב קוד פיבונאצ'י
 - הוסיף 1 לתחילת מילת קוד
 - הפוך את המילים
- תכונות:
- קוד פיבונאצ'י הוא קוד שלם (הוכחנו בעזרת קראפט)
 - יש F_k מילות קוד באורך k+1 (כולל ה1 שמוסיפים בסוף) כאשר $F_1 = F_2 = 1$
 - מספר n המקיים $F_{k+1} \leq j \leq F_{k+2}$ צריך k+1 סיביות - F_{k+2}

אלגוריתמים לפיתירת הבעיה (של החסרון שאין מילה באורך 1):

- אלגוריתם 1:
 - ניקח רק את המילים שמתחילות ב1 (לאחר שהפכנו)
 - נוריד את הסיבית האחרונה (שהיא 1 כי מסתיים ב11)
- אלגוריתם 2 (שקולה לראשונה):
 - נוריד את הסיבית האחרונה (שהיא 1 כי מסתיים ב11)
 - נוסיף 10 להתחלה של כל מילת קוד
 - לבסוף נוסיף את המילה 1 כמילת הקוד הראשונה

תכונות של האלגוריתמים:

- מכילות מילת קוד אחת בלבד באורך 1.
- פרט למילה הראשונה, ישנן F_{k-2} מילות קוד באורך k.
- ניתן להוכיח שתמיד יותר טוב מקוד C_δ וקוד C_γ .

Shannon Code

$$H(P) \leq L \leq H(P) + 1 \quad \text{אז:} \quad \log_2 \frac{1}{p} \leq \left\lceil \log_2 \frac{1}{p} \right\rceil \leq \log_2 \frac{1}{p} + 1$$

Coding blocks of symbols

$$H(p(s_i, s_j)) = 2H(p(s))$$

החולשה של שיטה זו שאנו יכולים לקבל עץ שאינו מלא.

encoding full binary tree

- בעץ מלא עם n עלים יש n-1 צמתים פנימיים
- צריך 2n-1 ביטים כדי לקודד עץ בינארי עם n עלים
- כל קוד חסר רישות ניתן לייצג ע"י עץ קידוד עץ:
- 8 סיביות ראשונות: מספר האלף בית
- $n \cdot 8$ סיביות הבאות: קידודי האסקי עבור כל אות
- (נניח שהאותיות המתקבלות מסודרות מהעלה השמאלי לעלה הכי ימני לפי הסדר)
- נעביר מחזרות המייצגת את העץ כך שכל סיבית היא אינדקטור אם צומת מסוים הוא פנימי (1) או עלה (0).

Shannon-Fano codes: אלגוריתם:

- סדר את ההסתברויות מהגדול לקטן.
- חלק את ההסתברויות לשני קבוצות כך שההפרש ביניהם יהיה מינימלי ככל הניתן (אם החלוקה אינה ודאית ויש התלבטות עבור הסתברות מסוימת אז ניקח אותה לקבוצה השמאלית שבה האורך קצר יותר)
- לקבוצה השמאלית ניתן את הספרה 0 ולימנית את הספרה 1.
- קוד שאינו תמידי אפקטיבי.

Huffman Code (Minimum Redundancy Coding)

אלגוריתם:

```
Huffman(S):
n = |S|
Q1 = S // min heap
For i=1 to n-1:
  Do allocate-node(Z) // new node
  z.left = x = Extract_min(Q1)
  z.right = y = Extract_min(Q1)
  z.weight = x.weight + y.weight
  insert(Q1, z)
return Extract_min(Q1)
```

משפט: בהינתן משקלים w_1, \dots, w_n הופמן מקצה אורכי קוד l_1, \dots, l_n כך ש $\sum_{i=1}^n w_i l_i$ הוא מינימלי.

למה 1: קוד אופטימלי עבור קובץ תמיד מיוצג ע"י עץ בינארי מלא כך שלכל צומת פנימית יש 2 ילדים.

למה 2: בעץ אופטימלי 2 ילדים עם המשקלים הכי נמוכים נמצאים ברמה הנמוכה ביותר.

למה 3: בעץ אופטימלי 2 ילדים עם המשקלים הכי נמוכים יכולים להיות אחים.

בהחלפת 0/1 עבור צמתים פנימיים הופמן יכול לתת 2^n קודים שונים.

Huffman (Canonical)

אלגוריתם לעץ הופמן קנוני משוך שמאלה (בהינתן קוד הופמן):
 בהינתן אותיות $\{i\}$ ואורך אופטימלי $\{l_i\}$

```
maxlength=max{l_i}
// find the number of codewords of each length
For l = 1 to maxlength
  num[l]=0
For i = 1 to n
  num[l_i]++
// store the first codeword
Firstcode[maxlength]=0
For l=maxlength-1 downto 1
  firstcode[l] = (firstcode[l+1]+num[l+1])/2
For l=1 to maxlength
  nextcode[l] = firstcode[l]
For i=1 to n
  codeword[i]=nextcode[l_i]
  symbol[l_i,nextcode[l_i]-firstcode[l_i]] = i
  nextcode[l_i]++
```

ייצוג עץ הופמן קנוני:

- העבר את האותיות לפי הסדר (מהאורך לקצר ובסדר לקסיקוגרפי)
- העבר טבלה שבה יש את רק את האותיות שבתחילת כל בלוק (כלומר, המילה הראשונה בכל אורך [מהארוכה אל הקצרה]) כך שהטבלה מכילה את האותיות הנ"ל ואת הקידוד שלה.

אלגוריתם פענוח הודעה בינארית (שימוש במטריצת symbol שבינו):

```
maxlength=max{l_i}
v=nextInputBit()
l = 1
while v < firstcode[l] do
  v = 2v+nextInputBit()
  l++
// the integer v is a valid codeword of l bits
return symbol[l,v-firstcode[l]]
// the index of the decoded symbol
```

הגדרה שקולה לעצי הופמן קנוני:

אלגוריתם לבניית עץ הופמן קנוני משוך ימינה.

- תחילה, נפעיל את אלגוריתם הופמן הרגיל
- לאחר מכן, נמייין את מילות הקוד שקיבלנו מהאלגוריתם הרגיל ע"פ האורכים של מילות הקוד
- לכל מילת קוד ניתן הסתברות 2^{-l_i} כאשר l_i היא אורך מילת הקוד
- לכל מילה נגדיר מספר $k = \sum_{j=1}^{i-1} 2^{-l_j}$. כלומר, עבור מילת קוד i אנו סוכמים את ההסתברויות של מילות הקוד שהאינדקס שלהן קטן ממש מהאינדקס i
- מילת הקוד החדשה תהיה l_i הספרות שאחרי הנקודה העשרונית של המספר k הנ"ל

Huffman (D-ary)

לעץ הופמן די-ארי עם n צמתים פנימיים יש $(D-1) \cdot n + 1$ עלים

אלגוריתם:

- חשב מה מספר העלים n_0 שצריך להוסיף כך שמתקיים $n + n_0 = (D-1) \cdot n' + 1$
- הפעל הופמן עבור D עלים בכל שלב
- תן ערך $D-1, \dots, 0$ לכל קשת (הספירה מתחילה מ0)
- התאם עלה להסתברות שלו (וזרוק את העלים עם 0)

דחיסת נתונים א' - דף נוסחאות

<div data-bbox="379 116 499 138" data-label="Section-Header"> <h2> <p>Move To Front</p> </h2> </div> <div data-bbox="424 138 499 161" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="67 161 499 183" data-label="Text"> <p>קודד את האינדקס של כל תו ואז העביר אותו לראש הרשימה.</p> </div> <div data-bbox="323 183 499 206" data-label="Text"> <p>בחזרה לאלגוריתם BWT:</p> </div> <div data-bbox="424 206 499 228" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="25 228 499 445" data-label="List-Group"> <ul style="list-style-type: none"> נפעיל על ההודעה את אלגוריתם BWT לאחר מכן נפעיל על העמודה האחרונה (שאנו מעבירים) את אלגוריתם Move To Front על הפלט של Move To Front נפעיל דחיסה סטטיסטית הופמן/אל/אס/RunLength/אריתמטי לבסוף נוכל להעביר את האינדקס שהאלגוריתם BWT פלט + הקידוד (הסופי) של העמודה <div data-bbox="199 398 328 445" data-label="Equation-Block"> $PPM(+EP) \\ H(X Y) \leq H(X)$ </div> </div> <div data-bbox="75 461 499 506" data-label="Text"> <p>הוכחנו בהרצאה כי אנטרופיה מותנת קטנה-שווה לאנטרופיה רגילה.</p> </div> <div data-bbox="319 506 499 528" data-label="Equation-Block"> $H(X Y) \leq H(X)$ </div> <div data-bbox="424 528 499 551" data-label="Text"> <p>כלומר, אלגוריתם:</p> </div> <div data-bbox="25 551 499 804" data-label="Text"> <p>נוסיף לטבלה תו \$ שערכו יהיה מספר האלף בית באותו הקשר ונקודד אותו על מנת לרדת רמה. אם הגענו לרמה 1- זה אומר שצריכים לקרוא תו חדש ולכן אנו מצפים לקידוד שלו. הנתו \$ יהיה חלק מהאלף בית וההסתברות לקודד תו חדש מהאלף בית תהיה אחידה. עבור קידוד של תו קיים ההסתברות לקידוד תהיה אחידה על פני יחס ההופעות של תו מתוך ההופעות של שאר התווים באותה רמה. בכל שלב נסמן בסוגריים מה ההסתברות לקידוד התו. מספר הסיביות הדרוש יהיה $\lceil \log_2 \frac{1}{p} \rceil$. חשוב לזכור: תמיד נעזר את כל הרמות לאחר קריאת תו.</p> </div> <div data-bbox="319 819 499 842" data-label="Text"> <p>Exclusion Principle (EP):</p> </div> <div data-bbox="54 842 499 864" data-label="Text"> <p>אלגוריתם לחישוב הסתברויות ברמה מסוימת $0 \leq k < k_{max}$:</p> </div> <div data-bbox="67 864 448 1207" data-label="Text"> <pre>estimate_prob(k) { sum=count[\$]; for (i=0; i<= Σ ; i++) if k=kmax count[i]=count[a_i previous k symbols] else if (k>=0) count[i]=(count[a_i previous k+1 symbols]=0)*count[a_i previous k symbols] else count[i]=(count[a_i]=0) Sum+=count[i] for (i=1; i<= Σ ; i++) p(encode a_i) = count[i]/sum p(encode \$) = count[\$]/sum }</pre> </div> <div data-bbox="319 1232 499 1254" data-label="Section-Header"> <h2> <p>Grammer Compression</p> </h2> </div> <div data-bbox="25 1254 499 1276" data-label="Text"> <p>דטרמיניסטי (אין הסתעפויות בגזירה) ואסור שיהיה מעגל בגזירות.</p> </div> <div data-bbox="424 1276 499 1299" data-label="Text"> <p>Sequiter:</p> </div> <div data-bbox="25 1299 499 1487" data-label="List-Group"> <ul style="list-style-type: none"> תווים עוקבים לא יופיעו יותר מפעם אחת חוק צריך להופיע לפחות פעמיים <p>אלגוריתם:</p> <p>כל עוד לא הגענו אל EOF:</p> <ul style="list-style-type: none"> אם ישנם 2 תווים צמודים שחוזרים על עצמם – צור כלל חדש אם ישנו חוק שמשמשמשים בו פעם אחת – הסר אותו והכנס את התוכן למופע היחיד שלו </div> <div data-bbox="440 1487 499 1509" data-label="Section-Header"> <h2> <p>Re-Pair</p> </h2> </div> <div data-bbox="424 1509 499 1532" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="25 1532 499 1677" data-label="List-Group"> <ul style="list-style-type: none"> מצא זוג תווים הכי נפוץ. לדוגמה, ab. צור כלל. לדוגמה, A → ab. החלף את כל מופעי הזוג בכלל שיצרת. לפי הדוגמה הנ"ל החלף את מופעי הזוג בכלל A. חזור ל1 כל עוד מצאת זוג שמופיע יותר מפעם אחת (בחר את הכי נפוץ). </div> <td data-bbox="504 96 1038 2172"> <div data-bbox="751 103 791 125" data-label="Section-Header"> <h2>LZ77</h2> </div> <div data-bbox="683 125 857 147" data-label="Text"> <p>(offset,length,symbol)</p> </div> <div data-bbox="909 147 1031 170" data-label="Text"> <p>אלגוריתם קידוד:</p> </div> <div data-bbox="555 170 983 304" data-label="Text"> <pre>p=1 // next char to be coded while (there is text to be coded): search for the longest match for S[p...] in s[p-w...p-1] // suppose match at pos m with len = 1 output the triple (p-m,l,s,p+1) p = p+l+1</pre> </div> <div data-bbox="751 311 791 333" data-label="Section-Header"> <h2>LZSS</h2> </div> <div data-bbox="729 333 1026 356" data-label="Text"> <p>(offset,length) [1] או בתווים בודדים [0].</p> </div> <div data-bbox="751 356 783 378" data-label="Section-Header"> <h2>LZS</h2> </div> <div data-bbox="627 378 909 400" data-label="Text"> <p>(offset1,length1), (offset2,length2)</p> </div> <div data-bbox="751 400 791 423" data-label="Section-Header"> <h2>LZ78</h2> </div> <div data-bbox="951 423 1031 445" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="533 445 1031 672" data-label="List-Group"> <ul style="list-style-type: none"> בנה מילון D שיכיל בהתחלה רק את המילה הריקה (שהאינדקס שלה הוא 0) הכנס תתי מחרוזת (של המחרוזות שנרצה לקודד) ותן לה את האינדקס הבא בתור שפנוי. נסה "לתפוס" את המילה הכי ארוכה שכבר במילון ונקודד את האינדקס שלה + הקידוד של התו הבא. במקביל, הוסף למילון את אותה הרישה בשרשרת התו הבא שמקודדים שיטה לדעת כמה ביטים צריך לקידוד מס': נשאל – "כמה ביטים צריך בשביל לקודד x מספרים" כאשר x הוא האינדקס באותה שורה </div> <div data-bbox="751 656 791 678" data-label="Section-Header"> <h2>LZW</h2> </div> <div data-bbox="595 678 1031 701" data-label="Text"> <p>אלגוריתם קידוד: (טבלה התחלתית גדולה ממש מגודל הא"ב)</p> </div> <div data-bbox="584 701 946 1081" data-label="Text"> <pre>Dictionary=single characters w=first char of input repeat { k=next char if(EOF) output code(w) else if (w · k)∈Dictionary w = w · k else output code(w) // output in num of bits corresponding to table size Dictionary.add(w · k) // Only enlarge the table when entering information and when it is full w = k }</pre> </div> <div data-bbox="903 1081 1031 1104" data-label="Text"> <p>אלגוריתם פענוח:</p> </div> <div data-bbox="584 1104 1031 1126" data-label="Text"> <p>חשוב: המפענת תמיד מגדיל את המילון ברגע שהמילון מתמלא</p> </div> <div data-bbox="584 1126 946 1514" data-label="Text"> <pre>Initialize table with single character strings OLD = first input code Output translation of OLD While not end of input stream { NEW = next input code If new is not in the string table S = translation of OLD S = S·C else S = translation of NEW Output S C = first character of S Enter Translation(OLD)·C to the string table // Increase the table as soon as the dictionary fills up OLD = NEW }</pre> </div> <div data-bbox="702 1536 834 1559" data-label="Section-Header"> <h2>Run length Code</h2> </div> <div data-bbox="716 1559 1031 1581" data-label="Text"> <p>זוג סדור (num,symbol) ע"פ רצף של תווים.</p> </div> <div data-bbox="842 1581 1031 1603" data-label="Text"> <p>Binary Run length Code</p> </div> <div data-bbox="746 1603 1031 1626" data-label="Text"> <p>escape codeword = ספרה מקסימלית.</p> </div> <div data-bbox="638 1626 898 1648" data-label="Section-Header"> <h2>The Burrows Wheeler Transform</h2> </div> <div data-bbox="951 1648 1031 1671" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="520 1671 1031 1879" data-label="List-Group"> <ul style="list-style-type: none"> בנה מטריצה מסדר n×n (כאשר הטקסט באורך n). כך שהשורות יהיו כל הטרנספורמציות הציקליות של הטקסט. כלומר, השורה הראשונה תתחיל באות הראשונה של הטקסט. השורה השנייה תתחיל מהאות השנייה של הטקסט וכו' (ממשיכים לשכתב את הטקסט בצורה מעגלית בסוף השורה) מייין את השורות בסדר לקסיקוגרפי (לפי האלף בית abcdefghijklmnopqrstuvwxyz) החזר את העמודה האחרונה במטריצה (מלמעלה למטה) + אינדקס של השורה המקורית של הטקסט (ספירה מ0). </div> <div data-bbox="978 1879 1031 1901" data-label="Text"> <p>פענוח:</p> </div> <div data-bbox="903 1901 1031 1924" data-label="Text"> <p>אלגוריתם פענוח:</p> </div> <div data-bbox="679 1924 847 2058" data-label="Text"> <pre>BWT(L, index): F = sort(L) M = index For (i=0; i < n; i++) T[i] = F[M] M = S[M]</pre> </div> <td data-bbox="1043 96 1568 2172"> <div data-bbox="1222 116 1378 138" data-label="Section-Header"> <h2>Huffman (Adaptive)</h2> </div> <div data-bbox="1430 138 1559 161" data-label="Text"> <p>Sibling property:</p> </div> <div data-bbox="1085 161 1559 351" data-label="List-Group"> <ul style="list-style-type: none"> משקל צומת שווה לסכום המשקל של שני ילדיו. הצמתים ניתנים לסידור ע"פ משקלם בסדר \geq כך שמתקיים שהקודקודים 2j-1 וגם 2j אינם. <p>אלגוריתם:</p> <ul style="list-style-type: none"> הכנס לעץ עם תדירות 0. קודד רק אז עדכן את העץ. <p>אלגוריתם העדכון: (כאשר x הוא העלה שהגענו אליו באחד משלבי הקידוד וכעת עלינו לעדכן)</p> </div> <div data-bbox="1062 351 1543 694" data-label="Text"> <pre>q = leaf(x) if (q is the 0-node) replace q by a parent 0-node with two 0-node children; q = left child; if (q is a sibling of a 0-node) interchange q with the highest numbered leaf of the same weight; increment q weight by 1; q = parent of q; while q!= root interchange q with the highest numbered node of the same weight; increment q weight by 1; q = parent of q; increment q weight by 1;</pre> </div> <div data-bbox="1222 694 1378 716" data-label="Section-Header"> <h2>Huffman (Skeleton)</h2> </div> <div data-bbox="1485 716 1559 739" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="1062 739 1559 842" data-label="List-Group"> <ul style="list-style-type: none"> צור קוד הופמן קנוני משוך ימינה הסר את כל תתי עצים השלמים שגובהם לכל הפחות 1. תן לצומת פנימי ערך 0. ותן לעלים אורך מילות הקוד מהשורה (הכללי) עד לעלים (של תת העץ שמחזקו). </div> <div data-bbox="1495 842 1559 864" data-label="Text"> <p>הגדרות:</p> </div> <div data-bbox="1062 864 1559 1296" data-label="Text"> <p>n_i = מספר מילות הקוד באורך i $\min\{i n_i > 0\} = m$. כלומר, זאת המילה עם האינדקס המינימלי. i ערך מספרי של מילת הקוד הראשונה באורך i $base(i)$ $0 = base(m)$ $2(base(i-1) + n_{i-1}) = base(i)$ $B_s(k)$ = התצוגה של המספר k בביטים בדיוק 0 (משמאל אם צריך). מילת הקוד באינדקס j באורך i (עבור $0 < j \leq n_i$) היא $B_i(base(i) + j)$ $seq(i)$ = מספר האינדקס של המילה הראשונה באורך i $0 = seq(m)$ $seq(i-1) + n_{i-1} = seq(i)$ $I(w)$ = הערך המספרי של הייצוג הבינארי w אם w באורך l אז מתקיים $w = B_l(I(w))$ מכאן: $I(w) = base(l) + B_l(I(w))$ = האינדקס הפנימי של מילת הקוד w בתוך הבלוק של מילות הקוד באורך l ולכן: $seq(l) + I(w) = base(l) + B_l(I(w))$ = האינדקס הכללי של מילת הקוד w (כללי) = ביחס לכלל מילות הקוד. אינדקס כללי: $diff(l) = I(w) - seq(l)$ כאשר $I(w) = base(l) + B_l(I(w))$ אלגוריתם פענוח:</p> </div> <div data-bbox="1062 1296 1543 1639" data-label="Text"> <pre>tree_pointer=root i=1 start=1 while i < length_of_string if string[i]=0 tree_pointer=left(tree_pointer) else tree_pointer=right(tree_pointer) if value(tree_pointer)>0 codeword=string[start...(start+value(tree_pointer)-1)] output=table[I(codeword)-diff[value(tree_pointer)]] tree_pointer=root start=start+value(tree_pointer) i=start else i++</pre> </div> <div data-bbox="1203 1639 1394 1662" data-label="Section-Header"> <h2>Arithmetic Code (Static)</h2> </div> <div data-bbox="1430 1662 1559 1684" data-label="Text"> <p>אלגוריתם קידוד:</p> </div> <div data-bbox="1110 1684 1485 1901" data-label="Text"> <pre>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]</pre> </div> <div data-bbox="1062 1901 1559 1924" data-label="List-Group"> <ul style="list-style-type: none"> גודל האינטרבל הוא כפל הסתברויות כל אות בהודעה שקידדנו. </div> <div data-bbox="1318 1924 1559 1946" data-label="Text"> <p>אלגוריתם פענוח (עבור מספר n):</p> </div> <div data-bbox="1078 1946 1525 2042" data-label="Text"> <pre>Find symbol whose range contains n Output the symbol Range = high(symbol) - low(symbol) Encoded = (encoded-low(symbol))/range</pre> </div> <div data-bbox="1190 2042 1410 2065" data-label="Section-Header"> <h2>Arithmetic Code (Adaptive)</h2> </div> <div data-bbox="1235 2065 1559 2087" data-label="Text"> <p>$N(a)$ = מספר התדירויות שהאות a הופיעה.</p> </div> <div data-bbox="1315 2087 1559 2110" data-label="Text"> <p>$P(a)$ = ההסתברות לקבל את a.</p> </div> <div data-bbox="1307 2110 1559 2132" data-label="Equation-Block"> $P(a) = \frac{N(a)+1}{N(a)+N(b)+N(c)+3}$ </div> <div data-bbox="1394 2132 1559 2154" data-label="Text"> <p>כאשר: $N(a)+N(b)+N(c)+3$</p> </div> <div data-bbox="1394 2154 1559 2177" data-label="Text"> <p>וכן לכל אות באלף בית.</p> </div> </td></td>	<div data-bbox="751 103 791 125" data-label="Section-Header"> <h2>LZ77</h2> </div> <div data-bbox="683 125 857 147" data-label="Text"> <p>(offset,length,symbol)</p> </div> <div data-bbox="909 147 1031 170" data-label="Text"> <p>אלגוריתם קידוד:</p> </div> <div data-bbox="555 170 983 304" data-label="Text"> <pre>p=1 // next char to be coded while (there is text to be coded): search for the longest match for S[p...] in s[p-w...p-1] // suppose match at pos m with len = 1 output the triple (p-m,l,s,p+1) p = p+l+1</pre> </div> <div data-bbox="751 311 791 333" data-label="Section-Header"> <h2>LZSS</h2> </div> <div data-bbox="729 333 1026 356" data-label="Text"> <p>(offset,length) [1] או בתווים בודדים [0].</p> </div> <div data-bbox="751 356 783 378" data-label="Section-Header"> <h2>LZS</h2> </div> <div data-bbox="627 378 909 400" data-label="Text"> <p>(offset1,length1), (offset2,length2)</p> </div> <div data-bbox="751 400 791 423" data-label="Section-Header"> <h2>LZ78</h2> </div> <div data-bbox="951 423 1031 445" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="533 445 1031 672" data-label="List-Group"> <ul style="list-style-type: none"> בנה מילון D שיכיל בהתחלה רק את המילה הריקה (שהאינדקס שלה הוא 0) הכנס תתי מחרוזת (של המחרוזות שנרצה לקודד) ותן לה את האינדקס הבא בתור שפנוי. נסה "לתפוס" את המילה הכי ארוכה שכבר במילון ונקודד את האינדקס שלה + הקידוד של התו הבא. במקביל, הוסף למילון את אותה הרישה בשרשרת התו הבא שמקודדים שיטה לדעת כמה ביטים צריך לקידוד מס': נשאל – "כמה ביטים צריך בשביל לקודד x מספרים" כאשר x הוא האינדקס באותה שורה </div> <div data-bbox="751 656 791 678" data-label="Section-Header"> <h2>LZW</h2> </div> <div data-bbox="595 678 1031 701" data-label="Text"> <p>אלגוריתם קידוד: (טבלה התחלתית גדולה ממש מגודל הא"ב)</p> </div> <div data-bbox="584 701 946 1081" data-label="Text"> <pre>Dictionary=single characters w=first char of input repeat { k=next char if(EOF) output code(w) else if (w · k)∈Dictionary w = w · k else output code(w) // output in num of bits corresponding to table size Dictionary.add(w · k) // Only enlarge the table when entering information and when it is full w = k }</pre> </div> <div data-bbox="903 1081 1031 1104" data-label="Text"> <p>אלגוריתם פענוח:</p> </div> <div data-bbox="584 1104 1031 1126" data-label="Text"> <p>חשוב: המפענת תמיד מגדיל את המילון ברגע שהמילון מתמלא</p> </div> <div data-bbox="584 1126 946 1514" data-label="Text"> <pre>Initialize table with single character strings OLD = first input code Output translation of OLD While not end of input stream { NEW = next input code If new is not in the string table S = translation of OLD S = S·C else S = translation of NEW Output S C = first character of S Enter Translation(OLD)·C to the string table // Increase the table as soon as the dictionary fills up OLD = NEW }</pre> </div> <div data-bbox="702 1536 834 1559" data-label="Section-Header"> <h2>Run length Code</h2> </div> <div data-bbox="716 1559 1031 1581" data-label="Text"> <p>זוג סדור (num,symbol) ע"פ רצף של תווים.</p> </div> <div data-bbox="842 1581 1031 1603" data-label="Text"> <p>Binary Run length Code</p> </div> <div data-bbox="746 1603 1031 1626" data-label="Text"> <p>escape codeword = ספרה מקסימלית.</p> </div> <div data-bbox="638 1626 898 1648" data-label="Section-Header"> <h2>The Burrows Wheeler Transform</h2> </div> <div data-bbox="951 1648 1031 1671" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="520 1671 1031 1879" data-label="List-Group"> <ul style="list-style-type: none"> בנה מטריצה מסדר n×n (כאשר הטקסט באורך n). כך שהשורות יהיו כל הטרנספורמציות הציקליות של הטקסט. כלומר, השורה הראשונה תתחיל באות הראשונה של הטקסט. השורה השנייה תתחיל מהאות השנייה של הטקסט וכו' (ממשיכים לשכתב את הטקסט בצורה מעגלית בסוף השורה) מייין את השורות בסדר לקסיקוגרפי (לפי האלף בית abcdefghijklmnopqrstuvwxyz) החזר את העמודה האחרונה במטריצה (מלמעלה למטה) + אינדקס של השורה המקורית של הטקסט (ספירה מ0). </div> <div data-bbox="978 1879 1031 1901" data-label="Text"> <p>פענוח:</p> </div> <div data-bbox="903 1901 1031 1924" data-label="Text"> <p>אלגוריתם פענוח:</p> </div> <div data-bbox="679 1924 847 2058" data-label="Text"> <pre>BWT(L, index): F = sort(L) M = index For (i=0; i < n; i++) T[i] = F[M] M = S[M]</pre> </div> <td data-bbox="1043 96 1568 2172"> <div data-bbox="1222 116 1378 138" data-label="Section-Header"> <h2>Huffman (Adaptive)</h2> </div> <div data-bbox="1430 138 1559 161" data-label="Text"> <p>Sibling property:</p> </div> <div data-bbox="1085 161 1559 351" data-label="List-Group"> <ul style="list-style-type: none"> משקל צומת שווה לסכום המשקל של שני ילדיו. הצמתים ניתנים לסידור ע"פ משקלם בסדר \geq כך שמתקיים שהקודקודים 2j-1 וגם 2j אינם. <p>אלגוריתם:</p> <ul style="list-style-type: none"> הכנס לעץ עם תדירות 0. קודד רק אז עדכן את העץ. <p>אלגוריתם העדכון: (כאשר x הוא העלה שהגענו אליו באחד משלבי הקידוד וכעת עלינו לעדכן)</p> </div> <div data-bbox="1062 351 1543 694" data-label="Text"> <pre>q = leaf(x) if (q is the 0-node) replace q by a parent 0-node with two 0-node children; q = left child; if (q is a sibling of a 0-node) interchange q with the highest numbered leaf of the same weight; increment q weight by 1; q = parent of q; while q!= root interchange q with the highest numbered node of the same weight; increment q weight by 1; q = parent of q; increment q weight by 1;</pre> </div> <div data-bbox="1222 694 1378 716" data-label="Section-Header"> <h2>Huffman (Skeleton)</h2> </div> <div data-bbox="1485 716 1559 739" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="1062 739 1559 842" data-label="List-Group"> <ul style="list-style-type: none"> צור קוד הופמן קנוני משוך ימינה הסר את כל תתי עצים השלמים שגובהם לכל הפחות 1. תן לצומת פנימי ערך 0. ותן לעלים אורך מילות הקוד מהשורה (הכללי) עד לעלים (של תת העץ שמחזקו). </div> <div data-bbox="1495 842 1559 864" data-label="Text"> <p>הגדרות:</p> </div> <div data-bbox="1062 864 1559 1296" data-label="Text"> <p>n_i = מספר מילות הקוד באורך i $\min\{i n_i > 0\} = m$. כלומר, זאת המילה עם האינדקס המינימלי. i ערך מספרי של מילת הקוד הראשונה באורך i $base(i)$ $0 = base(m)$ $2(base(i-1) + n_{i-1}) = base(i)$ $B_s(k)$ = התצוגה של המספר k בביטים בדיוק 0 (משמאל אם צריך). מילת הקוד באינדקס j באורך i (עבור $0 < j \leq n_i$) היא $B_i(base(i) + j)$ $seq(i)$ = מספר האינדקס של המילה הראשונה באורך i $0 = seq(m)$ $seq(i-1) + n_{i-1} = seq(i)$ $I(w)$ = הערך המספרי של הייצוג הבינארי w אם w באורך l אז מתקיים $w = B_l(I(w))$ מכאן: $I(w) = base(l) + B_l(I(w))$ = האינדקס הפנימי של מילת הקוד w בתוך הבלוק של מילות הקוד באורך l ולכן: $seq(l) + I(w) = base(l) + B_l(I(w))$ = האינדקס הכללי של מילת הקוד w (כללי) = ביחס לכלל מילות הקוד. אינדקס כללי: $diff(l) = I(w) - seq(l)$ כאשר $I(w) = base(l) + B_l(I(w))$ אלגוריתם פענוח:</p> </div> <div data-bbox="1062 1296 1543 1639" data-label="Text"> <pre>tree_pointer=root i=1 start=1 while i < length_of_string if string[i]=0 tree_pointer=left(tree_pointer) else tree_pointer=right(tree_pointer) if value(tree_pointer)>0 codeword=string[start...(start+value(tree_pointer)-1)] output=table[I(codeword)-diff[value(tree_pointer)]] tree_pointer=root start=start+value(tree_pointer) i=start else i++</pre> </div> <div data-bbox="1203 1639 1394 1662" data-label="Section-Header"> <h2>Arithmetic Code (Static)</h2> </div> <div data-bbox="1430 1662 1559 1684" data-label="Text"> <p>אלגוריתם קידוד:</p> </div> <div data-bbox="1110 1684 1485 1901" data-label="Text"> <pre>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]</pre> </div> <div data-bbox="1062 1901 1559 1924" data-label="List-Group"> <ul style="list-style-type: none"> גודל האינטרבל הוא כפל הסתברויות כל אות בהודעה שקידדנו. </div> <div data-bbox="1318 1924 1559 1946" data-label="Text"> <p>אלגוריתם פענוח (עבור מספר n):</p> </div> <div data-bbox="1078 1946 1525 2042" data-label="Text"> <pre>Find symbol whose range contains n Output the symbol Range = high(symbol) - low(symbol) Encoded = (encoded-low(symbol))/range</pre> </div> <div data-bbox="1190 2042 1410 2065" data-label="Section-Header"> <h2>Arithmetic Code (Adaptive)</h2> </div> <div data-bbox="1235 2065 1559 2087" data-label="Text"> <p>$N(a)$ = מספר התדירויות שהאות a הופיעה.</p> </div> <div data-bbox="1315 2087 1559 2110" data-label="Text"> <p>$P(a)$ = ההסתברות לקבל את a.</p> </div> <div data-bbox="1307 2110 1559 2132" data-label="Equation-Block"> $P(a) = \frac{N(a)+1}{N(a)+N(b)+N(c)+3}$ </div> <div data-bbox="1394 2132 1559 2154" data-label="Text"> <p>כאשר: $N(a)+N(b)+N(c)+3$</p> </div> <div data-bbox="1394 2154 1559 2177" data-label="Text"> <p>וכן לכל אות באלף בית.</p> </div> </td>	<div data-bbox="1222 116 1378 138" data-label="Section-Header"> <h2>Huffman (Adaptive)</h2> </div> <div data-bbox="1430 138 1559 161" data-label="Text"> <p>Sibling property:</p> </div> <div data-bbox="1085 161 1559 351" data-label="List-Group"> <ul style="list-style-type: none"> משקל צומת שווה לסכום המשקל של שני ילדיו. הצמתים ניתנים לסידור ע"פ משקלם בסדר \geq כך שמתקיים שהקודקודים 2j-1 וגם 2j אינם. <p>אלגוריתם:</p> <ul style="list-style-type: none"> הכנס לעץ עם תדירות 0. קודד רק אז עדכן את העץ. <p>אלגוריתם העדכון: (כאשר x הוא העלה שהגענו אליו באחד משלבי הקידוד וכעת עלינו לעדכן)</p> </div> <div data-bbox="1062 351 1543 694" data-label="Text"> <pre>q = leaf(x) if (q is the 0-node) replace q by a parent 0-node with two 0-node children; q = left child; if (q is a sibling of a 0-node) interchange q with the highest numbered leaf of the same weight; increment q weight by 1; q = parent of q; while q!= root interchange q with the highest numbered node of the same weight; increment q weight by 1; q = parent of q; increment q weight by 1;</pre> </div> <div data-bbox="1222 694 1378 716" data-label="Section-Header"> <h2>Huffman (Skeleton)</h2> </div> <div data-bbox="1485 716 1559 739" data-label="Text"> <p>אלגוריתם:</p> </div> <div data-bbox="1062 739 1559 842" data-label="List-Group"> <ul style="list-style-type: none"> צור קוד הופמן קנוני משוך ימינה הסר את כל תתי עצים השלמים שגובהם לכל הפחות 1. תן לצומת פנימי ערך 0. ותן לעלים אורך מילות הקוד מהשורה (הכללי) עד לעלים (של תת העץ שמחזקו). </div> <div data-bbox="1495 842 1559 864" data-label="Text"> <p>הגדרות:</p> </div> <div data-bbox="1062 864 1559 1296" data-label="Text"> <p>n_i = מספר מילות הקוד באורך i $\min\{i n_i > 0\} = m$. כלומר, זאת המילה עם האינדקס המינימלי. i ערך מספרי של מילת הקוד הראשונה באורך i $base(i)$ $0 = base(m)$ $2(base(i-1) + n_{i-1}) = base(i)$ $B_s(k)$ = התצוגה של המספר k בביטים בדיוק 0 (משמאל אם צריך). מילת הקוד באינדקס j באורך i (עבור $0 < j \leq n_i$) היא $B_i(base(i) + j)$ $seq(i)$ = מספר האינדקס של המילה הראשונה באורך i $0 = seq(m)$ $seq(i-1) + n_{i-1} = seq(i)$ $I(w)$ = הערך המספרי של הייצוג הבינארי w אם w באורך l אז מתקיים $w = B_l(I(w))$ מכאן: $I(w) = base(l) + B_l(I(w))$ = האינדקס הפנימי של מילת הקוד w בתוך הבלוק של מילות הקוד באורך l ולכן: $seq(l) + I(w) = base(l) + B_l(I(w))$ = האינדקס הכללי של מילת הקוד w (כללי) = ביחס לכלל מילות הקוד. אינדקס כללי: $diff(l) = I(w) - seq(l)$ כאשר $I(w) = base(l) + B_l(I(w))$ אלגוריתם פענוח:</p> </div> <div data-bbox="1062 1296 1543 1639" data-label="Text"> <pre>tree_pointer=root i=1 start=1 while i < length_of_string if string[i]=0 tree_pointer=left(tree_pointer) else tree_pointer=right(tree_pointer) if value(tree_pointer)>0 codeword=string[start...(start+value(tree_pointer)-1)] output=table[I(codeword)-diff[value(tree_pointer)]] tree_pointer=root start=start+value(tree_pointer) i=start else i++</pre> </div> <div data-bbox="1203 1639 1394 1662" data-label="Section-Header"> <h2>Arithmetic Code (Static)</h2> </div> <div data-bbox="1430 1662 1559 1684" data-label="Text"> <p>אלגוריתם קידוד:</p> </div> <div data-bbox="1110 1684 1485 1901" data-label="Text"> <pre>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]</pre> </div> <div data-bbox="1062 1901 1559 1924" data-label="List-Group"> <ul style="list-style-type: none"> גודל האינטרבל הוא כפל הסתברויות כל אות בהודעה שקידדנו. </div> <div data-bbox="1318 1924 1559 1946" data-label="Text"> <p>אלגוריתם פענוח (עבור מספר n):</p> </div> <div data-bbox="1078 1946 1525 2042" data-label="Text"> <pre>Find symbol whose range contains n Output the symbol Range = high(symbol) - low(symbol) Encoded = (encoded-low(symbol))/range</pre> </div> <div data-bbox="1190 2042 1410 2065" data-label="Section-Header"> <h2>Arithmetic Code (Adaptive)</h2> </div> <div data-bbox="1235 2065 1559 2087" data-label="Text"> <p>$N(a)$ = מספר התדירויות שהאות a הופיעה.</p> </div> <div data-bbox="1315 2087 1559 2110" data-label="Text"> <p>$P(a)$ = ההסתברות לקבל את a.</p> </div> <div data-bbox="1307 2110 1559 2132" data-label="Equation-Block"> $P(a) = \frac{N(a)+1}{N(a)+N(b)+N(c)+3}$ </div> <div data-bbox="1394 2132 1559 2154" data-label="Text"> <p>כאשר: $N(a)+N(b)+N(c)+3$</p> </div> <div data-bbox="1394 2154 1559 2177" data-label="Text"> <p>וכן לכל אות באלף בית.</p> </div>
--	---	---