

①

## 3 ח'סות נתונים א' / צנה שפירא

לספ' מבחן, יום 4-5 מלכות.

### נושא 1 - מבוא

#### שאלות של בחינת נתונים

- הצביר מוצג רק יותר בתשתית מונעלת.
- שכיח לאני התגובה של אפליקציות מסוימות.
- שכיח יכולות מחשב ע' שאינה מידע רק יותר בליכרון פנא' לא בליכרון משני.

נושא פ:

מבוא

מידע

קידום סט

קידום טון

קידום הסאן

קידום איתאט

קידום מילוני

L277

L278

#### מציאת בחינה

גם מציאת בחינה מורכבת מ-3 שלבים:

מידע - עמידה והנחיות על מענה הדוק, כאן מהו הא' הדוק.

איטל סטטסטי-על הדוק, וכאן כיצד כמות האיתות מתבטאת.

קידום - בחירת אלמנטים שקידום הדוק, כאן מוצג א' בינארי.

בכך בחינה יש לקודד ומענה. המקודד מקבל את הדוק המקורי ומקודד אותו עידי' ביט'.

המענה מקבל את הוצ' הביט' ומענה אתם חזרה עקובל המקורי או קובל צורה זו.

יש שני סוגים של בחינה: lossless שאינה מאבדת אידע בק שהדוקל האשומר לזה עקובל, בעידי' חשיפה בקבצי.

לסל. י- lossless שאבדת קרת מהאידע אך שאינה את העידר, משמשת בעידר בחינת תמונות וקבצי סאונד.

#### ס'מאניס (הדורות)

• את קבוצת הא' בדוקל נסמן  $S = [s_1, s_2, \dots, s_n]$

• את ההסתברות שהא' בדוקל נסמן  $P = [p_1, p_2, \dots, p_n]$ . נניח כי  $p_1 \geq p_2 \geq \dots \geq p_n$ , וזש  $\sum_{i=1}^n p_i = 1$

• את הקידום של א' נסמן  $C = [c_1, c_2, \dots, c_n]$  ואת אורך הקידום של  $c$  נסמן  $|c|$ .

• התוחלת של אורך קידום עבור א' אחת נחשב  $E(C, P) = \sum_{i=1}^n p_i |c_i|$ . אורך זה הוא קודד של עידי' בחינה.

הבחיסה. ככל ש-  $E(C, P)$  נמוך יותר כך הבחינה טובה יותר.

#### 3) ס'מאניס קידומי

• Prefix-free: זהו קידום שבו אין קידום של מילה שהוא רישא של קידום מילה אחרת. הדבר הנות,

עייתי' עשאוה קידום חסר רישות היא באמצעות ע' שבו כל כניה שאעלר היא ס ופנ"ה "מער היא 1. כע'.

כל העליס הם תנ"ס. בכל פעמח של תו נש"ע מהשיט עז עעלר. קיז חסר רישות אנו אוריז את אנות הבחי



②

זוגות נשואים:  
 111, 111, 111, 111  
 נדע. שאם לא  
 כמילאם ס' נעק יו  
 כענין יח' אק קי' צו  
 לא חסר יח' יח'.

בדרך אחת ניסיה שתחב'ר את הקיבוץ עצומתו האקורית.

משפט: כל קובץ שהוא חסר תועלת הוא גם UD, אך ההיפך לא נכון.

- Complete Code: זהו קידוש שמעור כס אחריית פגארי (סאן צוקא טעריה בארצות הקידוש) נהן

סמלנות באופן "חידושי". אך נ"צ ציצי צמק קולג' אני הרף יהיה בהכרח שם, שהוא קולג' טבו ספ

דואר סביצ'יט י"ש שט' בליק אק א'נו פהכרח לעם. אפיר קיז שבו א' לעם, י"ז ע' אחיות מלך

באופן "חיד" או תבואה דה"ר מכוללת כנראה, במסגרת נ"מ יהיה נמצאים ונ"מ את הקידוד.

צואה עקיצה עא עס: 101, 100, 011, 010, 001, 000. \* a, b, c, d, e, f. ניתן לראות שעבור אחריות בינארית

שנה יש בסוף 11 ע"א יהיה ניתן לעצמנו אימה, וכן שאת הקידוש  $100, 101$   $\begin{smallmatrix} e \\ f \end{smallmatrix}$  היה אפשר לקרר 8-11.10

שנה טובה 5787

• Instantaneous Code: נחו קידוד שבו בעצמם הם אחראים סינאקרת ניתן לדעת את מישה

[illegible]

שמירה הסת"מ רק האמר מל"ע"ס ע"ס הש"ק עמ"ל הבאה, וע"ן א"ט ק"י"ב מ"ב. ק"י"ב מ"ב א"ט

זמר, אגנס, ק'צ'צ' א"צ' זמר, ק'צ'צ' תסר, י'ש'ת, וכן, ק'צ'צ' חסר, י'ש'ת, זמר, ק'צ'צ' א"צ'.

במחצית 1 שחיות  
16 קידיו 3 מה"ס  
יור כח התנאים.

אמרינו הוא עצמו קיצי' שאק"ם את כל ארבעת התנאים שם שהוא גם בעל תחושת איפה הכל קטנה טעם

Sardinas - Patterson פת'אדק - UD | מחנך | ח

נאמז אעזארימס פונ'טא' אטור בהיתתן קידוצ אחויר האק היא טו אז עא, עכנ' ת'אור האל אזרימ

נספח של מושיק. נזכר בסוף ע-6 היא דבורה כה הוי"מ העידות.

• Suffix dangling - עפוי שט' קיינדיג אים'א, ב, א; אק א זיטש על אזי טאר הייט'ס פ-ב (שאנק פ-א) דק ס'טא וואו

• עבור שני קבוצות אינסופיות  $S, T$  נגדיר  $\text{card } S \leq \text{card } T$  אם קיימת פונקציה  $f: S \rightarrow T$  חד-חד-ערכית.  $S^{-1}T = \{a \in T \mid a \in S\}$ . כלומר קבוצת

Sardinas - Patterson §

דאס היסטארישע דאטא צווישן דעם יאר 1945-1946.

Si בזה קבוצה של אינסוף שמיאל את כל הסופות הקבוצה

$j = 1$

$S_1 = C^{-1}C - \{E\}$

// זכור! קבוצה של מילים שחבר את כל המילים הקטנות.

שעל זה יש לי אסטרטגיה וזוהי סוגית תלות של  $C$  עם עצמה.

```
while (true) {
```

while  $(\pm f u c) \neq$   
 $S_{i+1} = C^{-1} S_i \cup S_i^{-1} C$  //  $S_i$  - סיסטם תצורות שיש בו  $S_{i+1}$  - סיסטם תצורות שיש בו  $C$  ובהם  $C$  ובהם  $C$

;

```
IF  $e \in S_1$  OR  $c \in S_1$  FOR  $c \in C$  PRINT NOT UD, EXIT
```

$\text{else if } \exists j < i \text{ s.t. } S_i = S_j \text{ Print UD, Exit}$  // הסימן ב- $S_i$ . היסטוריה של  
כל הערכות והמיקומים האפשריים.

מטרת האדמ'יתם היא ש'מאצ'א ס'מא תע'ניה ב' - Si שיה'א ע'א ד'יצ'א'ל ע'ל מ'י'ס'ה ב'ע'א. ע'א נ'אצ'א

ס'טא מעלו'ק גלייב'יק היק'דיג אונד און ס'פאט גלייב'יק שטיק'דיג און ד.

פאמילי נאם  
טאטא באזיר 1  
שקופ'אר 21-22.



26/10

(3)

# 1) מושג הקידוד של שטון (Shannon)

באזור עולם, האטרה בדחיסת נתונים היא עמקוא, עביר דובל שבו התפלגות  $P$  מסומלת, קידוד שעבור אירק מילת קוד מאוצרת היא מימלית. אק הקידוד אכן מקיים תנאי זה נאמר שהוא "קידוד בעל יתריות מימלית", כלומר שהמידע המיותר בו מימלית. בניסוח נורמלי יותר נאמר ש- $C$  הוא בעל יתריות מימלית אם מתקיים  $E(C, P) \leq E(C^*, P)$  עבור כל קידוד  $C^*$  אחר.

חסם תחתון זה הוא אכן מימלית עבור כל קוד שבו לכל תו קטקסט יש אירק קוד קבוע אלא גשתה.

מושג הקידוד של שטון נותן לנו חסם תחתון  $E(C, P)$  כדי להסביר חסם זה יש להגדיר כמות לעמידה במתא אינפורמציה: שטון הגדיר שכמות האינפורמציה החזקה שאכלה  $S_i$  עם הסתברות  $P_i$  מחושבת עכ

$$I(S_i) = -\log_2 P_i$$

תכונות של אינפורמציה - בכל ש- $P_i$  מתקרב ע-1 כך האינפורמציה  $I(S_i)$  מתקרב ע-0. בכל ש- $P_i$  מתקרב ע-0 כך האינפורמציה  $I(S_i)$  מתקרב ע- $\infty$ .

אינפורמציה של רצף תווים  $S_1, S_2, \dots, S_n$  מחושבת:  $I(S_1, S_2, \dots, S_n) = I(S_1) + I(S_2) + \dots + I(S_n)$

אנתרופיה של התפלגות: איצת האינפורמציה הק אירק: בימים שטון הוכיח שהקידוד יהאוב ביותר הוא כזה שטמן עכס תו  $S_i$  באירק  $I(S_i)$ . בדויה כוילע'ת יותר, עבור התפלגות  $P$  של תווים, נגדיר  $H(P)$  עהיות האירק המאוצל של האינפורמציה עבור כל התווים,  $H(P) = -\sum_{i=1}^n P_i \cdot \log_2 P_i$ . עכס קידוד  $C$  מתקיים:  $H(P) \leq E(C, P)$  כלומר  $H(P)$  הוא חסם תחתון ע- $E(C, P)$ , עכן תאוצ נידה לתת עתל  $S_i$  קידוד באירק  $I(S_i)$ .  $H(P)$  נקרא "אנתרופיה של התפלגות  $P$ ".

## 2) אי-שוויון קראפט-מקמילן (Kraft-McMillan)

עכ חסם שטון ענעל אק קימלת התפלגות של תווים  $P$ , כך שכל תו  $S_i$  ההסתברות שלו היא  $P_i = 2^{-k_i}$  עכיה  $M \in \mathbb{N}$ , אזי  $I(S_i) = k_i$  מסביר שכל. בהתפלגות כזו ניתן עידור בקלות (למצ בהאטק) קידוד שהוא בדיוק שיה עחס שטון. הסבר זה הוא המוטבציה עסונקצית קראפט עכ כל קידוד  $C$  האולגרת כך:  $K(C) = \sum_{i=1}^n 2^{-|C_i|}$ . הקרק של  $K(C)$  יכול עכמז אלתו עטנעלם רבים חשויים בדחיסת נתונים: (1) אי שוויון קראפט-מקמילן: קידוד  $C$  הוא מימי אק ירק אק  $K(C) \leq 1$ .

(2) אק עביר קידוד  $C$  מתקיים  $K(C) \leq 1$  אזי קיים קידוד  $C^*$  המקיים:  $E(C^*, P) = E(C, P)$ ,  $|C^*| = |C|$ , ו- $C^* \leq C$ . הוא חסר חיות, כלומר קיים קידוד טקול טחז חסר חיות. מסיבה זו תאוצ נידה שהקידוד שעלנו יהיה חס חיות, שכן הוא עכ אורצ מליאת הדחיסה ויש ען יתרונות רבים נלמאר ענעל.

מסקנה א- (2) ו- (3) חסר חיות  $(K(C) \leq 1 \Rightarrow C$  חיות)

(3) אק עביר קידוד  $C$  מתקיים  $K(C) > 1$  אזי הוא עוידא עכ חסר חיות, עכ מימי ונדא  $\infty$ . (4) ו- (5) נאצת עקיות 8-9.

(5) שוויון מקמילן: קידוד  $C$  הוא שלם (ניתן עידעו בעל עכא שלא טמן עהיט 4 בו עכז עעים) אק  $K(C) = 1$ .

הוכחת עמטעל (1) ו- (4) נאצת עקיות 8-9.



26/10

בס"ב

(4)

האסקנה מכל הצטננים ע"פ היא שצמוד התנאים  $P$  של תנ"ס, גבי עקום קיפוז  $C$  שהוא יחס:

היות, א"פ,  $UD$  שלם בעל יחיות אינמל'ת, אסמ'ק ש"ק"פ שני תנא'ס:

$$(1) \quad K(C) \leq 1$$

$$(2) \quad E(C, P) \text{ אינמל'ת.}$$

אם בנוסף נרצה קיב שלם נצטרף  $K(C) = 1$ .



(5)

## נושא 2 - מצב

(א) הדברה

כאמור לעיל, כל מידת דחיסה מוכנית מ-3 מרכיבים: אורך, איסוף סטטיסטי ואלגוריתם קידוד.

המקרה של נוסף בארכיטקטורה הראשונה ונכנס מספר סוגי מידע. המידע בין המידע הוא הכיבדק המכיל את המידע.

את הקובץ / המידע שאנו מוציאים מכאן, וכתוצאה מכך יש לנו ה-Header, הנקרא גם "Prelude".

שניתן לתת עמנו (decoder) כדי שיוכל לפרש את הדחיסה. עבור כל מידע נחשב את

האנטינטיב  $H(p)$  ואת האורך הממוצע של קידוד  $E(c,p)$ . נבחר את היחס ביניהם.

נסתכל על סוגי מידע של מידעיות מידע.

מ נמצא שיש

סוגי מידע כאלו.

(1) Zero-order: מידע שבהם ההסתברות שנתונים מידע קידוד הן בלי תלות. נקרא גם "אינדיקטור מסדר 0".

(2) First-order: מידע שיש יתרון בין מידע קידוד. נתונים הסתברות על שתי מידע קידוד צמודות ביחס.

מספר המידע של הדחיסה. נקרא גם "מקור מסדר 1".

(ב) מידע דחיסה סטטית

המידע של מידע של הדחיסה שמוצא בדחיסה, נוסף להכנת שיש הדחיסה את כל התווים בלי

ה-ASCII, סה"כ 256 תווים. כל תו יש אורך הסתברות  $P_i = \frac{1}{256}$ . קידוד כל תו יעשה על ידי

של טבלת ASCII שמתן על תו קידוד באורך 8 ביט.

עבור מידע זה האנטינטיב היא:  $H(p) = - \sum_{i=1}^{256} \frac{1}{256} \log_2 \left( \frac{1}{256} \right) = 8$  [ביטים]. נזכור שהאנטינטיב היא

חסר תחתון באורך כל קידוד אפשרי. עכשיו הקידוד שמיאנו עני טבלת ASCII, שהוא Fixed-code

ועכשיו יתן לנו  $E(c,p) = 8$  והוא אכן הטוב ביותר עבור מידע זה.

(ג) מידע דחיסה סטטית ממוצע

זהו מידע שכן מוצא את הדחיסה שמוצא בדחיסה. הוא נותן הסתברות אכן וכן עליון שמוצא בדחיסה

M. אנחנו כל תו אורך הסתברות אחידה. נסמן את אורך הא"ב  $|S|$ . עבור כל תו  $i$  ההסתברות

של תו  $i$  היא  $P_i = \frac{1}{|S|}$ . האנטינטיב עבור מידע זה היא:

$$H(p) = - \sum_{i=1}^{|S|} \frac{1}{|S|} \log_2 \left( \frac{1}{|S|} \right)$$

באורך זה צריך גם header המכיל את כל  $|S|$  בקידוד ASCII וכן את מספר התווים  $|S|$ . סה"כ

אורך ה-header יהיה 8 ביטים עבור  $|S|$  ועוד  $|S| \cdot 8$  ביטים עבור הקידוד של כל תו.

כל תו  $i$  יקבל  
קידוד באורך  $\log_2(|S|)$   
אנחנו בעצם  $|S|$   
על חלקה של 2 אנשי  
מספר זאת נא לשק  
העצם.

ומה מספר התווים  
בהדחה.

מכאן נחשב אורך קידוד תו ממוצע ה-M:

$$E(c,p) = H(p) + \frac{|S| \cdot 8}{|M|}$$

צנארה: עבור הדחה באורך 25 תווים שמיאנו 25 תווים שונים נקבע:

$$E(c,p) = - \sum_{i=1}^{25} \frac{1}{25} \log_2 \left( \frac{1}{25} \right) + \frac{25 \cdot 8}{128} = 6.27$$

מידע זה טוב יותר מדחיסה סטטית. שיטת טבלה היא שהיא ה-25 טבלה להיות מספר יק 11 בין אין



(6)

### ב) מודל בחיסה סטט'ר למחצה עם הסתברות עצמית

בהו מודל שמוצא את הקובץ ועל מידת קיבוץ נ"ל ניתן את ההסתברות עצמית מידה אקראית

מהינדסה/קובץ היא זמן  $P_i = \frac{V_i}{M}$ , כאשר  $V_i$  היא מספר ההיסטור של  $S$  בהינדסה.

האנטי-רופיה ציור מודל זה היא:  $H(P) = -\sum_{i=1}^{|S|} \frac{V_i}{M} \log_2 \left( \frac{V_i}{M} \right)$ . במודל זה נספר שהצדד הידועים

קצחים יותר משאר המודלים שמאדני. החיסיון הוא שה- header של הקובץ גדולים צריך להיות

גדול שכן הוא גם צריך להיות את התדירות של מידת קיבוץ  $S$ .  $E(C, P) = H(P) + \frac{\text{גודל header}}{M}$

בואו: עבור הינדסה באורך 25 תווים שמכילה 25 תווים שונים, וכל ערכי תדירות של נ"ל

$$E(C, P) = -\sum_{i=1}^{25} \frac{V_i}{25} \log_2 \left( \frac{V_i}{25} \right) + \frac{1+25 \cdot 8 + 25 \cdot 1}{25} = 6.63$$

תו צריך עמדות 4 ביט' נקבע:

$V_i$  נקרא משבחה  
במצב 2 שונים  
17

ניתן לראות שקצת בחית טוב מהצדקה במודל הקודם. אולם בהינדסה אנוכית מודל זה הידע יותר טוב.

### ה) מודל סדרה - ראשון

נ"ל כה שלושת המודלים היו מסוג zero-order שכן לא היה תלות בהסתברות בין שני תווים.

במודל First-order מחלקים את ההינדסה לכל שני תווים סמוכים  $S$  ו- $S$ , ואזי כל מידת קיבוץ נ"ל:

תקבע הסתברות  $P_{ij} = \frac{V_{ij}}{M}$ , כאשר  $V_{ij}$  הוא מספר ההיסטור של  $S$  ו- $S$  בהינדסה.

בשיטה זו מודל ה- header מאבד גדול אך האנטי-רופיה הטלה יותר. עכ"ן מודל זה שימושי רק

בהינדסה/קובץ גדולים מאוד.